
Unlearning Inversion Attacks for Graph Neural Networks

Jiahao Zhang¹, Yilong Wang¹, Zhiwei Zhang¹, Xiaorui Liu², Suhang Wang¹

¹The Pennsylvania State University, ²North Carolina State University
{jiahao.zhang, yvw5769, zbz5349, szw494}@psu.edu, xliu96@ncsu.edu

Abstract

Graph unlearning methods aim to efficiently remove the impact of sensitive data from trained GNNs without full retraining, assuming that deleted information cannot be recovered. In this work, we challenge this assumption by introducing the *graph unlearning inversion* attack: given only black-box access to an unlearned GNN and partial graph knowledge, can an adversary reconstruct the removed edges? We identify two key challenges: varying probability-similarity thresholds for unlearned versus retained edges, and the difficulty of locating unlearned edge endpoints, and address them with **TrendAttack**. First, we derive and exploit the *confidence pitfall*, a theoretical and empirical pattern showing that nodes adjacent to unlearned edges exhibit a large drop in model confidence. Second, we design an adaptive prediction mechanism that applies different similarity thresholds to unlearned and other membership edges. Our framework flexibly integrates existing membership inference techniques and extends them with trend features. Experiments on four real-world datasets demonstrate that TrendAttack significantly outperforms state-of-the-art GNN membership inference baselines, exposing a critical privacy vulnerability in current graph unlearning methods.

1 Introduction

Graph-structured data is prevalent in numerous real-world applications, such as recommender systems [73, 28, 94], social media platforms [23, 63, 46], and financial transaction networks [20, 51, 75]. Graph Neural Networks (GNNs) [27, 69, 79] have emerged as powerful tools for modelling such data, leveraging their ability to capture both node attributes and graph topology. The effectiveness of GNNs relies on the message-passing mechanism [26, 24], which iteratively propagates information between nodes and their neighbors. This enables GNNs to generate rich node representations, facilitating key tasks like node classification [36, 45], link prediction [96, 76], and graph classification [38, 25].

Despite their success, GNNs raise significant concerns about privacy risks due to the sensitive nature of graph data [62, 18]. Real-world datasets often contain private information, such as purchasing records in recommender systems [95, 83] or loan histories in financial networks [60, 72]. During training, GNNs inherently encode such sensitive information into their model parameters. When these trained models are shared via model APIs, privacy breaches may occur. These risks have led to regulations like the GDPR [52], CCPA [56], and PIPEDA [57], which enforce the *right to be forgotten*, allowing users to request the removal of their personal data from systems and models. This demand has necessitated the development of methods to remove the influence of specific data points from trained GNNs, a process known as *graph unlearning* [12, 10, 80].

A straightforward way of graph unlearning is to retrain the model from scratch on a cleaned training graph. However, this approach is computationally infeasible for large-scale graphs, such as purchase networks in popular e-commerce platforms [71, 34] and social networks on social media [2, 13], which may involve billions of nodes and edges. To address this, a wide range of recent graph unlearning

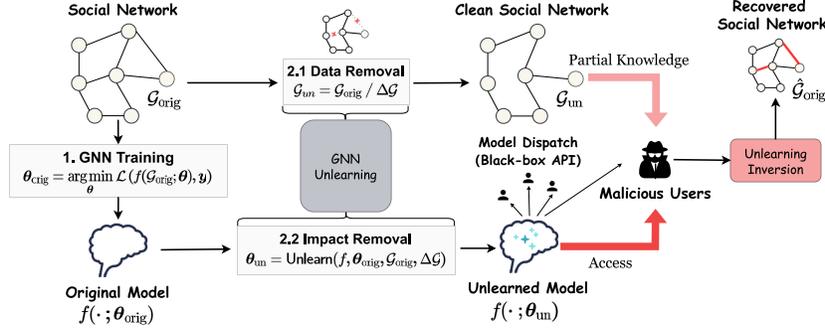


Figure 1: **Illustration of the unlearning inversion attack.** Considering an online social network $\mathcal{G}_{\text{orig}}$, where a user requests the deletion of sensitive friendship information, resulting in a cleaned graph \mathcal{G}_{un} and updated model parameters θ_{un} . The GNN model may be shared with third-parties via black-box APIs. If an attacker, leveraging the model API and auxiliary information about \mathcal{G}_{un} , can reconstruct the removed knowledge $\Delta\mathcal{G}$ through an unlearning inversion attack, sensitive relationships may be exposed, severely compromising user privacy.

methods focus on efficient parameter manipulation techniques [80, 81, 78], which approximate the removal of data by adjusting model parameters based on their influence. These techniques are often regarded as privacy-preserving, assuming that the unlearned data cannot be reconstructed.

In this paper, we challenge the assumption that existing graph unlearning methods are robust to privacy attacks. Specifically, we investigate a novel and important problem of whether unlearned information can be recovered through an emerging family of membership inference attacks, namely *unlearning inversion attacks* [32] (Figure 1). This attack is important because it reveals critical privacy vulnerabilities in graph unlearning methods and shows a concerning scenario where “forgotten” user information on the Web can be reconstructed through carefully designed attacks. Our central research question is: *Can third-parties exploit model APIs of unlearned GNNs to recover sensitive information that was meant to be forgotten?*

Although several pioneering studies have examined the privacy risks of unlearned ML models [9, 67, 32, 3], conducting an unlearning inversion attack on GNNs presents unique technical challenges. Specifically, we argue that the following prior works may not fully resolve the proposed research question: (i) **General-purpose unlearning inversion attacks** [32, 3] often assume attacker’s access to model parameters [3] or probability predictions [32] before and after unlearning. However, this assumption may not hold in real-world settings, as the pre-unlearning model is highly sensitive and may not be accessible to third parties. (ii) **Training a link prediction model** directly on \mathcal{G}_{un} is an intuitive way to recover unlearned edges from the cleaned graph. However, due to API query limitation policies of Twitter¹ and TikTok², such direct link prediction techniques are less practical in real-world settings, where only partial access to \mathcal{G}_{un} is available. (iii) **Membership inference attack (MIA)** [30, 29, 58, 103] methods designed for GNNs are another possible alternative. However, they are tailored to extract knowledge that is clearly “memorized” by the model, which is insufficient for unlearned GNNs that retain only minimal residual information of such sensitive knowledge.

In this work, we study the link-level unlearning inversion attack problem for black-box unlearned GNN models, aiming to accurately recover unlearned links using black-box GNN outputs and partial knowledge of the unlearned graph \mathcal{G}_{un} . We identify two key technical challenges for such attacks: (i) the output probabilities for two unlearned nodes may not exhibit sufficiently high similarity, which may require different prediction thresholds when inferring their membership compared to inferring other edges; (ii) it is highly non-trivial to determine whether a node is connected to an unlearned edge, since we only have access to the post-unlearning model, with no reference to pre-unlearning models. In response to these challenges, we propose a novel link-level unlearning inversion attack for black-box unlearned GNN models, namely **TrendAttack**. To address challenge (i), we design an adaptive prediction mechanism that applies different thresholds to infer two types of training graph edges: unlearned edges and other membership edges, enhancing TrendAttack’s flexibility to

¹<https://developer.x.com/en/docs/x-api>

²<https://developers.tiktok.com/doc/research-api-faq>

accommodate varying similarity levels. For challenge (ii), we identify a key phenomenon called the *confidence pitfall*, which enables the distinction between nodes connected to unlearned edges and others, using only black-box model outputs. This phenomenon describes how the model’s confidence in nodes near unlearned edges tends to decrease, and is supported by both empirical and theoretical evidence, as detailed in Section 5.1. By jointly incorporating the adaptive treatment of different edge types and the confidence trend pattern, TrendAttack achieves strong membership inference performance for both unlearned and other membership edges.

Our **main contributions** are: (i) We study a novel problem of graph unlearning inversion attack, pointing out the vulnerability of existing graph unlearning; (ii) We identify a simple yet effective pattern, the *confidence pitfall*, which distinguishes nodes connected to unlearned edges, supported by both empirical and theoretical evidence; (iii) We introduce a novel unlearning inversion attack **TrendAttack** that leverages confidence trend features in an adaptive membership inference framework, which can accurately identify unlearned edges and be flexibly integrated to existing GNN MIA methods; and (iv) Comprehensive evaluation on four real-world datasets shows that our method consistently outperforms state-of-the-art GNN MIA baselines in membership inference accuracy.

2 Related Works

Graph Unlearning. Graph Unlearning enables the efficient removal of unwanted data’s influence from trained graph ML models [10, 61, 22]. This removal process balances model utility, unlearning efficiency, and removal guarantees, following two major lines of research: retrain-based unlearning and approximate unlearning. Retrain-based unlearning partitions the original training graph into disjoint subgraphs, training independent submodels on them, enabling unlearning through retraining on a smaller subset of the data. Specifically, GraphEraser [10] pioneered the first retraining-based unlearning framework for GNNs, utilizing balanced clustering methods for subgraph partitioning and ensembling submodels in prediction with trainable fusion weights to enhance model utility. Subsequently, many studies [70, 43, 93, 42] have made significant contributions to improving the Pareto front of utility and efficiency in these methods, employing techniques such as data condensation [42] and enhanced clustering [43, 93]. Approximate unlearning efficiently updates model parameters to remove unwanted data. Certified graph unlearning [12] provides an important early exploration of approximate unlearning in SGC [79], with provable unlearning guarantees. GraphGuard [78] introduces a comprehensive system to mitigate training data misuse in GNNs, featuring a significant gradient ascent unlearning method as one of its core components. GIF [80] presents a novel influence function-based unlearning approach tailored to graph data, considering feature, node, and edge unlearning settings. Recent innovative works have further advanced the scalability [59, 39, 89, 86, 92] and model utility [41, 97] of approximate unlearning methods. In this paper, we explore the privacy vulnerabilities of graph unlearning by proposing a novel membership inference attack tailored to unlearned GNN models, introducing a new defense frontier that graph unlearning should consider from a security perspective.

Membership Inference Attack for GNNs. Membership Inference Attack (MIA) is a privacy attack targeting ML models, aiming to distinguish whether a specific data point belongs to the training set [65, 31]. Recently, MIA has been extended to graph learning, where a pioneering work [21] explored the feasibility of membership inference in classical graph embedding models. Subsequently, interest has shifted towards attacking graph neural networks (GNNs), with several impactful and innovative studies revealing GNNs’ privacy vulnerabilities in node classification [30, 29, 58, 103] and graph classification tasks [77], covering cover node-level [30, 58], link-level [29], and graph-level [77, 103] inference risks. Building on this, GroupAttack [91] presents a compelling advancement in link-stealing attacks [29] on GNNs, theoretically demonstrating that different edge groups exhibit varying risk levels and require distinct attack thresholds, while a label-only attack has been proposed to target node-level privacy vulnerabilities [15] with a stricter setting. Another significant line of research involves graph model inversion attacks, which aim to reconstruct the graph structure using model gradients from white-box models [102] or approximated gradients from black-box models [101].

Despite the impressive contributions of previous MIA studies in graph ML models, existing approaches overlook GNNs containing unlearned sensitive knowledge and do not focus on recovering such knowledge from unlearned GNN models. Additional related works are in Appendix A.

3 Preliminaries

In this section, we present the notations used in this paper and give preliminaries on graph unlearning.

Notations. In this paper, bold uppercase letters (e.g., \mathbf{X}) denote matrices, bold lowercase letters (e.g., \mathbf{x}) denote column vectors, and normal letters (e.g., x) indicate scalars. Let $\mathbf{e}_u \in \mathbb{R}^d$ be the column vector with the u -th element as 1 and all others as 0. We use $\|$ to denote concatenating two vectors. We define the weighted inner product with a PSD matrix \mathbf{H} as $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{H}} := \mathbf{x}^\top \mathbf{H} \mathbf{y}$. We use $\mathcal{A} \setminus \mathcal{B} := \{x : x \in \mathcal{A}, x \notin \mathcal{B}\}$ to denote the set difference between sets \mathcal{A} and \mathcal{B} . Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ denote a graph, where $\mathcal{V} = \{v_1, \dots, v_n\}$ is the node set and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ is the edge set. The node feature of node $v_i \in \mathcal{V}$ is denoted by $\mathbf{x}_i \in \mathbb{R}^d$, and the feature matrix for all nodes is denoted by $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$. The adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$ encodes the edge set, where $A_{i,j} = 1$ if $(v_i, v_j) \in \mathcal{E}$, and $A_{i,j} = 0$ otherwise. Specifically, $\mathbf{D} \in \mathbb{R}^{n \times n}$ is the degree matrix, where the diagonal elements $D_{i,i} = \sum_{j=1}^n A_{i,j}$. We use $\mathcal{N}(v)$ to denote the neighborhood of node $v_i \in \mathcal{V}$, and use $\widehat{\mathcal{N}}(v)$ to denote a subset of $\mathcal{N}(v)$. Specifically, $\mathcal{N}^{(k)}(v_i)$ represents all nodes in v_i 's k -hop neighborhood.

Semi-supervised Node Classification. We focus on a semi-supervised node classification task in a transductive setting, which is common in real-world applications [36, 82]. In this setting, the training graph \mathcal{G} includes a small subset of labeled nodes $\mathcal{V}_L = \{u_1, \dots, u_{|\mathcal{V}_L|}\} \subseteq \mathcal{V}$, where each node is annotated with a label $y \in \mathcal{Y}$. The remaining nodes are unlabeled and belong to the subset \mathcal{V}_U , where $\mathcal{V}_U \cap \mathcal{V}_L = \emptyset$. The test set \mathcal{V}_T is a subset of the unlabeled nodes, represented as $\mathcal{V}_T \subseteq \mathcal{V}_U$. We denote the GNN output for a specific target node $v \in \mathcal{V}$ as $f_{\mathcal{G}}(v; \theta)$, where \mathcal{G} is the graph used for neighbor aggregation and θ is the model parameters.

Graph Unlearning. Graph unlearning aims to remove the impact of some undesirable training data from trained GNN models under limited computational overhead [10, 80]. Specifically, consider the original training graph $\mathcal{G}_{\text{orig}} := (\mathcal{V}_{\text{orig}}, \mathcal{E}_{\text{orig}})$ including both desirable data and undesirable data. Normally, the parameters θ_{orig} of the GNN model trained on the original graph is given as:

$$\theta_{\text{orig}} := \arg \min_{\theta} \sum_{v \in \mathcal{V}_L} \mathcal{L}(f_{\mathcal{G}_{\text{orig}}}(v; \theta), y_v), \quad (1)$$

where $\mathcal{V}_L \subseteq \mathcal{V}_{\text{orig}}$ is the set of labeled nodes, and \mathcal{L} is an loss function (e.g., cross-entropy [104]).

Let the undesirable knowledge be a subgraph $\Delta \mathcal{G} := (\Delta \mathcal{V}, \Delta \mathcal{E})$ of the original graph, where $\Delta \mathcal{V} \subseteq \mathcal{V}_{\text{orig}}$ and $\Delta \mathcal{E} \subseteq \mathcal{E}_{\text{orig}}$. The unlearned graph is defined as $\mathcal{G}_{\text{un}} := (\mathcal{V} \setminus \Delta \mathcal{V}, \mathcal{E} \setminus \Delta \mathcal{E})$, which excludes the undesirable knowledge. The goal of graph unlearning is to obtain parameters θ_{un} using an efficient algorithm UNLEARN (e.g., gradient ascent [78, 105], or influence function computation [80, 81]), such that θ_{un} closely approximates the retrained parameters θ_{re} from the cleaned graph \mathcal{G}_{un} , while being significantly more efficient than retraining from scratch. Formally, the unlearning process is defined as:

$$\theta_{\text{un}} := \text{UNLEARN}(f, \theta_{\text{orig}}, \mathcal{G}_{\text{orig}}, \Delta \mathcal{G}) \approx \arg \min_{\theta} \sum_{v \in \mathcal{V}_L / \Delta \mathcal{V}} \mathcal{L}(f_{\mathcal{G}_{\text{orig}} / \Delta \mathcal{G}}(v; \theta), y_v), \quad (2)$$

where the right optimization problem denotes retrain from scratch on the unlearned graph.

We focus on the **edge unlearning** setting [81], where $\Delta \mathcal{V} = \emptyset$, i.e., only edges are removed. This setting captures practical scenarios such as users requesting the removal of private friendship links from social media or the deletion of sensitive purchase records from recommender systems.

4 Problem Formulation

4.1 Threat Model

Attacker's Goal. The adversary aims to recover links in the original training graph $\mathcal{G}_{\text{orig}}$, i.e., $\mathcal{E}_{\text{orig}}$. It includes both the unlearned edges $\Delta \mathcal{E}$ and the remaining membership edges $\mathcal{E}_{\text{orig}} \setminus \Delta \mathcal{E}$. As both types of edges can reveal private user information, with unlearned edges typically being more sensitive, the attacker's goal is to accurately infer both. Specifically, given any pair of nodes $v_i, v_j \in \mathcal{V}$, the attacker aims to determine whether the edge (v_i, v_j) existed in $\mathcal{E}_{\text{orig}}$, i.e., whether $(v_i, v_j) \in \mathcal{E}_{\text{orig}}$. We leave the study of node-level and feature-level unlearning inversion as future work.

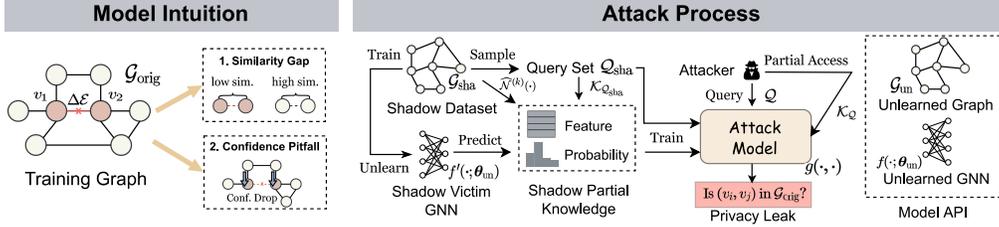


Figure 2: Illustration of the proposed TrendAttack.

Attacker’s Knowledge and Capability. We consider a black-box setting, motivated by the widespread deployment of machine learning models as a service via APIs [65, 68]. The attacker can query the unlearned GNN to obtain output probabilities $f_{\mathcal{G}_{\text{un}}}(\cdot; \theta_{\text{un}})$ for target nodes in \mathcal{V} . Additionally, the attacker has partial access to the unlearned graph \mathcal{G}_{un} . For any pair of target nodes $v_i, v_j \in \mathcal{V}$, the attacker has access to the following information: **(i)** The model output probabilities $f_{\mathcal{G}_{\text{un}}}(v_i; \theta_{\text{un}})$ and $f_{\mathcal{G}_{\text{un}}}(v_j; \theta_{\text{un}})$; **(ii)** The input features \mathbf{x}_i and \mathbf{x}_j ; and **(iii)** A subset of the k -hop neighborhood of v_i and v_j , i.e., $\hat{\mathcal{N}}^{(k)}(v_i)$ and $\hat{\mathcal{N}}^{(k)}(v_j)$. Despite extending the typical probability-only black-box attack setting in (i), this level of access is realistic in practice, since (ii) can be obtained from users’ public profiles on online social platforms, and (iii) can be retrieved by querying public friend lists or social connections. Furthermore, due to the wide availability of public social network datasets [17, 1, 55], we assume that the attacker also has access to a shadow dataset \mathcal{G}^{sha} with a distribution similar to $\mathcal{G}_{\text{orig}}$. This shadow graph can be used to support the unlearning inversion attack.

4.2 Graph Unlearning Inversion

With the threat model presented in Section 4.1, we now formalize the graph unlearning inversion attack under concrete conditions and attack objectives. We begin by defining the query set and the adversary’s partial knowledge.

Definition 4.1 (Query Set). *The query set is a set of Q node pairs of interest, defined as $\mathcal{Q} := \{(v_{i_k}, v_{j_k}) : v_{i_k}, v_{j_k} \in \mathcal{V}\}_{k=1}^Q$.*

Definition 4.2 (Partial Knowledge of Query Set). *Let $f_{\mathcal{G}_{\text{un}}}(\cdot; \theta_{\text{un}})$ be an unlearned GNN model where the parameters θ_{un} are obtained from Eq. (2), and $\Delta\mathcal{E}$ denotes the unlearned edge set. For a given query set \mathcal{Q} , the adversary’s partial knowledge on the unlearned graph \mathcal{G}_{un} is defined as a tuple $\mathcal{K}_{\mathcal{Q}} := (\mathcal{P}_{\mathcal{Q}}, \mathcal{F}_{\mathcal{Q}})$, where: (i) Probability set $\mathcal{P}_{\mathcal{Q}} := \{f_{\mathcal{G}_{\text{un}}}(v_i; \theta_{\text{un}}) : v_i \in \hat{\mathcal{N}}^{(k)}(v_j), v_j \in \mathcal{Q}\}$ represents the model output probabilities for nodes in the k -hop neighborhood of the query nodes; and (ii) Feature set $\mathcal{F}_{\mathcal{Q}} := \{\mathbf{x}_i : v_i \in \mathcal{Q}\}$ contains the input features of nodes in the query pairs.*

With the above definition, the graph unlearning inversion problem is defined as

Definition 4.3 (Graph Unlearning Inversion Problem). *Let \mathcal{Q} be the adversary’s node pairs of interest (Definition 4.1), with partial knowledge $\mathcal{K}_{\mathcal{Q}}$ on the unlearned graph \mathcal{G}_{un} (Definition 4.2). Suppose the adversary also has access to a shadow graph \mathcal{G}^{sha} drawn from a distribution similar to $\mathcal{G}_{\text{orig}}$. The goal of the graph unlearning inversion attack is to predict, for each $(v_i, v_j) \in \mathcal{Q}$, whether $(v_i, v_j) \in \mathcal{E}_{\text{orig}}$ (label 1) or $(v_i, v_j) \notin \mathcal{E}_{\text{orig}}$ (label 0).*

More specifically, the query set \mathcal{Q} can be divided into three disjoint subsets: (i) $\mathcal{Q}_{\text{un}}^+ := \{(v_i, v_j) : (v_i, v_j) \in \Delta\mathcal{E}\}$, representing the positive unlearned edges; (ii) $\mathcal{Q}_{\text{me}}^+ := \{(v_i, v_j) : (v_i, v_j) \in \mathcal{E}_{\text{orig}} \setminus \Delta\mathcal{E}\}$, representing the remaining membership edges in the original graph; and (iii) $\mathcal{Q}^- := \{(v_i, v_j) : (v_i, v_j) \notin \mathcal{E}_{\text{orig}}\}$, denoting non-member (negative) pairs. The goal of this work is to distinguish both $\mathcal{Q}_{\text{un}}^+$ and $\mathcal{Q}_{\text{me}}^+$ from \mathcal{Q}^- , enabling accurate membership inference on both positive subsets. This highlights a key difference between our work and prior MIA methods for GNNs [30, 29, 91], which focus solely on distinguishing $\mathcal{Q}_{\text{me}}^+$ from \mathcal{Q}^- and may fall short of inferring $\mathcal{Q}_{\text{un}}^+$.

5 Proposed Method

In this section, we first present the key motivation behind our method, focusing on the adaptive threshold and confidence pitfalls. We then introduce our proposed TrendAttack framework. An illustration of design motivation and attack process of TrendAttack is in Figure 2.

5.1 Design Motivation

As discussed in Section 4.2, a primary challenge that differentiates graph unlearning inversion attacks from traditional MIA on GNNs is the need to distinguish both unlearned edges in $\mathcal{Q}_{\text{un}}^+$ and other membership edges in $\mathcal{Q}_{\text{mem}}^+$ from non-member node pairs in \mathcal{Q}^- . To address this, we identify two key observations: one explaining why existing MIA methods fall short, and another guiding how to improve upon them using insights from confidence trends.

Probability Similarity Gap. Existing MIA methods infer the presence of a link between v_i and v_j by evaluating the similarity between their input features $\text{sim}(\mathbf{x}_i, \mathbf{x}_j)$ [29] and/or the similarity between their predicted probabilities $\text{sim}(\mathbf{p}_i, \mathbf{p}_j)$ [29, 58], where \mathbf{p}_i is shorthand for the GNN output $f_{\mathcal{G}_{\text{un}}}(v_i; \theta_{\text{un}})$. In the link unlearning setting, since input features remain unchanged after removing $\Delta\mathcal{E}$ from $\mathcal{G}_{\text{orig}}$ and deriving θ_{un} from θ_{orig} , feature similarity $\text{sim}(\mathbf{x}_i, \mathbf{x}_j)$ provides a consistent signal for identifying both unlearned and other membership edges. This suggests that existing feature-based MIA techniques remain applicable.

However, it is essential to carefully examine how probability similarity $\text{sim}(\mathbf{p}_i, \mathbf{p}_j)$ behaves across different edge types. For a removed edge $(v_i, v_j) \in \mathcal{Q}_{\text{un}}^+$, which is excluded from the training graph (i.e., $(v_i, v_j) \notin \mathcal{G}_{\text{un}}$) and unlearned from the parameters θ_{un} , the link’s influence on model predictions tends to be weaker than that of other membership edges in $\mathcal{Q}_{\text{mem}}^+$. Due to approximation errors that are inherent in practical unlearning techniques [12, 80], some residual signal from (v_i, v_j) may still persist in the model’s predictions. These observations lead to the following hypothesis:

Claim 5.1 (Probability Similarity Gap). *The average similarity between predicted probabilities across the three query subsets follows the order:*

$$\text{ProbSim}(\mathcal{Q}^-) < \text{ProbSim}(\mathcal{Q}_{\text{un}}^+) < \text{ProbSim}(\mathcal{Q}_{\text{mem}}^+),$$

where $\text{ProbSim}(\mathcal{Q}_0) := |\mathcal{Q}_0|^{-1} \sum_{(v_i, v_j) \in \mathcal{Q}_0} \text{sim}(\mathbf{p}_i, \mathbf{p}_j)$ denotes the average probability similarity over a query subset \mathcal{Q}_0 .

This claim is empirically verified by our preliminary study in Appendix E.

Confidence Pitfall. Our analysis of the probability similarity gap (Claim 5.1) suggests that unlearned edges $\mathcal{Q}_{\text{un}}^+$ and remaining membership edges $\mathcal{Q}_{\text{mem}}^+$ require different probability-similarity thresholds for accurate inference. The challenge, however, is that the attacker cannot directly tell whether a given query (v_i, v_j) belongs to $\mathcal{Q}_{\text{un}}^+$ or $\mathcal{Q}_{\text{mem}}^+$. To resolve this problem, we present the following intuition:

Claim 5.2 (Confidence Pitfall). *The average model confidence of nodes appearing in unlearned edges is lower than that of other nodes, i.e., $\text{AvgConf}(\mathcal{V}_{\mathcal{Q}_{\text{un}}^+}) < \text{AvgConf}(\mathcal{V} \setminus \mathcal{V}_{\mathcal{Q}_{\text{un}}^+})$, where $\mathcal{V}_{\mathcal{Q}_0} := \{v_i : (v_i, v_j) \in \mathcal{Q}_0 \text{ or } (v_j, v_i) \in \mathcal{Q}_0\}$ is the set of all nodes involved in query subset \mathcal{Q}_0 , and $\text{AvgConf}(\mathcal{V}_0) := \frac{1}{|\mathcal{V}_0|} \sum_{v_i \in \mathcal{V}_0} \max_{\ell} \mathbf{p}_{i, \ell}$ is the average model confidence within a node set \mathcal{V}_0 .*

For this claim, we provide a detailed theoretical analysis in Appendix C, and then show the empirical evidence and preliminary experiments supporting the claim in Appendix E.

5.2 The Proposed TrendAttack

Based on the two key design principles established in Claim 5.1 and Claim 5.2 for enabling graph unlearning inversion, we now instantiate these principles through concrete model components, resulting in a simple, flexible, and effective inversion framework, named TrendAttack.

Shadow Victim Model Training. To construct an attack model $g(v_i, v_j)$ that accurately predicts membership information given partial knowledge of the unlearned graph $\mathcal{K}_{\mathcal{Q}} = (\mathcal{P}_{\mathcal{Q}}, \mathcal{F}_{\mathcal{Q}})$, we first simulate the victim model’s behavior using a shadow dataset \mathcal{G}^{sha} and a shadow victim GNN $f'_{\mathcal{G}_{\text{un}}^{\text{sha}}}(\cdot; \theta'_{\text{un}})$. This model is trained on a node classification task and then unlearns a small subset of edges $\Delta\mathcal{G}^{\text{sha}}$. Specifically, the pre- and post-unlearning parameters are obtained similarly via Eq. (1) and Eq (2) in Section 3.

This training setup explicitly models unlearning behavior, in contrast to prior MIA approaches that rely solely on the original model $f'_{\mathcal{G}_{\text{orig}}^{\text{sha}}}(\cdot; \theta'_{\text{orig}})$ and do not account for the effects of unlearning. Using this shadow victim model, we construct a shadow query set \mathcal{Q}_{sha} and its associated partial

knowledge $\mathcal{K}_{\mathcal{Q}_{\text{sha}}} = (\mathcal{P}_{\mathcal{Q}_{\text{sha}}}, \mathcal{F}_{\mathcal{Q}_{\text{sha}}})$ from the features, connectivity, and outputs of f' on \mathcal{G}^{sha} , which are then used to train the attack model.

Attack Model. It is well established that membership between v_i and v_j can be inferred from the similarity between their features and output probabilities [58, 29]. Since prior MIA methods have developed a variety of similarity computation frameworks for this task, we adopt a general formulation that computes a scalar similarity score between v_i and v_j as $\phi([\mathbf{x}_i \parallel \mathbf{p}_i], [\mathbf{x}_j \parallel \mathbf{p}_j])$.

This formulation is flexible and covers several existing MIA methods. For example, when $\phi([\mathbf{x}_i \parallel \mathbf{p}_i], [\mathbf{x}_j \parallel \mathbf{p}_j]) = \mathbf{h}^\top \cdot \text{MLP}(\mathbf{p}_i, \mathbf{p}_j)$, the model recovers MIA-GNN [58]. It can also recover the StealLink attack [29] by incorporating manually defined similarity features into ϕ . This flexible structure allows our attack model to incorporate any existing MIA method as the backbone, ensuring it performs at least as well as prior approaches.

In addition, as indicated by Claim 5.1 and Claim 5.2, a key challenge in graph unlearning inversion is to distinguish nodes associated with unlearned edges from others and to apply an adaptive similarity threshold. Therefore, relying solely on $\phi(\cdot, \cdot)$ may be insufficient to capture this complexity. Thus, to explicitly capture confidence trends and address Claim 5.2, we define scalar-valued confidence trend features for each node $v_i \in \mathcal{V}$:

$$\tau_i^{(0)} := \text{Conf}(v_i), \quad \tau_i^{(k)} := \sum_{v_j \in \mathcal{N}^{(1)}(v_i)} \tilde{A}_{i,j} \tau_i^{(k-1)} \quad (k \geq 1), \quad (3)$$

where $\tilde{\mathbf{A}} := \mathbf{D}^{-0.5} \mathbf{A} \mathbf{D}^{-0.5}$ is the normalized adjacency matrix and $\tilde{A}_{i,j}$ denotes its (i, j) -th entry.

We define the confidence difference between orders as $\Delta \tau_i^{(k)} := \tau_i^{(k)} - \tau_i^{(k-1)}$ for $k \geq 1$. Based on Claim 5.2, the sign of these differences (e.g., between zeroth and first order, and first and second order) serves as a useful signal for identifying nodes that are endpoints of unlearned edges. Thus, we define the following binary-valued trend feature for each node v_i :

$$\tilde{\tau}_i := \left[\mathbf{1}\{\Delta \tau_i^{(1)} < 0\}, \mathbf{1}\{\Delta \tau_i^{(1)} > 0\}, \mathbf{1}\{\Delta \tau_i^{(2)} < 0\}, \mathbf{1}\{\Delta \tau_i^{(2)} > 0\} \right]. \quad (4)$$

The final attack model is defined as:

$$g(v_i, v_j) := \underbrace{\sigma(\phi([\mathbf{x}_i \parallel \mathbf{p}_i], [\mathbf{x}_j \parallel \mathbf{p}_j]))}_{\text{MIA}} + \underbrace{\mathbf{h}^\top [\tilde{\tau}_i \parallel \tilde{\tau}_j]}_{\text{Trend}}. \quad (5)$$

In the equation above, the MIA term estimates membership from similarity in features and probabilities, and the trend term compensates for the similarity gap by distinguishing node types as indicated by Claim 5.1. The sigmoid function $\sigma(\cdot)$ maps the output to $[0, 1]$, with 0 indicating non-member and 1 indicating member. This design allows the attack model to incorporate an adaptive threshold based on node-specific trend features, while leveraging any existing MIA framework as its base.

Unlearning-Aware Attack Model Training. We train the attack model g on the shadow dataset \mathcal{G}^{sha} using the outputs of the shadow victim model f' . Given the shadow query set \mathcal{Q}_{sha} with known membership labels, we optimize a link-prediction loss:

$$\mathcal{L}_{\text{attack}}(\mathcal{Q}_{\text{sha}}) := - \sum_{(v_i, v_j) \in \mathcal{Q}_{\text{sha}}^+} \log g(v_i, v_j) - \sum_{(v_i, v_j) \in \mathcal{Q}_{\text{sha}}^-} \log(1 - g(v_i, v_j)). \quad (6)$$

Performing TrendAttack. After training, we transfer g to the target unlearned graph \mathcal{G}_{un} by computing, for each query pair $(v_i, v_j) \in \mathcal{Q}$, the feature-probability similarity and trend features from the real model output $f_{\mathcal{G}_{\text{un}}}(\cdot; \boldsymbol{\theta}_{\text{un}})$ and partial graph knowledge \mathcal{K} . We then evaluate $g(v_i, v_j)$ on predicting whether $(v_i, v_j) \in \mathcal{E}_{\text{orig}}$. By combining both similarity and trend signals learned on the shadow graph, g effectively supports unlearning inversion on the real unlearned graph.

6 Experiments

In this section, we describe our experimental setup and present the main empirical results.

6.1 Experiment Settings

Datasets. We evaluate our attack method on four standard graph ML benchmark datasets: Cora [87], Citeseer [87], Pubmed [87], and LastFM-Asia [87]. To construct the shadow dataset \mathcal{G}_{sha} and the real

Unlearn method	Attack	Cora			Citeseer			Pubmed			LastFM-Asia		
		Unlearned	Original	All									
GIF	GraphSAGE	0.5356	0.5484	0.5420	0.5275	0.5216	0.5246	0.6503	0.6457	0.6480	0.6914	0.6853	0.6884
	NCN	0.7403	0.7405	0.7404	0.6750	0.6872	0.6811	0.6661	0.6718	0.6690	0.7283	0.7273	0.7278
	MIA-GNN	0.7547	0.7916	0.7732	0.7802	0.8245	0.8023	0.7028	0.7902	0.7465	0.5955	0.5744	0.5850
	StealLink	0.7841	0.8289	0.8065	0.7369	0.8404	0.7887	0.8248	0.8964	0.8606	<u>0.8472</u>	<u>0.9037</u>	<u>0.8755</u>
	GroupAttack	0.7982	0.8053	0.8018	0.7771	0.7618	0.7695	0.6497	0.6554	0.6525	0.7858	0.7850	0.7854
	TrendAttack-MIA	<u>0.8240</u>	<u>0.8448</u>	<u>0.8344</u>	<u>0.8069</u>	<u>0.8078</u>	<u>0.8073</u>	<u>0.8950</u>	<u>0.9171</u>	<u>0.9060</u>	0.7795	0.7649	0.7722
TrendAttack-SL	0.8309	0.8527	0.8418	0.8410	0.8430	0.8420	0.9524	0.9535	0.9529	0.9078	0.9134	0.9106	
CEU	GraphSAGE	0.5356	0.5484	0.5420	0.5275	0.5216	0.5246	0.6503	0.6457	0.6480	0.6914	0.6853	0.6884
	NCN	0.7403	0.7405	0.7404	0.6750	0.6872	0.6811	0.6661	0.6718	0.6690	0.7283	0.7273	0.7278
	MIA-GNN	0.7458	0.7810	0.7634	0.7718	0.8248	0.7983	0.6626	0.6561	0.6593	0.6004	0.5811	0.5908
	StealLink	0.7901	<u>0.8486</u>	0.8193	0.7643	0.8450	0.8046	0.8467	0.9088	0.8777	<u>0.8416</u>	<u>0.9021</u>	<u>0.8719</u>
	GroupAttack	0.7941	0.7976	0.7958	0.7557	0.7458	0.7508	0.6388	0.6430	0.6409	0.7845	0.7817	0.7831
	TrendAttack-MIA	0.8194	0.8333	0.8263	0.7933	0.8041	0.7987	<u>0.8982</u>	<u>0.9184</u>	<u>0.9083</u>	0.7676	0.7576	0.7626
TrendAttack-SL	0.8467	0.8612	0.8539	0.8514	<u>0.8400</u>	0.8457	0.9550	0.9579	0.9565	0.9037	0.9088	0.9062	
GA	GraphSAGE	0.5356	0.5484	0.5420	0.5275	0.5216	0.5246	0.6503	0.6457	0.6480	0.6914	0.6853	0.6884
	NCN	0.7403	0.7405	0.7404	0.6750	0.6872	0.6811	0.6661	0.6718	0.6690	0.7283	0.7273	0.7278
	MIA-GNN	0.7676	0.8068	0.7872	0.7798	0.8353	0.8076	0.7242	0.8039	0.7641	0.6200	0.6057	0.6129
	StealLink	0.7862	0.8301	0.8082	0.7479	<u>0.8431</u>	0.7955	0.8203	0.8898	0.8550	<u>0.8342</u>	<u>0.8947</u>	<u>0.8644</u>
	GroupAttack	0.7945	0.8042	0.7993	0.7662	0.7563	0.7613	0.6458	0.6493	0.6475	0.7746	0.7760	0.7753
	TrendAttack-MIA	<u>0.8193</u>	0.8397	<u>0.8295</u>	<u>0.8080</u>	0.8249	<u>0.8165</u>	<u>0.8932</u>	<u>0.9158</u>	<u>0.9045</u>	0.7255	0.7099	0.7177
TrendAttack-SL	0.8270	<u>0.8382</u>	0.8326	0.8628	0.8614	0.8621	0.9531	0.9537	0.9534	0.9041	0.9119	0.9080	

Table 1: **Main Comparison Results.** We present the AUC scores for attack methods across different edge groups. The best results are highlighted in **bold**, while the second-best results are underlined.

attack dataset $\mathcal{Q}_{\text{orig}}$, which should share similar distributions, we use METIS to partition the entire training graph into two balanced subgraphs, one for shadow and one for attack. Following the setting in GIF [80], we use 90% of the nodes in each subgraph for training and the remaining for testing. All edges between the two subgraphs are removed, and there are no shared nodes, simulating real-world scenarios where shadow and attack datasets are disconnected.

Victim Model. We adopt a two-layer GCN [36] as the victim model, trained for node classification. Our GCN implementation follows the standard settings used in GIF [80] for consistency and reproducibility. After training the GCN on both the shadow and attack datasets, we perform edge unlearning on 5% of randomly selected edges using standard graph unlearning methods, including GIF [80], CEU [81], and Gradient Ascent (GA) [78]. We follow the official settings from each method’s paper and codebase to ensure faithful reproduction.

Baselines. We do not compare with unlearning inversion attacks [32] that require access to pre-unlearning models, which is unrealistic under our setting. To demonstrate that unlearning inversion cannot be solved by naive link prediction, we evaluate a simple GraphSAGE model and a state-of-the-art link prediction method, NCN [74]. For membership inference attacks (MIAs) under the same black-box assumption as ours, we consider three widely used methods: StealLink [29], MIA-GNN [58], and GroupAttack [91]. We follow their official hyperparameter settings from their respective papers and repositories. All experiments are run five times, and we report the mean and the standard error.

Evaluation Metrics. We evaluate our attack on the attack dataset using a specific query set \mathcal{Q} . We randomly select 5% of the edges as unlearned edges $\mathcal{Q}_{\text{un}}^+$ (label 1), and another 5% of remaining edges as regular member edges $\mathcal{Q}_{\text{mem}}^+$ (label 1). We then sample an equal number of non-existent (negative) edges \mathcal{Q}^- (label 0) such that $|\mathcal{Q}^-| = |\mathcal{Q}_{\text{un}}^+| + |\mathcal{Q}_{\text{mem}}^+|$. We use AUC as the primary evaluation metric. Since both unlearned edges and regular member edges are important for membership inference, we compute the overall AUC on the full query set $\mathcal{Q} = \mathcal{Q}_{\text{un}}^+ \cup \mathcal{Q}_{\text{mem}}^+ \cup \mathcal{Q}^-$.

More experimental details, including model parameters, baselines, and datasets, are in Appendix F.

6.2 Comparison Experiments

In this study, we present a comprehensive comparison of all our baselines mentioned in Section 6.1, and the results are shown in Table 1. Specifically, to demonstrate that our method better captures unlearned edges, we evaluate AUC across three groups: *Unlearned* ($\mathcal{Q}_{\text{un}}^+ \cup \mathcal{Q}^-$), *Original* ($\mathcal{Q}_{\text{mem}}^+ \cup \mathcal{Q}^-$), and *All* (the entire \mathcal{Q}). A small gap between Unlearned and Original AUCs indicates a better balance in the model’s predictions. The full table with standard deviation can be found in Appendix G.

We consider two variants of our method, TrendAttack-MIA and TrendAttack-SL, which adopt MIA-GNN [58] and StealLink [29] as their respective backbone models. From the table, we make the following observations: (i) Compared with their MIA prototypes, both variants of TrendAttack significantly improve the gap between Unlearned and Original AUCs, as well as the overall AUC. This demonstrates the effectiveness of our proposed trend-based attack and unlearning-aware training

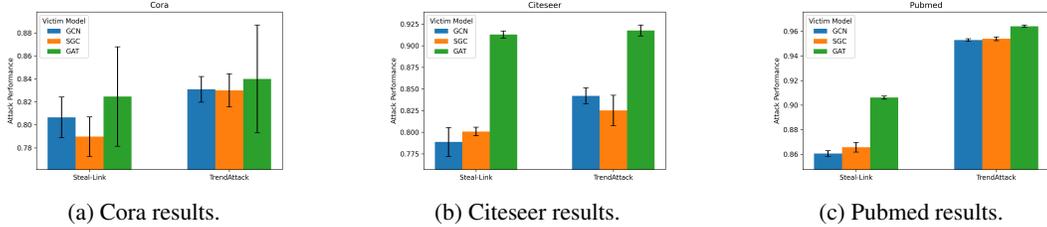


Figure 3: Ablation study on the impact of victim models.

design. (ii) Overall, our proposed attack, especially TrendAttack-SL, achieves the best performance on most datasets and unlearning methods. The only failure case is on Citeseer + CEU + Original, where the gap to StealLink is small. This does not undermine our contribution, as our primary focus is on the unlearned sets, and this case is a special exception. TrendAttack-MIA achieves the second-best results across many settings, while its relatively lower performance is mainly due to the weak backbone model (MIA-GNN), not the attack framework itself. (iii) Among all baselines, attack-based methods consistently outperform link prediction methods, showing that link prediction alone cannot effectively solve the unlearning inversion problem. Among the attack baselines, StealLink is the strongest, while MIA-GNN performs poorly as it only uses output probabilities and ignores features.

6.3 Ablation Studies

Impact of Victim Models. From Table 1, we observe that the proposed attack remains stable against unlearning methods. In this study, we further investigate whether TrendAttack’s performance maintains stability with respect to changes in victim models. Specifically, we fix the unlearning method to GIF and use TrendAttack-SL as our model variant. We compare it against the best-performing baseline, StealLink, with results shown in Figure 3. From the figure, we see that our proposed model consistently outperforms the baseline, demonstrating stability across different victim models. An interesting observation is that performance significantly improves on GAT compared to other baselines, suggesting that GAT may be more vulnerable to model attacks.

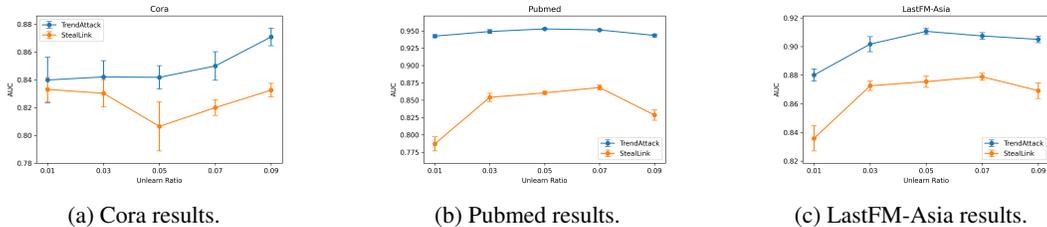


Figure 4: Ablation study on the impact of unlearning ratio.

Impact of Unlearning Ratio. In this study, we examine whether the victim model’s randomly selected unlearning edge ratio affects the overall attack AUC. Specifically, we fix the unlearning method to GIF and use TrendAttack-SL as our model variant, with results presented in Figure 4. From the figure, we find that compared to our most important baseline, StealLink, TrendAttack consistently performs better with less fluctuation. This indicates that our model remains relatively stable across different unlearning ratios, making it more universally applicable to various attack settings.

Due to space limitations, more experiments can be found in Appendix G.

7 Conclusion

In this work, we study a novel and challenging problem of graph unlearning inversion attack, which aims to recover unlearned edges using black-box GNN outputs and partial knowledge of the unlearned graph. To tackle this problem, we identify two key intuitions, probability gap and confidence pitfall, which motivate a simple yet effective attack framework, TrendAttack, revealing a potential privacy risk of current graph unlearning methods. Despite the effectiveness of TrendAttack, it has several limitations. First, our method focuses solely on link-level inference, and does not consider to node-level or feature-level membership inference. Second, while we provide empirical evidence for the probability gap intuition in the appendix, we do not offer a formal theoretical justification. These limitations highlight important directions for future research.

Appendix

List of Contents

In this appendix, we provide the following additional information:

- **Section A:** Additional Related Works.
- **Section B:** Model Details.
- **Section C:** Theoretical Motivation of the Confidence Pitfall.
- **Section D:** Missing Proofs in Section C.
- **Section E:** Preliminary Experiments.
- **Section F:** Experimental Settings.
- **Section G:** Additional Experiments.
- **Section H:** Impact Statement.

A Additional Related Works

In this section, we present additional related works for this paper. We first review prior works on general-purpose machine unlearning and the privacy vulnerabilities of unlearning. Next, we provide a comprehensive list of related works on graph unlearning and membership inference attacks (MIA) on GNNs, supplementing our earlier discussion in Section 2.

A.1 Machine Unlearning

Machine unlearning aims to remove the influence of specific training samples from a trained model, balancing utility, removal guarantees, and efficiency [4, 50]. Unlike full retraining, unlearning seeks practical alternatives to efficiently revoke data. These methods can help removing data effect in simple learning settings like statistical query learning [6], and recently have broad applications in LLMs [88, 47, 49], generative models [40, 99, 100], e-commerce [7, 43, 98], and graph learning [10, 80, 19].

Existing machine unlearning methods broadly fall into two categories: exact and approximate unlearning. Exact unlearning methods, such as SISA [4], partition data into shards, train sub-models in each shard independently, and merge them together, which allows targeted retraining for removing certain data from a shard. ARCANE [85] improves this by framing unlearning as one-class classification and caching intermediate states, enhancing retraining efficiency. These ideas have been extended to other domains, including graph learning [10, 70] and ensemble methods [5]. However, exact methods often degrade performance due to weak sub-models and straightforward ensembling.

Approximate unlearning methods update model parameters to emulate removal effects with better efficiency-performance trade-offs. For instance, Jia et al. [33] uses pruning and sparsity-regularized fine-tuning to approximate exact removal. Tarun et al. [66] propose a model-agnostic class removal technique using error-maximizing noise and a repair phase. Liu et al. [48] tackle adversarially trained models with a closed-form Hessian-guided update, approximated efficiently without explicit inversion.

A.2 Privacy Vulnerabilities of Machine Unlearning

A growing line of research investigates the privacy vulnerabilities of machine unlearning. Chen et al. [9] show that if an adversary has access to both the pre- and post-unlearning black-box model outputs, they can infer the membership status of a target sample. A recent work [67] highlights the limitations of approximate unlearning by arguing that unlearning is only well-defined at the algorithmic level, and parameter-level manipulations alone may be insufficient. Unlearning inversion attacks [32] go further by reconstructing both features and labels of training samples using two model versions, while Bertran et al. [3] provide an in-depth analysis of such reconstruction attacks in regression settings.

These findings motivate our investigation into unlearned GNNs, where we aim to capture residual signals left by the unlearned data. While prior attacks mostly target general-purpose models and focus on red-teaming under strong assumptions, such as access to both pre- and post-unlearning models [9, 32, 3], our work considers a more practical black-box setting motivated by real-world social network scenarios. Furthermore, our attack explicitly accounts for the unique structural and relational properties of graphs, in contrast to prior work primarily focused on i.i.d. data.

A.3 Graph Unlearning

Graph Unlearning enables the efficient removal of unwanted data’s influence from trained graph ML models [10, 61, 22]. This removal process balances model utility, unlearning efficiency, and removal guarantees, following two major lines of research: retrain-based unlearning and approximate unlearning.

Retrain-based unlearning partitions the original training graph into disjoint subgraphs, training independent submodels on them, enabling unlearning through retraining on a smaller subset of the data. Specifically, GraphEraser [10] pioneered the first retraining-based unlearning framework for GNNs, utilizing balanced clustering methods for subgraph partitioning and ensembling submodels in prediction with trainable fusion weights to enhance model utility. GUIDE [70] presents an important follow-up by extending the retraining-based unlearning paradigm to inductive graph learning settings. Subsequently, many studies [70, 43, 93, 42] have made significant contributions to improving the Pareto front of utility and efficiency in these methods, employing techniques such as data condensation [42] and enhanced clustering [43, 93].

Approximate unlearning efficiently updates model parameters to remove unwanted data. Certified graph unlearning [12] provides an important early exploration of approximate unlearning in SGC [79], with provable unlearning guarantees. GraphGuard [78] introduces a comprehensive system to mitigate training data misuse in GNNs, featuring a significant gradient ascent unlearning method as one of its core components. These gradient ascent optimization objectives can be approximated using Taylor expansions, inspiring several beautiful influence function-based methods [37, 80, 81]. Specifically, GIF [80] presents a novel influence function-based unlearning approach tailored to graph data, considering feature, node, and edge unlearning settings, while CEU [81] focuses on edge unlearning and establishes theoretical bounds for removal guarantees. Building on this, IDEA [19] proposes a theoretically sound and flexible approximate unlearning framework, offering removal guarantees across node, edge, and feature unlearning settings. Recent innovative works have further advanced the scalability [59, 39, 89, 86, 92] and model utility [41, 97] of approximate unlearning methods, pushing the boundaries of promising applications. Additionally, a noteworthy contribution is GNNDelete [11], which unlearns graph knowledge by employing an intermediate node embedding mask between GNN layers. This method offers a unique solution that differs from both retraining and model parameter modification approaches.

In this paper, we explore the privacy vulnerabilities of graph unlearning by proposing a novel membership inference attack tailored to unlearned GNN models, introducing a new defence frontier that graph unlearning should consider from a security perspective.

A.4 Membership Inference Attack for GNNs

Membership Inference Attack (MIA) is a privacy attack targeting ML models, aiming to distinguish whether a specific data point belongs to the training set [65, 31, 16]. Such attacks have profoundly influenced ML across various downstream applications, including computer vision [8, 14], NLP [53, 54], and recommender systems [95, 90]. Recently, MIA has been extended to graph learning, where a pioneering work [21] explored the feasibility of membership inference in classical graph embedding models. Subsequently, interest has shifted towards attacking graph neural networks (GNNs), with several impactful and innovative studies revealing GNNs’ privacy vulnerabilities in node classification [30, 29, 58, 103] and graph classification tasks [77, 44], covering cover node-level [30, 58], link-level [29], and graph-level [77, 103] inference risks. Building on this, GroupAttack [91] presents a compelling advancement in link-stealing attacks [29] on GNNs, theoretically demonstrating that different edge groups exhibit varying risk levels and require distinct attack thresholds, while a label-only attack has been proposed to target node-level privacy vulnerabilities [15] with a stricter setting. Another significant line of research involves graph model inversion attacks, which aim to

reconstruct the graph structure using model gradients from white-box models [102] or approximated gradients from black-box models [101].

Despite the impressive contributions of previous MIA studies in graph ML models, existing approaches overlook GNNs containing unlearned sensitive knowledge and do not focus on recovering such knowledge from unlearned GNN models.

B Model Details

In this section, we present the detailed computation of the TrendAttack, considering both the training and attack processes.

B.1 Shadow Training Algorithm for TrendAttack

Algorithm 1 Attack Model Training

Input: Shadow dataset $\mathcal{G}_{\text{orig}}^{\text{sha}}$, GNN architecture f'
Output: Attack model $g(\cdot, \cdot)$

- 1: // *Train and unlearn the victim model*
- 2: Train the shadow victim model on $\mathcal{G}_{\text{orig}}^{\text{sha}}$ with Eq. (1) to obtain θ'_{orig} ▷ Train victim model
- 3: Randomly select some unlearned edges $\Delta\mathcal{E}^{\text{sha}}$ from original edges $\mathcal{E}_{\text{orig}}^{\text{sha}}$
- 4: Unlearn $\Delta\mathcal{E}^{\text{sha}}$ with Eq. (2) to obtain θ'_{un} ▷ Unlearn victim model
- 5: $\mathcal{E}_{\text{un}}^{\text{sha}} \leftarrow \mathcal{E}_{\text{orig}}^{\text{sha}} \setminus \Delta\mathcal{E}^{\text{sha}}$ ▷ Remove the unlearned edges from $\mathcal{G}_{\text{orig}}^{\text{sha}}$
- 6: // *Construct the query set \mathcal{Q}_{sha}*
- 7: Randomly select some existing edges from $\mathcal{E}_{\text{un}}^{\text{sha}}$ to obtain $\mathcal{Q}_{\text{mem}}^+$ ▷ $\mathcal{Q}_{\text{mem}}^+ \subseteq \mathcal{E}_{\text{un}}^{\text{sha}}$
- 8: Randomly select some negative edges \mathcal{Q}^- ▷ $\mathcal{Q}^- \cap \mathcal{E}_{\text{un}}^{\text{sha}} = \emptyset$
- 9: Initialize the unlearned set of edges $\mathcal{Q}_{\text{un}}^+ \leftarrow \Delta\mathcal{E}^{\text{sha}}$
- 10: Initialize the entire query set $\mathcal{Q}_{\text{sha}} \leftarrow \mathcal{Q}_{\text{un}}^+ \cup \mathcal{Q}_{\text{mem}}^+ \cup \mathcal{Q}^-$
- 11: // *Attack model training*
- 12: Obtain the partial knowledge for the query set $\mathcal{K}_{\mathcal{Q}_{\text{sha}}} = (\mathcal{P}_{\mathcal{Q}_{\text{sha}}}, \mathcal{F}_{\mathcal{Q}_{\text{sha}}})$
- 13: Train the attack model g on \mathcal{Q}_{sha} with $\mathcal{K}_{\mathcal{Q}_{\text{sha}}}$ and $\mathcal{L}_{\text{attack}}$ in Eq. (6)
- 14: **return** Attack model g

In this algorithm, we first train and unlearn the shadow victim model f' (lines 1–5), and then construct the query set \mathcal{Q}_{sha} (lines 6–10), which includes three different types of edges. Next, we train the attack model with the link prediction objective $\mathcal{L}_{\text{attack}}$ in Eq. (6), with \mathcal{Q}_{sha} as the pairwise training data (lines 11–13). Then, we return the attack model g for future attacks (line 14).

B.2 The Attack Process of TrendAttack

In this algorithm, we first request the partial knowledge $\mathcal{K}_{\mathcal{Q}}$ from the unlearned graph \mathcal{G}_{un} and the black-box victim model $f_{\mathcal{G}_{\text{un}}}(\cdot; \theta_{\text{un}})$ for our links of interest \mathcal{Q} (lines 1–11). This process is highly flexible and can effortlessly incorporate many different levels of the attacker’s knowledge. For instance, the feature knowledge $\mathcal{F}_{\mathcal{Q}}$ (line 5) is optional if the model API owner does not respond to node feature requests on the unlearned graph \mathcal{G}_{un} .

Moreover, for the probability knowledge $\mathcal{P}_{\mathcal{Q}}$ (lines 6–10), the more we know about the neighborhood $\hat{\mathcal{N}}^{(0)}, \dots, \hat{\mathcal{N}}^{(k)}$ of nodes of interest $\mathcal{V}_{\mathcal{Q}}$, the more accurately we can construct our trend features. This approach is fully adaptive to any level of access to the unlearned graph and model API. In the strictest setting, where we only have access to the node of interest itself, we have $k = 0$, and our model perfectly recovers previous MIA methods with no performance loss or additional knowledge requirements. An empirical study on the impact of trend feature orders can be found in Figure 7.

After requesting the partial knowledge, we compute the trend features and then predict the membership information (lines 12–18). In the end, we return our attack results (line 19).

Algorithm 2 Attack Process

Input: Unlearned graph \mathcal{G}_{un} , black-box unlearned model $f_{\mathcal{G}_{\text{un}}}(\cdot; \boldsymbol{\theta}_{\text{un}})$, query set \mathcal{Q} , trained attack model $g(\cdot, \cdot)$, trend feature order k

Output: Membership predictions on the query set $\hat{\mathbf{y}}$

- 1: // Request the partial knowledge $\mathcal{K}_{\mathcal{Q}} = (\mathcal{P}_{\mathcal{Q}}, \mathcal{F}_{\mathcal{Q}})$
- 2: Obtain the nodes of interest $\mathcal{V}_{\mathcal{Q}} \leftarrow \{v_i : (v_i, v_j) \in \mathcal{Q} \text{ or } (v_j, v_i) \in \mathcal{Q}\}$
- 3: $\mathcal{P}_{\mathcal{Q}} \leftarrow \emptyset, \mathcal{F}_{\mathcal{Q}} \leftarrow \emptyset$
- 4: **for** $v_i \in \mathcal{V}_{\mathcal{Q}}$ **do**
- 5: $\mathcal{F}_{\mathcal{Q}} \leftarrow \mathcal{F}_{\mathcal{Q}} \cup \{(v_i, \mathbf{x}_i)\}$ ▷ Feature knowledge
- 6: Request the available neighborhood $\hat{\mathcal{N}}^{(0)}(v_i), \dots, \hat{\mathcal{N}}^{(k)}(v_i)$ from \mathcal{G}_{un}
- 7: **for** $v_j \in \bigcup_{r=0}^k \hat{\mathcal{N}}^{(r)}(v_i)$ **do** ▷ Probability knowledge
- 8: $\mathbf{p}_j \leftarrow f_{\mathcal{G}_{\text{un}}}(v_j; \boldsymbol{\theta}_{\text{un}})$
- 9: $\mathcal{P}_{\mathcal{Q}} \leftarrow \mathcal{P}_{\mathcal{Q}} \cup \{(v_j, \mathbf{p}_j)\}$
- 10: **end for**
- 11: **end for**
- 12: // Membership inference
- 13: **for** $v_i \in \mathcal{V}_{\mathcal{Q}}$ **do** ▷ Build trend features
- 14: Compute trend features $\tilde{\boldsymbol{\tau}}_i$ for node v_i with Eq. (3) and Eq. (4)
- 15: **end for**
- 16: **for** $(v_i, v_j) \in \mathcal{Q}$ **do** ▷ Predict with attack model g
- 17: Compute the model prediction $\hat{y}_{i,j} \leftarrow g(v_i, v_j)$ with Eq. (5)
- 18: **end for**
- 19: **return** Membership predictions $\hat{\mathbf{y}}$

C Theoretical Motivation of the Confidence Pitfall

In this section, we present our main theoretical results, offering a robust theoretical foundation for our model design, with a specific focus on the confidence pitfall intuition in Section 5.1.

Our analysis of the probability similarity gap (Claim 5.1) suggests that unlearned edges $\mathcal{Q}_{\text{un}}^+$ and remaining membership edges $\mathcal{Q}_{\text{mem}}^+$ require different probability-similarity thresholds for accurate inference. The challenge, however, is that the attacker cannot directly tell whether a given query (v_i, v_j) belongs to $\mathcal{Q}_{\text{un}}^+$ or $\mathcal{Q}_{\text{mem}}^+$.

To address this, we analyze how removing a specific edge (v_i, v_j) affects model outputs, using the widely used analytical framework of influence functions [37, 80, 81]. This analysis separates two effects: (i) the immediate impact of dropping the edge from the training graph $\mathcal{G}_{\text{orig}}$, and (ii) the subsequent adjustment of model parameters by the unlearning procedure. We focus on a linear GCN model, which, despite its simplicity, captures the core behavior of many GNN architectures (see Remark D.2 in Appendix D).

Specifically, we first present the influence result of all the unlearned edges $\Delta\mathcal{E}$ as follows:

Theorem C.1 (Closed-form Edge Influence, Informal). *Let $f(\mathbf{C}, \mathbf{X}; \mathbf{w}^*)$ be a linear GCN with propagation matrix $\mathbf{C} \in \mathbb{R}^{n \times n}$ and parameters \mathbf{w}^* obtained by least-squares on labels $\mathbf{y} \in \mathbb{R}^d$. If we denote the matrix of all node predictions by \mathbf{P} and let $\boldsymbol{\Xi} := \sum_{(v_i, v_j) \in \Delta\mathcal{E}} \mathbf{e}_i \mathbf{e}_j^\top$ be the adjacency matrix of the removed edges, then unlearning $\Delta\mathcal{E}$ produces:*

- *Weight influence:*

$$\mathcal{I}(\mathbf{w}^*) = \mathbf{H}^{-1}(\mathbf{X}^\top \boldsymbol{\Xi}^\top \mathbf{y} - \mathbf{X}^\top \boldsymbol{\Xi}^\top \mathbf{C} \mathbf{X} \mathbf{w}^* - \mathbf{X}^\top \mathbf{C}^\top \boldsymbol{\Xi} \mathbf{X} \mathbf{w}^*),$$

- *Model output influence:*

$$\mathcal{I}(\mathbf{P}) = \boldsymbol{\Xi} \mathbf{X} \mathbf{w}^* + \mathbf{C} \mathbf{X} \mathbf{H}^{-1}(\mathbf{X}^\top \boldsymbol{\Xi}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{C}^\top \boldsymbol{\Xi} \mathbf{X} \mathbf{w}^* - \mathbf{X}^\top \boldsymbol{\Xi}^\top \mathbf{C} \mathbf{X} \mathbf{w}^*),$$

where \mathbf{H} is the Hessian of the least squares loss evaluated at \mathbf{w}^* .

Proof. This follows directly from Theorem D.8 and Theorem D.9. See Appendix D for proofs. \square

While Theorem C.1 characterizes the combined effect of all unlearned edges, we also need to understand the effect of a single edge on a single node’s output. This finer-grained analysis reveals why nodes adjacent to unlearned edges exhibit a distinctive drop in confidence compared to more distant nodes.

Corollary C.2 (Single-Edge to Single-Node Influence, Informal). *Under the same setup as Theorem C.1, the influence of an undirected edge $(v_i, v_j) \in \Delta\mathcal{E}$ on the output \mathbf{p}_k of node $v_k \in \mathcal{V}$ can be decomposed as follows:*

$$\begin{aligned} \mathcal{I}(\mathbf{p}_k) = & \underbrace{\mathbf{1}\{v_k = v_i\} \cdot (\mathbf{x}_i^\top \mathbf{w}^*) + \mathbf{1}\{v_k = v_j\} \cdot (\mathbf{x}_j^\top \mathbf{w}^*)}_{\text{edge influence}} \\ & - \underbrace{\langle (\mathbf{x}_j^\top \mathbf{w}^*)\mathbf{z}_i + (\mathbf{x}_i^\top \mathbf{w}^*)\mathbf{z}_j, \mathbf{z}_k \rangle_{\mathbf{H}^{-1}}}_{\text{magnitude weight influence}} + \underbrace{\langle (y_j - \mathbf{z}_j^\top \mathbf{w}^*)\mathbf{x}_i + (y_i - \mathbf{z}_i^\top \mathbf{w}^*)\mathbf{x}_j, \mathbf{z}_k \rangle_{\mathbf{H}^{-1}}}_{\text{error weight influence}}, \end{aligned}$$

where $\mathbf{Z} := \mathbf{C}\mathbf{X}$ and \mathbf{z}_l is the l -th row of \mathbf{Z} , and $\mathbf{1}\{\cdot\}$ denotes the indicator function.

Proof. Please see Corollary D.11 in Appendix D. □

In the corollary above, the third error weight influence term is governed by the empirical residuals at the optimal weight, $(y_i - \mathbf{z}_i^\top \mathbf{w}^*)$ and $(y_j - \mathbf{z}_j^\top \mathbf{w}^*)$. Under the assumption of a well-trained (i.e., learnable) graph ML problem, these residuals become negligibly small, making the error weight influence term nearly zero.

Meanwhile, the first edge influence term activates only when v_k coincides with one of the endpoints v_i or v_j , producing a direct and large perturbation. The second magnitude weight influence term, which depends on the inner-product similarity $\langle (\mathbf{x}_j^\top \mathbf{w}^*)\mathbf{z}_i + (\mathbf{x}_i^\top \mathbf{w}^*)\mathbf{z}_j, \mathbf{z}_k \rangle_{\mathbf{H}^{-1}}$, is also relatively larger when $v_k \in \{v_i, v_j\}$ and smaller otherwise. As a result, for $v_k = v_i$ or v_j , the sum of a substantial positive edge influence and a significant negative magnitude weight influence yields a dramatic net effect on the model output. Therefore, the influence on the endpoints of unlearned edges may be more significant than on other nodes. Based on this, we reasonably assume their confidence will drop and make the following claim:

Claim C.3 (Confidence Pitfall, Restatement of Claim 5.2). *The average model confidence of nodes appearing in unlearned edges is lower than that of other nodes, i.e.,*

$$\text{AvgConf}(\mathcal{V}_{\mathcal{Q}_{\text{un}}^+}) < \text{AvgConf}(\mathcal{V} \setminus \mathcal{V}_{\mathcal{Q}_{\text{un}}^+}),$$

where $\mathcal{V}_{\mathcal{Q}_0} := \{v_i : (v_i, v_j) \in \mathcal{Q}_0 \text{ or } (v_j, v_i) \in \mathcal{Q}_0\}$ is the set of all nodes involved in query subset \mathcal{Q}_0 , and $\text{AvgConf}(\mathcal{V}_0) := \frac{1}{|\mathcal{V}_0|} \sum_{v_i \in \mathcal{V}_0} \max_{\ell} \mathbf{p}_{i,\ell}$ is the average model confidence within a node set \mathcal{V}_0 .

Further empirical evidence supporting this claim can be found in Appendix E.

D Missing Proofs in Section C

In this section, we provide formal definitions for all concepts introduced in Section C and supplement the missing technical proofs. We begin by describing the architecture and training process of linear GCN, and then compute the edge influence on model weights. Next, we analyze how edge influence affects the model output, considering both the direct impact from the edge itself and the indirect influence mediated through model weights.

D.1 Linear GCN

In this analysis, we consider a simple but effective variant of GNNs, linear GCN, which is adapted from the classical SGC model [79].

Definition D.1 (Linear GCN). *Let $\mathbf{C} \in [0, 1]^{n \times n}$ be the propagation matrix, $\mathbf{X}\mathbf{b} \in \mathbb{R}^{n \times d}$ be the input feature matrix, and $\mathbf{w} \in \mathbb{R}^d$ denote the learnable weight vector. Linear GCN computes the model output as follows:*

$$f_{\text{LGN}}(\mathbf{C}, \mathbf{X}; \mathbf{w}) := \mathbf{C}\mathbf{X}\mathbf{w}.$$

Remark D.2 (Universality of Linear GCN). *The propagation matrix \mathbf{C} is adaptable to any type of convolution matrices, recovering multiple different types of GNNs, including but not limited to:*

- *One-layer GCN:* $\mathbf{C}_{I-GCN} := (\mathbf{D} + \mathbf{I})^{-0.5}(\mathbf{A} + \mathbf{I})(\mathbf{D} + \mathbf{I})^{-0.5}$;
- *k -layer SGC:* $\mathbf{C}_{k-SGC} := \mathbf{C}_{I-GCN}^k$;
- *Infinite layer PPNP:* $\mathbf{C}_{PPNP} := (\mathbf{I} - (1 - \alpha)\mathbf{C}_{I-GCN})^{-1}$;
- *k -layer APPNP:* $\mathbf{C}_{k-APPNP} := (1 - \alpha)^k \mathbf{C}_{I-GCN}^k + \alpha \sum_{l=0}^{k-1} (1 - \alpha)^l \mathbf{C}_{I-GCN}^l$;
- *One-layer GIN:* $\mathbf{C}_{GIN} := \mathbf{A} + \mathbf{I}$.

To analyze the output of a single node, we state the following basic result without proof.

Fact D.3 (Single Node Output). *The output of Linear GCN for a single node $v_i \in \mathcal{V}$ is*

$$f_{\text{LGN}}(\mathbf{C}, \mathbf{X}; \mathbf{w})_i = (\mathbf{C}\mathbf{X}\mathbf{w})_i = \sum_{j=1}^n \mathbf{C}_{j,i}(\mathbf{x}_j^\top \mathbf{w}).$$

In the training process of Linear GCNs, we consider a general least squares problem [64], which applies to both regression and binary classification tasks.

Definition D.4 (Training of Linear GCNs). *Let $\mathbf{y} \in \mathbb{R}^d$ denote the label vector; training the linear GCN in Definition D.1 is equivalent to minimize the following loss function:*

$$\mathcal{L}(\mathbf{C}, \mathbf{X}, \mathbf{w}, \mathbf{y}) := \frac{1}{2} \|\mathbf{y} - \mathbf{C}\mathbf{X}\mathbf{w}\|_2^2.$$

Proposition D.5 (Closed-form Solution for Linear GCN Training). *The following optimization problem:*

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{C}, \mathbf{X}, \mathbf{w}, \mathbf{y})$$

has a closed-form solution

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{C}^\top \mathbf{C} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{C}^\top \mathbf{y}.$$

Proof. The loss function \mathcal{L} is convex. Following the first-order optimality condition, we can set the gradient to zero and solve for \mathbf{w} , which trivially yields the desired result. \square

D.2 Edge Influence on Model Weight

In this section, we examine how a specific set of edges $\Delta\mathcal{E}$ influences the retrained model weights \mathbf{w}^* . We first define the perturbation matrix, which is the adjacency matrix for an edge subset $\Delta\mathcal{E}$, and can be used to perturb the original propagation matrix \mathbf{C} for influence evaluation.

Definition D.6 (Perturbation Matrix). *For an arbitrary edge subset $\Delta\mathcal{E}$, the corresponding perturbation matrix $\Xi^{\Delta\mathcal{E}}$ is defined as:*

$$\Xi^{\Delta\mathcal{E}} := \sum_{(v_i, v_j) \in \Delta\mathcal{E}} \mathbf{e}_i \mathbf{e}_j^\top.$$

For notation simplicity, we sometimes ignore $\Delta\mathcal{E}$ and use Ξ as a shorthand notation for Ξ in this paper.

Next, we begin with a general case of edge influence that does not address the linear GCN architecture.

Lemma D.7 (Edge Influence on Model Weight, General Case). *Let $\Xi^{\Delta\mathcal{E}} \in \mathbb{R}^{n \times n}$ be a perturbation matrix as defined in Definition D.6. Let $\epsilon \in \mathbb{R}$ be a perturbation magnitude and $\mathbf{w}^*(\epsilon)$ be the optimal model weight after perturbation. Considering an edge perturbation $\mathbf{C}(\epsilon) := \mathbf{C} + \epsilon \Xi^{\Delta\mathcal{E}}$ on the original propagation matrix \mathbf{C} , the sensitivity of the model weight \mathbf{w}^* can be characterized by:*

$$\mathcal{I}_\epsilon(\mathbf{w}^*) := \left. \frac{d\mathbf{w}^*(\epsilon)}{d\epsilon} \right|_{\epsilon=0} = -\mathbf{H}^{-1} \left. \frac{\partial}{\partial \epsilon} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{C} + \epsilon \Xi^{\Delta\mathcal{E}}, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) \right|_{\epsilon=0},$$

where $\mathbf{H} := \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{C}, \mathbf{X}, \mathbf{w}^*, \mathbf{y})$.

Proof. For notation simplicity, we use Ξ as a shorthand notation for $\Xi^{\Delta\epsilon}$ in this proof. Since $\mathbf{w}^*(\epsilon)$ is the minimizer of the perturbed loss function $\mathcal{L}(\mathbf{C}(\epsilon), \mathbf{X}, \mathbf{w}, \mathbf{y})$, we have:

$$\mathbf{w}^*(\epsilon) = \arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{C} + \epsilon\Xi, \mathbf{X}, \mathbf{w}, \mathbf{y}).$$

Examining the first-order optimality condition of the minimization problem, we have:

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{C} + \epsilon\Xi, \mathbf{X}, \mathbf{w}^*(\epsilon), \mathbf{y}) = 0. \quad (7)$$

Let us define the change in parameters as:

$$\Delta \mathbf{w} := \mathbf{w}^*(\epsilon) - \mathbf{w}^*.$$

Therefore, since $\mathbf{w}^*(\epsilon) \rightarrow \mathbf{w}^*$ as $\epsilon \rightarrow 0$, we expand Eq. (7) with multi-variate Taylor Series at the local neighborhood of $(\mathbf{w}^*, 0)$ to approximate the value of $\nabla_{\mathbf{w}} \mathcal{L}$ at $(\mathbf{w}^* + \Delta \mathbf{w}, \epsilon)$:

$$\begin{aligned} 0 &= \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{C} + \epsilon\Xi, \mathbf{X}, \mathbf{w}^*(\epsilon), \mathbf{y}) \\ &= \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{C} + \epsilon\Xi, \mathbf{X}, \mathbf{w}^* + \Delta \mathbf{w}, \mathbf{y}) \\ &= \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{C} + \epsilon\Xi, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) + \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{C} + \epsilon\Xi, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) \Delta \mathbf{w} \\ &\quad + \epsilon \frac{\partial}{\partial \epsilon} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{C} + \epsilon\Xi, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) \Big|_{\epsilon=0} + o(\|\Delta \mathbf{w}\| + |\epsilon|) \\ &= \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{C}, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) + \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{C}, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) \Delta \mathbf{w} \\ &\quad + \epsilon \frac{\partial}{\partial \epsilon} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{C} + \epsilon\Xi, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) \Big|_{\epsilon=0} + o(\|\Delta \mathbf{w}\| + |\epsilon|) \\ &= \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{C}, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) \Delta \mathbf{w} + \epsilon \frac{\partial}{\partial \epsilon} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{C} + \epsilon\Xi, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) \Big|_{\epsilon=0} \\ &\quad + o(\|\Delta \mathbf{w}\| + |\epsilon|) \end{aligned} \quad (8)$$

where the first equality follows from Eq. (7), the second equality follows from the definition of $\Delta \mathbf{w}$, the third equality follows from Taylor Series, the fourth equality follows from $\epsilon \rightarrow 0$, and the last equality follows from the fact that \mathbf{w}^* is the minimizer of loss function L .

Let the Hessian matrix at \mathbf{w}^* be $\mathbf{H} := \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{C}, \mathbf{X}, \mathbf{w}^*, \mathbf{y})$. Ignoring the remainder term and rearrange Eq. (8), we can conclude that:

$$\Delta \mathbf{w} = -\epsilon \mathbf{H}^{-1} \frac{\partial}{\partial \epsilon} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{C} + \epsilon\Xi, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) \Big|_{\epsilon=0}.$$

Dividing both sides by ϵ and taking the limit $\epsilon \rightarrow 0$, we have:

$$\frac{d\mathbf{w}^*(\epsilon)}{d\epsilon} \Big|_{\epsilon=0} = -\mathbf{H}^{-1} \frac{\partial}{\partial \epsilon} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{C} + \epsilon\Xi, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) \Big|_{\epsilon=0}.$$

This completes the proof. \square

Afterward, we extend the general case of edge influence on model weights to the Linear GCN framework, providing a closed-form solution for the edge-to-weight influence in this model.

Theorem D.8 (Edge Influence on Model Weight, Linear GCN Case). *Given an arbitrary edge perturbation matrix Ξ as defined in Definition D.6, the influence function for the optimal model weight \mathbf{w}^* is:*

$$\mathcal{I}_{\epsilon}(\mathbf{w}^*) := \frac{d\mathbf{w}^*(\epsilon)}{d\epsilon} \Big|_{\epsilon=0} = \mathbf{H}^{-1} (\mathbf{X}^{\top} \Xi^{\top} \mathbf{y} - \mathbf{X}^{\top} \Xi^{\top} \mathbf{C} \mathbf{X} \mathbf{w}^* - \mathbf{X}^{\top} \mathbf{C}^{\top} \Xi \mathbf{X} \mathbf{w}^*).$$

Proof. We start from the general result in Lemma D.7 which states that

$$\mathcal{I}_{\epsilon}(\mathbf{w}^*) := \frac{d\mathbf{w}^*(\epsilon)}{d\epsilon} \Big|_{\epsilon=0} = -\mathbf{H}^{-1} \frac{\partial}{\partial \epsilon} \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{C} + \epsilon\Xi, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) \Big|_{\epsilon=0},$$

where $\mathbf{H} = \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{C}, \mathbf{X}, \mathbf{w}^*, \mathbf{y})$ is the Hessian of the loss with respect to \mathbf{w} evaluated at \mathbf{w}^* .

Recalling the loss function in Definition D.4, when we introduce the perturbation $\mathbf{C}(\epsilon) = \mathbf{C} + \epsilon\mathbf{\Xi}$, the loss becomes

$$\mathcal{L}(\mathbf{C} + \epsilon\mathbf{\Xi}, \mathbf{X}, \mathbf{w}, \mathbf{y}) = \frac{1}{2} \|\mathbf{y} - (\mathbf{C} + \epsilon\mathbf{\Xi})\mathbf{X}\mathbf{w}\|_2^2.$$

Thus, the gradient with respect to \mathbf{w} is

$$\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{C} + \epsilon\mathbf{\Xi}, \mathbf{X}, \mathbf{w}, \mathbf{y}) = -\mathbf{X}^\top (\mathbf{C} + \epsilon\mathbf{\Xi})^{-1} \mathbf{P} [\mathbf{y} - (\mathbf{C} + \epsilon\mathbf{\Xi})\mathbf{X}\mathbf{w}].$$

Evaluating at $\mathbf{w} = \mathbf{w}^*$ yields

$$\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{C} + \epsilon\mathbf{\Xi}, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) = -\mathbf{X}^\top (\mathbf{C} + \epsilon\mathbf{\Xi})^{-1} \mathbf{P} [\mathbf{y} - (\mathbf{C} + \epsilon\mathbf{\Xi})\mathbf{X}\mathbf{w}^*].$$

We now differentiate this expression with respect to ϵ and evaluate at $\epsilon = 0$. Writing the gradient as the sum of two terms, we have:

$$\begin{aligned} \nabla_{\mathbf{w}}\mathcal{L}(\mathbf{C} + \epsilon\mathbf{\Xi}, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) &= -\underbrace{\mathbf{X}^\top (\mathbf{C} + \epsilon\mathbf{\Xi})^{-1} \mathbf{P} \mathbf{y}}_{:=\mathbf{T}_1} + \underbrace{\mathbf{X}^\top (\mathbf{C} + \epsilon\mathbf{\Xi})^{-1} \mathbf{P} (\mathbf{C} + \epsilon\mathbf{\Xi})\mathbf{X}\mathbf{w}^*}_{:=\mathbf{T}_2} \\ &= -\mathbf{T}_1 + \mathbf{T}_2. \end{aligned}$$

Now we differentiate each term with respect to ϵ . Specifically, for the first term \mathbf{T}_1 , we have the following result:

$$\begin{aligned} \frac{d\mathbf{T}_1}{d\epsilon} &= \frac{d}{d\epsilon} \left[-\mathbf{X}^\top (\mathbf{C} + \epsilon\mathbf{\Xi})^{-1} \mathbf{P} \mathbf{y} \right] \\ &= -\mathbf{X}^\top \mathbf{\Xi}^\top \mathbf{y}. \end{aligned}$$

For the second term, we can conclude by basic matrix algebra that:

$$\begin{aligned} \frac{d\mathbf{T}_2}{d\epsilon} &= \frac{d}{d\epsilon} \left[\mathbf{X}^\top (\mathbf{C} + \epsilon\mathbf{\Xi})^{-1} \mathbf{P} (\mathbf{C} + \epsilon\mathbf{\Xi})\mathbf{X}\mathbf{w}^* \right] \\ &= \mathbf{X}^\top (\mathbf{C}^\top \mathbf{\Xi} + \mathbf{\Xi}^\top \mathbf{C} + 2\epsilon\mathbf{\Xi}^\top \mathbf{\Xi})\mathbf{X}\mathbf{w}^* \end{aligned}$$

Combining both terms, we obtain

$$\left. \frac{\partial}{\partial \epsilon} \nabla_{\mathbf{w}}\mathcal{L}(\mathbf{C} + \epsilon\mathbf{\Xi}, \mathbf{X}, \mathbf{w}^*, \mathbf{y}) \right|_{\epsilon=0} = -\mathbf{X}^\top \mathbf{\Xi}^\top \mathbf{y} + \mathbf{X}^\top (\mathbf{C}^\top \mathbf{\Xi} + \mathbf{\Xi}^\top \mathbf{C})\mathbf{X}\mathbf{w}^*.$$

Substituting back into the expression for $\mathcal{I}_\epsilon(\mathbf{w}^*)$ gives:

$$\begin{aligned} \mathcal{I}_\epsilon(\mathbf{w}^*) &= -\mathbf{H}^{-1} \left(-\mathbf{X}^\top \mathbf{\Xi}^\top \mathbf{y} + \mathbf{X}^\top (\mathbf{C}^\top \mathbf{\Xi} + \mathbf{\Xi}^\top \mathbf{C})\mathbf{X}\mathbf{w}^* \right) \\ &= \mathbf{H}^{-1} \left(\mathbf{X}^\top \mathbf{\Xi}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{C}^\top \mathbf{\Xi} \mathbf{X} \mathbf{w}^* - \mathbf{X}^\top \mathbf{\Xi}^\top \mathbf{C} \mathbf{X} \mathbf{w}^* \right). \end{aligned}$$

This finishes the proof. □

D.3 Edge Influence on Model Output

In this section, we compute the influence function for edges on the final model output, considering both direct edge influence and the effect of model parameters on the output. First, we calculate the influence on all model outputs.

Theorem D.9 (Edge Influence on Model Output, Linear GCN Case). *Let $\mathbf{\Xi}^{\Delta\mathcal{E}} \in \mathbb{R}^{n \times n}$ be a perturbation matrix as defined in Definition D.6. Let $\epsilon \in \mathbb{R}$ be a perturbation magnitude and $\mathbf{w}^*(\epsilon)$ be the optimal model weight after perturbation. Considering an edge perturbation $\mathbf{C}(\epsilon) := \mathbf{C} + \epsilon\mathbf{\Xi}^{u,v}$ on the original propagation matrix \mathbf{C} , the sensitivity of the Linear GCN model output $f_{\text{LGN}}(\mathbf{C} + \epsilon\mathbf{\Xi}, \mathbf{X}; \mathbf{w}^*(\epsilon))$ as defined in Definition D.1 can be characterized by:*

$$\begin{aligned} \mathcal{I}_\epsilon(f_{\text{LGN}}) &:= \left. \frac{df_{\text{LGN}}(\mathbf{C} + \epsilon\mathbf{\Xi}, \mathbf{X}; \mathbf{w}^*(\epsilon))}{d\epsilon} \right|_{\epsilon=0} \\ &= \mathbf{\Xi} \mathbf{X} \mathbf{w}^* + \mathbf{C} \mathbf{X} \mathbf{H}^{-1} \left(\mathbf{X}^\top \mathbf{\Xi}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{C}^\top \mathbf{\Xi} \mathbf{X} \mathbf{w}^* - \mathbf{X}^\top \mathbf{\Xi}^\top \mathbf{C} \mathbf{X} \mathbf{w}^* \right), \end{aligned}$$

where $\mathbf{H} = \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{C}, \mathbf{X}, \mathbf{w}^*, \mathbf{y})$ is the Hessian of the loss with respect to \mathbf{w} evaluated at \mathbf{w}^* .

Proof. The edge perturbation $\epsilon \Xi$ influences f_{LGN} from two different aspects: propagation matrix and optimal model weights. Thus, the influence function can be derived by chain rule as follows:

$$\begin{aligned}
\mathcal{I}_\epsilon(f_{\text{LGN}}) &= \left. \frac{df_{\text{LGN}}(\mathbf{C} + \epsilon \Xi, \mathbf{X}; \mathbf{w}^*(\epsilon))}{d\epsilon} \right|_{\epsilon=0} \\
&= \left(\frac{\partial f_{\text{LGN}}(\mathbf{C} + \epsilon \Xi, \mathbf{X}; \mathbf{w}^*(\epsilon))}{\partial(\mathbf{C} + \epsilon \Xi)} \cdot \frac{\partial(\mathbf{C} + \epsilon \Xi)}{\partial \epsilon} \right) \Big|_{\epsilon=0} \\
&\quad + \left(\frac{\partial f_{\text{LGN}}(\mathbf{C} + \epsilon \Xi, \mathbf{X}; \mathbf{w}^*(\epsilon))}{\partial \mathbf{w}^*(\epsilon)} \cdot \frac{\partial \mathbf{w}^*(\epsilon)}{\partial \epsilon} \right) \Big|_{\epsilon=0} \\
&= \Xi \mathbf{X} \mathbf{w}^* + \left(\frac{\partial f_{\text{LGN}}(\mathbf{C} + \epsilon \Xi, \mathbf{X}; \mathbf{w}^*(\epsilon))}{\partial \mathbf{w}^*(\epsilon)} \cdot \frac{\partial \mathbf{w}^*(\epsilon)}{\partial \epsilon} \right) \Big|_{\epsilon=0} \\
&= \Xi \mathbf{X} \mathbf{w}^* + \mathbf{C} \mathbf{X} \left(\frac{\partial \mathbf{w}^*(\epsilon)}{\partial \epsilon} \right) \Big|_{\epsilon=0} \\
&= \Xi \mathbf{X} \mathbf{w}^* + \mathbf{C} \mathbf{X} \mathbf{H}^{-1} \left(\mathbf{X}^\top \Xi^\top \mathbf{y} - \mathbf{X}^\top \mathbf{C}^\top \Xi \mathbf{X} \mathbf{w}^* - \mathbf{X}^\top \Xi^\top \mathbf{C} \mathbf{X} \mathbf{w}^* \right).
\end{aligned}$$

where the first equality follows from the definition of the influence function, the second equality follows from the chain rule, the third and fourth equality follow from basic matrix calculus, and the last equality follows from Theorem D.8. \square

Next, we present two immediate corollaries of this theorem and determine the influence of a single edge on a single node's output.

Corollary D.10 (Single Edge Influence on Model Output, Linear GCN Case). *Let $\mathbf{Z} := \mathbf{C} \mathbf{X}$ and \mathbf{z}_l is the l -th row of \mathbf{Z} . Considering the influence of one specific edge (v_i, v_j) on undirected graphs (i.e., $\Xi = \mathbf{e}_i \mathbf{e}_j^\top + \mathbf{e}_j \mathbf{e}_i^\top$), the influence function for the model output is:*

$$\mathcal{I}_\epsilon(f_{\text{LGN}}) = (\mathbf{x}_j^\top \mathbf{w}^*) \mathbf{e}_i + (\mathbf{x}_i^\top \mathbf{w}^*) \mathbf{e}_j + (\mathbf{q} \otimes \mathbf{C} \mathbf{X}) \text{vec}(\mathbf{H}^{-1}),$$

where $\mathbf{H} = \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{C}, \mathbf{X}, \mathbf{w}^*, \mathbf{y})$ is the Hessian of the loss with respect to \mathbf{w} evaluated at \mathbf{w}^* and

$$\mathbf{q} := (y_j - \mathbf{z}_j^\top \mathbf{w}^*) \mathbf{x}_i + (y_i - \mathbf{z}_i^\top \mathbf{w}^*) \mathbf{x}_j - ((\mathbf{x}_j^\top \mathbf{w}^*) \mathbf{z}_i + (\mathbf{x}_i^\top \mathbf{w}^*) \mathbf{z}_j).$$

Proof. This follows from basic algebra and the fact that $\text{vec}(\mathbf{A} \mathbf{X} \mathbf{B}) = (\mathbf{B} \otimes \mathbf{A}) \text{vec}(\mathbf{X})$. \square

Corollary D.11 (Single Edge Influence on Single Node's Model Output, Linear GCN Case). *Let $\mathbf{Z} := \mathbf{C} \mathbf{X}$ and \mathbf{z}_l is the l -th row of \mathbf{Z} . Considering the influence of one specific edge (v_i, v_j) on undirected graphs (i.e., $\Xi = \mathbf{e}_i \mathbf{e}_j^\top + \mathbf{e}_j \mathbf{e}_i^\top$), the influence function for the model output for specific node v_k is:*

$$\begin{aligned}
\mathcal{I}_\epsilon(f_{\text{LGN}})_k &= \mathbf{1}\{v_k = v_i\} \cdot (\mathbf{x}_v^\top \mathbf{w}^*) + \mathbf{1}\{v_k = v_j\} \cdot (\mathbf{x}_u^\top \mathbf{w}^*) + \mathbf{q}^\top \mathbf{H}^{-1} \mathbf{z}_k \\
&= \underbrace{\mathbf{1}\{v_k = v_i\} \cdot (\mathbf{x}_i^\top \mathbf{w}^*) + \mathbf{1}\{v_k = v_j\} \cdot (\mathbf{x}_j^\top \mathbf{w}^*)}_{\text{edge influence}} \\
&\quad - \underbrace{\langle (\mathbf{x}_j^\top \mathbf{w}^*) \mathbf{z}_i + (\mathbf{x}_i^\top \mathbf{w}^*) \mathbf{z}_j, \mathbf{z}_k \rangle_{\mathbf{H}^{-1}}}_{\text{magnitude weight influence}} + \underbrace{\langle (y_j - \mathbf{z}_j^\top \mathbf{w}^*) \mathbf{x}_i + (y_i - \mathbf{z}_i^\top \mathbf{w}^*) \mathbf{x}_j, \mathbf{z}_k \rangle_{\mathbf{H}^{-1}}}_{\text{error weight influence}},
\end{aligned}$$

where $\mathbf{H} = \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{C}, \mathbf{X}, \mathbf{w}^*, \mathbf{y})$ is the Hessian of the loss with respect to \mathbf{w} evaluated at \mathbf{w}^* .

Proof. This follows from basic algebra, the definition of z , and the definition of inner product w.r.t. a PSD matrix $\langle \cdot, \cdot \rangle_{\mathbf{A}}$. \square

E Preliminary Experiments

In this section, we present empirical evidence from preliminary experiments to support our main claims, probability similarity gap (Claim 5.1) and confidence pitfall (Claim 5.2), in this paper.

Dataset	Unlearn Method	ProbSim(\mathcal{Q}^-)	ProbSim($\mathcal{Q}_{\text{un}}^+$)	ProbSim($\mathcal{Q}_{\text{mem}}^+$)
Cora	GIF	0.1979 ± 0.3147	0.6552 ± 0.3457	0.8001 ± 0.2689
	CEU	0.1997 ± 0.3172	0.6553 ± 0.3502	0.8122 ± 0.2528
	GA	0.1955 ± 0.3119	0.6511 ± 0.3486	0.8101 ± 0.2625
Citeseer	GIF	0.2689 ± 0.2997	0.6446 ± 0.3046	0.8250 ± 0.2116
	CEU	0.2492 ± 0.3145	0.6248 ± 0.3317	0.8251 ± 0.2311
	GA	0.2695 ± 0.3071	0.6397 ± 0.3200	0.8340 ± 0.2051
Pubmed	GIF	0.5833 ± 0.3507	0.7450 ± 0.2752	0.8843 ± 0.1682
	CEU	0.5799 ± 0.3677	0.7494 ± 0.2856	0.8954 ± 0.1741
	GA	0.5997 ± 0.3687	0.7548 ± 0.2865	0.9031 ± 0.1677
LastFM-Asia	GIF	0.1608 ± 0.3413	0.8529 ± 0.3278	0.9308 ± 0.2234
	CEU	0.1631 ± 0.3433	0.8564 ± 0.3224	0.9330 ± 0.2129
	GA	0.1638 ± 0.3451	0.8439 ± 0.3411	0.9238 ± 0.2373

Table 2: **Average black-box probability similarity between different types of edges.** We consider three types of edges: negative edges (\mathcal{Q}^-), unlearned edges ($\mathcal{Q}_{\text{un}}^+$), and other membership edges ($\mathcal{Q}_{\text{mem}}^+$). The similarity measure is based on JS Divergence.

E.1 Probability Similarity Gap

To validate that different types of edges require different levels of similarity thresholds (Claim 5.1), we present an important preliminary study in this section. Specifically, we adopt a GCN backbone and define a specific query set $\mathcal{Q} = \mathcal{Q}_{\text{un}}^+ \cup \mathcal{Q}_{\text{mem}}^+ \cup \mathcal{Q}^-$, which consists of 5% of edges marked as unlearned ($\mathcal{Q}_{\text{un}}^+$), 5% of other membership edges ($\mathcal{Q}_{\text{mem}}^+$), and 10% of negative edges (\mathcal{Q}^-). For any two nodes $(v_i, v_j) \in \mathcal{Q}$, we compute the following similarity metric based on Jensen-Shannon divergence:

$$\phi(\mathbf{p}_i, \mathbf{p}_j) = 1 - \frac{1}{2} [\text{KL}(\mathbf{p}_i \parallel \mathbf{m}) + \text{KL}(\mathbf{p}_j \parallel \mathbf{m})],$$

where \mathbf{p}_i and \mathbf{p}_j are the predictive probability distributions of v_i and v_j from the unlearned victim model, and $\mathbf{m} = \frac{1}{2}(\mathbf{p}_i + \mathbf{p}_j)$. The higher $\phi(\mathbf{p}_i, \mathbf{p}_j)$ is, the more similar the model’s black-box predictions are on the two nodes.

We report the mean and standard deviation of the probability similarity for each group in Table 2. The key observations are:

(i) There exists a clear and consistent gap in the probability similarity among the three edge types, aligning with Claim 5.1:

$$\text{ProbSim}(\mathcal{Q}^-) < \text{ProbSim}(\mathcal{Q}_{\text{un}}^+) < \text{ProbSim}(\mathcal{Q}_{\text{mem}}^+).$$

This supports the need for an adaptive prediction mechanism that distinguishes between unlearned and other membership edges, potentially improving prediction accuracy. Accordingly, our model in Eq. (5) incorporates a learnable transformation on trend features to adjust predicted similarities obtained from MIA methods.

(ii) The variance within each group is relatively large, indicating that although a similarity gap exists, simple probability similarity computation may not be sufficient for membership inference. For instance, on the Cora dataset, the average probability similarity in \mathcal{Q}^- is around 0.195, while its standard variance is nearly 0.31. This effect arises because, although most edges in \mathcal{Q}^- have similarity near 0, there are outliers with high similarity levels. To address this large variance, we incorporate a broad range of probability-based similarity features, alongside node feature similarities used in prior MIA methods. These additional features enhance similarity representation and help stabilize membership inference.

E.2 Confidence Pitfall

To demonstrate that the confidence of nodes connected to unlearned edges drops significantly compared to other nodes, we conduct a preliminary experiment. Specifically, we use GCN as the

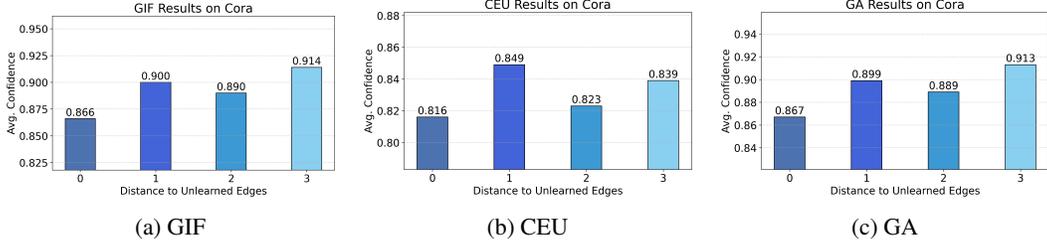


Figure 5: The relation between average model confidence and distance to unlearned edges on Cora.

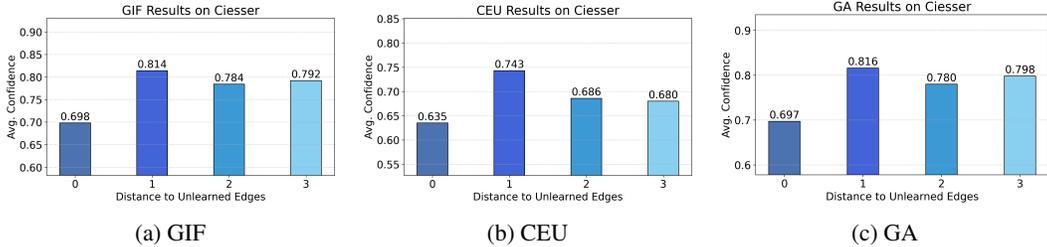


Figure 6: The relation between average model confidence and distance to unlearned edges on Citeseer.

victim model. We randomly select 5% of the edges as unlearned edges $\Delta\mathcal{E}$ and apply three unlearning methods, GIF, CEU, and GA, to unlearn the trained GCN. The selected edges are also removed from the graph.

We then perform inference on all nodes using the unlearned GNN and compute the average confidence for nodes grouped by their shortest path distance to the endpoints of unlearned edges.

The results on two datasets are shown in Figure 5 and Figure 6. We summarize the following observation:

For all datasets, nodes directly connected to unlearned edges (distance = 0) exhibit substantially lower confidence than their immediate neighbors (distance ≥ 1). This aligns with our theoretical analysis in Section C, which shows that the influence of edge unlearning on the endpoint is significantly greater than on other nodes. This empirical finding suggests that the confidence gap between a node and its neighborhood is a strong indicator of its connection to an unlearned edge. Consequently, an attacker can adjust the similarity threshold accordingly to enhance membership inference accuracy.

F Experimental Settings

Model Parameters. All victim models use an embedding size of 16 and are trained for 100 epochs. For learning rates and weight decays, we follow the settings of GIF [80]. All models are optimized using the Adam optimizer [35]. For all neural attack models, including StealLink [29], MIA-GNN [58], and TrendAttack, we follow the settings of StealLink [29], using a hidden dimension of 64 and a 2-layer MLP to encode features. All membership inference models are trained using learning rate = 0.01, weight decay = 0.0001, and the Adam optimizer [35]. They are trained with the binary cross-entropy loss.

Unlearning Method Settings. We now supplement the missing details of our unlearning methods.

- **GIF** [80]: The number of estimation iterations is 100 and the damping factor is 0. On the Cora and Citeseer datasets, the scale factor λ is set to 500, following the official implementation. On the Pubmed and LastFM-Asia datasets, we search λ in $\{10^1, 10^2, 10^3, 10^4, \dots\}$, following the original search space. This ensures that the unlearned models maintain utility and produce meaningful outputs. We select the smallest λ that results in a meaningful model and find that $\lambda = 10^3$ works well for Pubmed and $\lambda = 10^4$ for LastFM-Asia.

- **CEU** [81]: For a fair comparison, we use the same parameter search space as GIF for all influence-function-related parameters, including iterations, damping factor, and λ . For the noise variance, we search in $\{0.1, 0.5, 0.01, 0.005, 0.001, 0.0005, 0.0001\}$ to ensure meaningful unlearning performance.
- **GA** [78]: We follow the official GraphGuard settings and use 1 gradient ascent epoch. We tune the unlearning magnitude α in $\{1, 0.5, 0.1, 0.05, 0.01, 0.005, 0.001\}$ and find that $\alpha = 0.5$ gives the best results.

Baselines. Details for our baselines are as follows:

- **GraphSAGE** [27]: We use a simple link prediction model trained on the shadow graph to perform link prediction on the unlearned graph. The embedding size is 64, with 2 layers. The negative-to-positive link ratio is 1:1.
- **NCN** [74]: This is a state-of-the-art link prediction method. We follow the official settings from their code repository. We search across the provided backbones: GCN [36], GIN [84], and SAGE [70]. The embedding dimension is set to the default value of 256. The batch size is 16, the learning rate is 0.01, and the number of training epochs is 100.
- **MIA-GNN** [58]: Since no official implementation is available, we follow the settings described in the paper. We re-implement the model using a 2-layer MLP feature extractor with a hidden dimension of 64. We use another 2-layer MLP with a dimension of 16 as the link predictor. Training uses binary cross-entropy loss and the same setup as StealLink, ensuring a fair comparison.
- **StealLink** [29]: This is a representative membership inference attack. We use its strongest variant, Attack-7, which has access to the shadow dataset and partial features and connectivity from the target dataset. We follow the official codebase and use all 8 distance metrics (e.g., cosine, Euclidean). For the probability metric, we search among entropy, KL divergence, and JS divergence, and report the best result. The reference model is a 2-layer MLP with an embedding size of 16.
- **GroupAttack** [91]: This is a recent membership inference attack that relies solely on hard labels and thresholding. We search the threshold hyperparameter α on the shadow dataset in the range $[0, 1]$ with a step size of 0.05.

Reproducibility. All experiments are conducted using Python 3.12.2, PyTorch 2.3.1+cu121, and PyTorch Geometric 2.5.3. The experiments run on a single NVIDIA RTX A6000 GPU. The server used has 64 CPUs, and the type is AMD EPYC 7282 16-Core CPU. All source code is provided in the supplementary materials.

G Additional Experiments

Unlearn method	Attack	Cora			Citeseer		
		Unlearned	Original	All	Unlearned	Original	All
GIF	GraphSAGE	0.5356 ± 0.0124	0.5484 ± 0.0180	0.5420 ± 0.0129	0.5275 ± 0.0452	0.5216 ± 0.0535	0.5246 ± 0.0474
	NCN	0.7403 ± 0.0309	0.7405 ± 0.0222	0.7404 ± 0.0237	0.6750 ± 0.0881	0.6872 ± 0.0799	0.6811 ± 0.0816
	MIA-GNN	0.7547 ± 0.0809	0.7916 ± 0.0939	0.7732 ± 0.0865	0.7802 ± 0.0251	0.8245 ± 0.0230	0.8023 ± 0.0210
	StealLink	0.7841 ± 0.0357	0.8289 ± 0.0576	0.8065 ± 0.0395	0.7369 ± 0.0463	0.8404 ± 0.0360	0.7887 ± 0.0373
	GroupAttack	0.7982 ± 0.0072	0.8053 ± 0.0190	0.8018 ± 0.0125	0.7771 ± 0.0092	0.7618 ± 0.0069	0.7695 ± 0.0051
	TrendAttack-MIA	<u>0.8240 ± 0.0209</u>	<u>0.8448 ± 0.0185</u>	<u>0.8344 ± 0.0189</u>	<u>0.8069 ± 0.0185</u>	0.8078 ± 0.0284	<u>0.8073 ± 0.0228</u>
	TrendAttack-SL	0.8309 ± 0.0250	0.8527 ± 0.0140	0.8418 ± 0.0188	0.8410 ± 0.0169	0.8430 ± 0.0276	0.8420 ± 0.0207
	CEU	GraphSAGE	0.5356 ± 0.0124	0.5484 ± 0.0180	0.5420 ± 0.0129	0.5275 ± 0.0452	0.5216 ± 0.0535
	NCN	0.7403 ± 0.0309	0.7405 ± 0.0222	0.7404 ± 0.0237	0.6750 ± 0.0881	0.6872 ± 0.0799	0.6811 ± 0.0816
	MIA-GNN	0.7458 ± 0.0761	0.7810 ± 0.0928	0.7634 ± 0.0840	0.7718 ± 0.0262	0.8248 ± 0.0264	0.7983 ± 0.0219
	StealLink	0.7901 ± 0.0224	0.8486 ± 0.0099	0.8193 ± 0.0068	0.7643 ± 0.0242	0.8450 ± 0.0257	0.8046 ± 0.0217
	GroupAttack	0.7941 ± 0.0086	0.7976 ± 0.0137	0.7958 ± 0.0105	0.7557 ± 0.0092	0.7458 ± 0.0131	0.7508 ± 0.0107
	TrendAttack-MIA	<u>0.8194 ± 0.0170</u>	0.8333 ± 0.0213	<u>0.8263 ± 0.0178</u>	<u>0.7933 ± 0.0206</u>	0.8041 ± 0.0261	0.7987 ± 0.0226
	TrendAttack-SL	0.8467 ± 0.0229	0.8612 ± 0.0113	0.8539 ± 0.0149	0.8514 ± 0.0214	0.8400 ± 0.0216	0.8457 ± 0.0208
GA	GraphSAGE	0.5356 ± 0.0124	0.5484 ± 0.0180	0.5420 ± 0.0129	0.5275 ± 0.0452	0.5216 ± 0.0535	0.5246 ± 0.0474
	NCN	0.7403 ± 0.0309	0.7405 ± 0.0222	0.7404 ± 0.0237	0.6750 ± 0.0881	0.6872 ± 0.0799	0.6811 ± 0.0816
	MIA-GNN	0.7676 ± 0.0584	0.8068 ± 0.0489	0.7872 ± 0.0531	0.7798 ± 0.0146	0.8353 ± 0.0175	0.8076 ± 0.0117
	StealLink	0.7862 ± 0.0402	0.8301 ± 0.0667	0.8082 ± 0.0475	0.7479 ± 0.0512	0.8431 ± 0.0305	0.7955 ± 0.0332
	GroupAttack	0.7945 ± 0.0133	0.8042 ± 0.0123	0.7993 ± 0.0122	0.7662 ± 0.0112	0.7563 ± 0.0080	0.7613 ± 0.0086
	TrendAttack-MIA	<u>0.8193 ± 0.0276</u>	0.8397 ± 0.0219	<u>0.8295 ± 0.0244</u>	0.8080 ± 0.0135	0.8249 ± 0.0331	0.8165 ± 0.0227
	TrendAttack-SL	0.8270 ± 0.0307	<u>0.8382 ± 0.0308</u>	0.8326 ± 0.0287	0.8628 ± 0.0131	0.8614 ± 0.0298	0.8621 ± 0.0209

Table 3: **Main Comparison Results on Cora and Citeseer (with variance).** We present the AUC scores for attack methods across different edge groups, now including mean \pm variance. The best results are highlighted in **bold**, while the second-best results are underlined.

Unlearn method	Attack	Pubmed			LastFM-Asia		
		Unlearned	Original	All	Unlearned	Original	All
GIF	GraphSAGE	0.6503 ± 0.0114	0.6457 ± 0.0134	0.6480 ± 0.0118	0.6914 ± 0.0620	0.6853 ± 0.0613	0.6884 ± 0.0611
	NCN	0.6661 ± 0.0321	0.6718 ± 0.0314	0.6690 ± 0.0312	0.7283 ± 0.0177	0.7273 ± 0.0120	0.7278 ± 0.0141
	MIA-GNN	0.7028 ± 0.0069	0.7902 ± 0.0041	0.7465 ± 0.0044	0.5955 ± 0.0498	0.5744 ± 0.0751	0.5850 ± 0.0614
	StealLink	0.8248 ± 0.0098	0.8964 ± 0.0027	0.8606 ± 0.0052	0.8472 ± 0.0099	0.9037 ± 0.0097	0.8755 ± 0.0089
	GroupAttack	0.6497 ± 0.0021	0.6554 ± 0.0049	0.6525 ± 0.0033	0.7858 ± 0.0044	0.7850 ± 0.0033	0.7854 ± 0.0027
	TrendAttack-MIA	0.8950 ± 0.0054	0.9171 ± 0.0039	0.9060 ± 0.0032	0.7795 ± 0.0690	0.7649 ± 0.0944	0.7722 ± 0.0814
	TrendAttack-SL	0.9524 ± 0.0026	0.9535 ± 0.0025	0.9529 ± 0.0024	0.9078 ± 0.0069	0.9134 ± 0.0039	0.9106 ± 0.0050
CEU	GraphSAGE	0.6503 ± 0.0114	0.6457 ± 0.0134	0.6480 ± 0.0118	0.6914 ± 0.0620	0.6853 ± 0.0613	0.6884 ± 0.0611
	NCN	0.6661 ± 0.0321	0.6718 ± 0.0314	0.6690 ± 0.0312	0.7283 ± 0.0177	0.7273 ± 0.0120	0.7278 ± 0.0141
	MIA-GNN	0.6626 ± 0.0237	0.6561 ± 0.0256	0.6593 ± 0.0243	0.6004 ± 0.0363	0.5811 ± 0.0568	0.5908 ± 0.0459
	StealLink	0.8467 ± 0.0168	0.9088 ± 0.0102	0.8777 ± 0.0134	0.8416 ± 0.0163	0.9021 ± 0.0084	0.8719 ± 0.0114
	GroupAttack	0.6388 ± 0.0052	0.6430 ± 0.0054	0.6409 ± 0.0053	0.7845 ± 0.0038	0.7817 ± 0.0010	0.7831 ± 0.0023
	TrendAttack-MIA	0.8982 ± 0.0038	0.9184 ± 0.0025	0.9083 ± 0.0018	0.7676 ± 0.0722	0.7576 ± 0.0986	0.7626 ± 0.0850
	TrendAttack-SL	0.9550 ± 0.0032	0.9579 ± 0.0028	0.9565 ± 0.0028	0.9037 ± 0.0062	0.9088 ± 0.0020	0.9062 ± 0.0040
GA	GraphSAGE	0.6503 ± 0.0114	0.6457 ± 0.0134	0.6480 ± 0.0118	0.6914 ± 0.0620	0.6853 ± 0.0613	0.6884 ± 0.0611
	NCN	0.6661 ± 0.0321	0.6718 ± 0.0314	0.6690 ± 0.0312	0.7283 ± 0.0177	0.7273 ± 0.0120	0.7278 ± 0.0141
	MIA-GNN	0.7242 ± 0.0099	0.8039 ± 0.0128	0.7641 ± 0.0112	0.6200 ± 0.0636	0.6057 ± 0.0884	0.6129 ± 0.0756
	StealLink	0.8203 ± 0.0128	0.8898 ± 0.0056	0.8550 ± 0.0080	0.8342 ± 0.0182	0.8947 ± 0.0040	0.8644 ± 0.0106
	GroupAttack	0.6458 ± 0.0076	0.6493 ± 0.0111	0.6475 ± 0.0093	0.7746 ± 0.0069	0.7760 ± 0.0031	0.7753 ± 0.0048
	TrendAttack-MIA	0.8932 ± 0.0052	0.9158 ± 0.0048	0.9045 ± 0.0037	0.7255 ± 0.0715	0.7099 ± 0.0778	0.7177 ± 0.0742
	TrendAttack-SL	0.9531 ± 0.0027	0.9537 ± 0.0034	0.9534 ± 0.0029	0.9041 ± 0.0054	0.9119 ± 0.0034	0.9080 ± 0.0040

Table 4: **Main Comparison Results on Pubmed and LastFM-Asia (with variance)**. We present the AUC scores for attack methods across different edge groups, now including mean \pm variance. The best results are highlighted in **bold**, while the second-best results are underlined.

Variance of Comparison Results. Due to space limitations, the comparison results in Table 1 in Section 6 do not include the standard variance from five repeated experiments. We now supplement all variance results in Table 3 and Table 4. We make the following observations regarding the stability of our proposed TrendAttack: Compared with the no-trend-feature counterparts MIA-GNN and StealLink, our TrendAttack-MIA and TrendAttack-SL exhibit smaller variances, indicating better stability. This highlights the effectiveness of the trend features, which not only improve attack performance on both edge groups but also reduce variance, enhancing model stability.

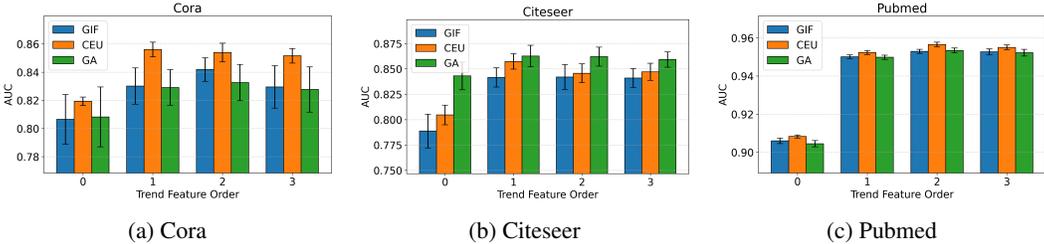


Figure 7: **Ablation study on the impact of trend feature order**. Overall attack AUC as a function of the trend feature order (0–3) for three unlearn methods across three datasets.

Impact of Trend Feature Orders. We investigate the influence of trend feature orders on the performance of the proposed TrendAttack framework. Following the experimental setup described in Section 6.2, we employ the most effective variant, TrendAttack-SL, for this ablation study. Specifically, we vary the trend feature order k (see Algorithm 2) to assess the trade-off between incorporating additional neighborhood information and achieving high attack performance. The results are presented in Figure 7, from which we make the following observations:

- (i) Incorporating trend features of orders 1, 2, or 3 significantly improves attack performance compared to the 0-th order (which corresponds to a degenerate form of TrendAttack that reduces to a simple StealLink attack). This highlights the effectiveness of our proposed trend feature design.
- (ii) The attack performance remains relatively stable across orders 1 to 3, indicating that the method is robust to the choice of trend order. Notably, lower-order features (e.g., order 1) still have strong performance while requiring less auxiliary neighborhood information, making them more practical in real-world inversion attack scenarios.

H Impact Statement

In this work, we propose a novel membership inference attack targeting unlearned GNNs, revealing critical privacy vulnerabilities in existing graph unlearning methods. Our goal is to raise awareness about privacy protection in Web services and to inspire future research on privacy-preserving graph machine learning. While our method introduces a new attack strategy, all experiments are conducted

on publicly available benchmark datasets. Given the gap between this pioneering study and real-world deployment scenarios, we do not foresee significant negative societal implications from this work.

References

- [1] Max Aliapoulios, Emmi Bevensee, Jeremy Blackburn, Barry Bradlyn, Emiliano De Cristofaro, Gianluca Stringhini, and Savvas Zannettou. A large open dataset from the parler social network. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 15, pages 943–951, 2021.
- [2] Timothy G Armstrong, Vamsi Ponnkanti, Dhruva Borthakur, and Mark Callaghan. Linkbench: a database benchmark based on the facebook social graph. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, pages 1185–1196, 2013.
- [3] Martin Bertran, Shuai Tang, Michael Kearns, Jamie H Morgenstern, Aaron Roth, and Steven Z Wu. Reconstruction attacks on machine unlearning: Simple models are vulnerable [35]. *Advances in Neural Information Processing Systems*, 37:104995–105016, 2024.
- [4] Lucas Bourtole, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *2021 IEEE symposium on security and privacy (SP)*, pages 141–159. IEEE, 2021.
- [5] Jonathan Brophy and Daniel Lowd. Machine unlearning for random forests. In *International Conference on Machine Learning*, pages 1092–1104. PMLR, 2021.
- [6] Yinzi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In *2015 IEEE symposium on security and privacy*, pages 463–480. IEEE, 2015.
- [7] Chong Chen, Fei Sun, Min Zhang, and Bolin Ding. Recommendation unlearning. In *Proceedings of the ACM web conference 2022*, pages 2768–2777, 2022.
- [8] Dingfan Chen, Ning Yu, Yang Zhang, and Mario Fritz. Gan-leaks: A taxonomy of membership inference attacks against generative models. In *Proceedings of the 2020 ACM SIGSAC conference on computer and communications security*, pages 343–362, 2020.
- [9] Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. When machine unlearning jeopardizes privacy. In *Proceedings of the 2021 ACM SIGSAC conference on computer and communications security*, pages 896–911, 2021.
- [10] Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. Graph unlearning. In *Proceedings of the 2022 ACM SIGSAC conference on computer and communications security*, pages 499–513, 2022.
- [11] Jiali Cheng, George Dasoulas, Huan He, Chirag Agarwal, and Marinka Zitnik. GNNDelete: A general strategy for unlearning in graph neural networks. In *The Eleventh International Conference on Learning Representations*, 2023.
- [12] Eli Chien, Chao Pan, and Olgica Milenkovic. Certified graph unlearning. In *NeurIPS 2022 Workshop: New Frontiers in Graph Learning*, 2022.
- [13] Avery Ching, Sergey Edunov, Maja Kabiljo, Dionysios Logothetis, and Sambavi Muthukrishnan. One trillion edges: Graph processing at facebook-scale. *Proceedings of the VLDB Endowment*, 8(12):1804–1815, 2015.
- [14] Christopher A Choquette-Choo, Florian Tramer, Nicholas Carlini, and Nicolas Papernot. Label-only membership inference attacks. In *International conference on machine learning*, pages 1964–1974. PMLR, 2021.
- [15] Mauro Conti, Jiaxin Li, Stjepan Picek, and Jing Xu. Label-only membership inference attack against node-level graph neural networks. In *Proceedings of the 15th ACM Workshop on Artificial Intelligence and Security*, pages 1–12, 2022.
- [16] Enyan Dai, Limeng Cui, Zhengyang Wang, Xianfeng Tang, Yinghan Wang, Monica Cheng, Bing Yin, and Suhang Wang. A unified framework of graph information bottleneck for robustness and membership privacy. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 368–379, 2023.

- [17] Enyan Dai, Yiwei Sun, and Suhang Wang. Ginger cannot cure cancer: Battling fake health news with a comprehensive data repository. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 14, pages 853–862, 2020.
- [18] Enyan Dai, Tianxiang Zhao, Huaisheng Zhu, Junjie Xu, Zhimeng Guo, Hui Liu, Jiliang Tang, and Suhang Wang. A comprehensive survey on trustworthy graph neural networks: Privacy, robustness, fairness, and explainability. *Machine Intelligence Research*, 21(6):1011–1061, 2024.
- [19] Yushun Dong, Binchi Zhang, Zhenyu Lei, Na Zou, and Jundong Li. Idea: A flexible framework of certified unlearning for graph neural networks. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 621–630, 2024.
- [20] Yingtong Dou, Zhiwei Liu, Li Sun, Yutong Deng, Hao Peng, and Philip S Yu. Enhancing graph neural network-based fraud detectors against camouflaged fraudsters. In *Proceedings of the 29th ACM international conference on information & knowledge management*, pages 315–324, 2020.
- [21] Vasisht Duddu, Antoine Boutet, and Virat Shejwalkar. Quantifying privacy leakage in graph embedding. In *MobiQuitous 2020-17th EAI International Conference on Mobile and Ubiquitous Systems: Computing, Networking and Services*, pages 76–85, 2020.
- [22] Bowen Fan, Yuming Ai, Xunkai Li, Zhilin Guo, Rong-Hua Li, and Guoren Wang. Opengu: A comprehensive benchmark for graph unlearning. *arXiv preprint arXiv:2501.02728*, 2025.
- [23] Wenqi Fan, Yao Ma, Qing Li, Yuan He, Eric Zhao, Jiliang Tang, and Dawei Yin. Graph neural networks for social recommendation. In *The world wide web conference*, pages 417–426, 2019.
- [24] Jiarui Feng, Yixin Chen, Fuhai Li, Anindya Sarkar, and Muhan Zhang. How powerful are k-hop message passing graph neural networks. *Advances in Neural Information Processing Systems*, 35:4776–4790, 2022.
- [25] Ziwang Fu, Feng Liu, Jiahao Zhang, Hanyang Wang, Chengyi Yang, Qing Xu, Jiayin Qi, Xiangling Fu, and Aimin Zhou. Sagn: semantic adaptive graph network for skeleton-based human action recognition. In *Proceedings of the 2021 International Conference on Multimedia Retrieval*, pages 110–117, 2021.
- [26] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International conference on machine learning*, pages 1263–1272. PMLR, 2017.
- [27] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30, 2017.
- [28] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, pages 639–648, 2020.
- [29] Xinlei He, Jinyuan Jia, Michael Backes, Neil Zhenqiang Gong, and Yang Zhang. Stealing links from graph neural networks. In *30th USENIX security symposium (USENIX security 21)*, pages 2669–2686, 2021.
- [30] Xinlei He, Rui Wen, Yixin Wu, Michael Backes, Yun Shen, and Yang Zhang. Node-level membership inference attacks against graph neural networks. *arXiv preprint arXiv:2102.05429*, 2021.
- [31] Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, Philip S Yu, and Xuyun Zhang. Membership inference attacks on machine learning: A survey. *ACM Computing Surveys (CSUR)*, 54(11s):1–37, 2022.

- [32] Hongsheng Hu, Shuo Wang, Tian Dong, and Minhui Xue. Learn what you want to unlearn: Unlearning inversion attacks against machine unlearning. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 3257–3275. IEEE, 2024.
- [33] Jinghan Jia, Jiancheng Liu, Parikshit Ram, Yuguang Yao, Gaowen Liu, Yang Liu, Pranay Sharma, and Sijia Liu. Model sparsity can simplify machine unlearning. *Advances in Neural Information Processing Systems*, 36:51584–51605, 2023.
- [34] Wei Jin, Haitao Mao, Zheng Li, Haoming Jiang, Chen Luo, Hongzhi Wen, Haoyu Han, Hanqing Lu, Zhengyang Wang, Ruirui Li, et al. Amazon-m2: A multilingual multi-locale shopping session dataset for recommendation and text generation. In *NeurIPS*, 2023.
- [35] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.
- [36] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- [37] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International conference on machine learning*, pages 1885–1894. PMLR, 2017.
- [38] John Boaz Lee, Ryan Rossi, and Xiangnan Kong. Graph classification using structural attention. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1666–1674, 2018.
- [39] Fan Li, Xiaoyang Wang, Dawei Cheng, Wenjie Zhang, Ying Zhang, and Xuemin Lin. Tcgu: Data-centric graph unlearning based on transferable condensation. *arXiv preprint arXiv:2410.06480*, 2024.
- [40] Guihong Li, Hsiang Hsu, Chun-Fu Chen, and Radu Marculescu. Machine unlearning for image-to-image generative models. In *The Twelfth International Conference on Learning Representations*, 2024.
- [41] Xunkai Li, Yulin Zhao, Zhengyu Wu, Wentao Zhang, Rong-Hua Li, and Guoren Wang. Towards effective and general graph unlearning via mutual evolution. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 13682–13690, 2024.
- [42] Yi Li, Shichao Zhang, Guixian Zhang, and Debo Cheng. Community-centric graph unlearning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 18548–18556, 2025.
- [43] Yuyuan Li, Chaochao Chen, Yizhao Zhang, Weiming Liu, Lingjuan Lyu, Xiaolin Zheng, Dan Meng, and Jun Wang. Ultrare: Enhancing receraser for recommendation unlearning via error decomposition. *Advances in Neural Information Processing Systems*, 36:12611–12625, 2023.
- [44] Minhua Lin, Enyan Dai, Junjie Xu, Jinyuan Jia, Xiang Zhang, and Suhang Wang. Stealing training graphs from graph neural networks. *arXiv preprint arXiv:2411.11197*, 2024.
- [45] Minhua Lin, Zhiwei Zhang, Enyan Dai, Zongyu Wu, Yilong Wang, Xiang Zhang, and Suhang Wang. Trojan prompt attacks on graph neural networks. *arXiv preprint arXiv:2410.13974*, 2024.
- [46] Chengyi Liu, Jiahao Zhang, Shijie Wang, Wenqi Fan, and Qing Li. Score-based generative diffusion models for social recommendations. *arXiv preprint arXiv:2412.15579*, 2024.
- [47] Chris Liu, Yaxuan Wang, Jeffrey Flanigan, and Yang Liu. Large language model unlearning via embedding-corrupted prompts. *Advances in Neural Information Processing Systems*, 37:118198–118266, 2024.
- [48] Junxu Liu, Mingsheng Xue, Jian Lou, Xiaoyu Zhang, Li Xiong, and Zhan Qin. Muter: Machine unlearning on adversarially trained models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4892–4902, 2023.

- [49] Sijia Liu, Yuanshun Yao, Jinghan Jia, Stephen Casper, Nathalie Baracaldo, Peter Hase, Yuguang Yao, Chris Yuhao Liu, Xiaojun Xu, Hang Li, et al. Rethinking machine unlearning for large language models. *Nature Machine Intelligence*, pages 1–14, 2025.
- [50] Zheyuan Liu, Guangyao Dou, Eli Chien, Chunhui Zhang, Yijun Tian, and Ziwei Zhu. Breaking the trilemma of privacy, utility, and efficiency via controllable machine unlearning. In *Proceedings of the ACM Web Conference 2024*, pages 1260–1271, 2024.
- [51] Zhiwei Liu, Yingdong Dou, Philip S Yu, Yutong Deng, and Hao Peng. Alleviating the inconsistency problem of applying graph neural network to fraud detection. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*, pages 1569–1572, 2020.
- [52] Alessandro Mantelero. The eu proposal for a general data protection regulation and the roots of the ‘right to be forgotten’. *Computer Law & Security Review*, 29(3):229–235, 2013.
- [53] Justus Mattern, Fatemehsadat Mireshghallah, Zhijing Jin, Bernhard Schoelkopf, Mrinmaya Sachan, and Taylor Berg-Kirkpatrick. Membership inference attacks against language models via neighbourhood comparison. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 11330–11343, 2023.
- [54] Matthieu Meeus, Shubham Jain, Marek Rei, and Yves-Alexandre de Montjoye. Did the neurons read your book? document-level membership inference for large language models. In *33rd USENIX Security Symposium (USENIX Security 24)*, pages 2369–2385, 2024.
- [55] Dan S Nielsen and Ryan McConville. Mumin: A large-scale multilingual multimodal fact-checked misinformation social network dataset. In *Proceedings of the 45th international ACM SIGIR conference on research and development in information retrieval*, pages 3141–3153, 2022.
- [56] California Office of the Attorney General. California consumer privacy act (ccpa), 2025.
- [57] Office of the Privacy Commissioner of Canada. Announcement: Privacy commissioner seeks federal court determination on key issue for Canadians’ online reputation, 2018.
- [58] Iyiola E Olatunji, Wolfgang Nejdl, and Megha Khosla. Membership inference attack on graph neural networks. In *2021 Third IEEE International Conference on Trust, Privacy and Security in Intelligent Systems and Applications (TPS-ISA)*, pages 11–20. IEEE, 2021.
- [59] Chao Pan, Eli Chien, and Olgica Milenkovic. Unlearning graph classifiers with limited data resources. In *Proceedings of the ACM Web Conference 2023*, pages 716–726, 2023.
- [60] Pedram Pedarsani and Matthias Grossglauser. On the privacy of anonymized networks. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1235–1243, 2011.
- [61] Anwar Said, Yuying Zhao, Tyler Derr, Mudassir Shabbir, Waseem Abbas, and Xenofon Koutsoukos. A survey of graph unlearning. *arXiv preprint arXiv:2310.02164*, 2023.
- [62] Sina Sajadmanesh and Daniel Gatica-Perez. Locally private graph neural networks. In *Proceedings of the 2021 ACM SIGSAC conference on computer and communications security*, pages 2130–2145, 2021.
- [63] Aravind Sankar, Yozen Liu, Jun Yu, and Neil Shah. Graph neural networks for friend ranking in large-scale social platforms. In *Proceedings of the Web Conference 2021*, pages 2535–2546, 2021.
- [64] Seiyun Shin, Ilan Shomorony, and Han Zhao. Efficient learning of linear graph neural networks via node subsampling. *Advances in Neural Information Processing Systems*, 36:55479–55501, 2023.
- [65] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pages 3–18. IEEE, 2017.

- [66] Ayush K Tarun, Vikram S Chundawat, Murari Mandal, and Mohan Kankanhalli. Fast yet effective machine unlearning. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [67] Anvith Thudi, Hengrui Jia, Iliia Shumailov, and Nicolas Papernot. On the necessity of auditable algorithmic definitions for machine unlearning. In *31st USENIX security symposium (USENIX Security 22)*, pages 4007–4022, 2022.
- [68] Stacey Truex, Ling Liu, Mehmet Emre Gursoy, Lei Yu, and Wenqi Wei. Demystifying membership inference attacks in machine learning as a service. *IEEE transactions on services computing*, 14(6):2073–2089, 2019.
- [69] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018.
- [70] Cheng-Long Wang, Mengdi Huai, and Di Wang. Inductive graph unlearning. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 3205–3222, 2023.
- [71] Jizhe Wang, Pipei Huang, Huan Zhao, Zhibo Zhang, Binqiang Zhao, and Dik Lun Lee. Billion-scale commodity embedding for e-commerce recommendation in alibaba. In *KDD*, 2018.
- [72] Shuang Wang, Muhammad Asif, Muhammad Farrukh Shahzad, and Muhammad Ashfaq. Data privacy and cybersecurity challenges in the digital transformation of the banking sector. *Computers & security*, 147:104051, 2024.
- [73] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. Neural graph collaborative filtering. In *Proceedings of the 42nd international ACM SIGIR conference on Research and development in Information Retrieval*, pages 165–174, 2019.
- [74] Xiyuan Wang, Haotong Yang, and Muhan Zhang. Neural common neighbor with completion for link prediction. In *The Twelfth International Conference on Learning Representations*, 2024.
- [75] Yilong Wang, Jiahao Zhang, Tianxiang Zhao, and Suhang Wang. Enhance gnns with reliable confidence estimation via adversarial calibration learning. *arXiv preprint arXiv:2503.18235*, 2025.
- [76] Yilong Wang, Tianxiang Zhao, Zongyu Wu, and Suhang Wang. Bridging source and target domains via link prediction for unsupervised domain adaptation on graphs. In *Proceedings of the Eighteenth ACM International Conference on Web Search and Data Mining*, pages 678–687, 2025.
- [77] Bang Wu, Xiangwen Yang, Shirui Pan, and Xingliang Yuan. Adapting membership inference attacks to gnn for graph classification: Approaches and implications. In *2021 IEEE International Conference on Data Mining (ICDM)*, pages 1421–1426. IEEE, 2021.
- [78] Bang Wu, He Zhang, Xiangwen Yang, Shuo Wang, Minhui Xue, Shirui Pan, and Xingliang Yuan. Graphguard: Detecting and counteracting training data misuse in graph neural networks. In *31st Annual Network and Distributed System Security Symposium (NDSS)*, 2024.
- [79] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph convolutional networks. In *International conference on machine learning*, pages 6861–6871. PMLR, 2019.
- [80] Jiancan Wu, Yi Yang, Yuchun Qian, Yongduo Sui, Xiang Wang, and Xiangnan He. Gif: A general graph unlearning strategy via influence function. In *Proceedings of the ACM Web Conference 2023*, pages 651–661, 2023.
- [81] Kun Wu, Jie Shen, Yue Ning, Ting Wang, and Wendy Hui Wang. Certified edge unlearning for graph neural networks. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 2606–2617, 2023.

- [82] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S Yu. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems*, 32(1):4–24, 2020.
- [83] Xin Xin, Jiyuan Yang, Hanbing Wang, Jun Ma, Pengjie Ren, Hengliang Luo, Xinlei Shi, Zhumin Chen, and Zhaochun Ren. On the user behavior leakage from recommender system exposure. *ACM Transactions on Information Systems*, 41(3):1–25, 2023.
- [84] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *ICLR*, 2019.
- [85] Haonan Yan, Xiaoguang Li, Ziyao Guo, Hui Li, Fenghua Li, and Xiaodong Lin. Arcane: An efficient architecture for exact machine unlearning. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 4006–4013, 7 2022.
- [86] Zhe-Rui Yang, Jindong Han, Chang-Dong Wang, and Hao Liu. Erase then rectify: A training-free parameter editing approach for cost-effective graph unlearning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 13044–13051, 2025.
- [87] Zhilin Yang, William Cohen, and Ruslan Salakhudinov. Revisiting semi-supervised learning with graph embeddings. In *International conference on machine learning*, pages 40–48. PMLR, 2016.
- [88] Yuanshun Yao, Xiaojun Xu, and Yang Liu. Large language model unlearning. *Advances in Neural Information Processing Systems*, 37:105425–105475, 2024.
- [89] Lu Yi and Zhewei Wei. Scalable and certifiable graph unlearning: Overcoming the approximation error barrier. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [90] Wei Yuan, Chaoqun Yang, Quoc Viet Hung Nguyen, Lizhen Cui, Tieke He, and Hongzhi Yin. Interaction-level membership inference attack against federated recommender systems. In *Proceedings of the ACM Web Conference 2023*, pages 1053–1062, 2023.
- [91] He Zhang, Bang Wu, Shuo Wang, Xiangwen Yang, Minhui Xue, Shirui Pan, and Xingliang Yuan. Demystifying uneven vulnerability of link stealing attacks against graph neural networks. In *International Conference on Machine Learning*, pages 41737–41752. PMLR, 2023.
- [92] He Zhang, Bang Wu, Xiangwen Yang, Xingliang Yuan, Xiaoning Liu, and Xun Yi. Dynamic graph unlearning: A general and efficient post-processing method via gradient transformation. In *Proceedings of the ACM on Web Conference 2025*, pages 931–944, 2025.
- [93] Jiahao Zhang. Graph unlearning with efficient partial retraining. In *Companion Proceedings of the ACM Web Conference 2024*, pages 1218–1221, 2024.
- [94] Jiahao Zhang, Rui Xue, Wenqi Fan, Xin Xu, Qing Li, Jian Pei, and Xiaorui Liu. Linear-time graph neural networks for scalable recommendations. In *Proceedings of the ACM Web Conference 2024*, pages 3533–3544, 2024.
- [95] Minxing Zhang, Zhaochun Ren, Zihan Wang, Pengjie Ren, Zhunmin Chen, Pengfei Hu, and Yang Zhang. Membership inference attacks against recommender systems. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, pages 864–879, 2021.
- [96] Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. *Advances in neural information processing systems*, 31, 2018.
- [97] Qiuchen Zhang, Carl Yang, Li Xiong, et al. Node-level contrastive unlearning on graph neural networks. *arXiv preprint arXiv:2503.02959*, 2025.
- [98] Yang Zhang, Zhiyu Hu, Yimeng Bai, Jiancan Wu, Qifan Wang, and Fuli Feng. Recommendation unlearning via influence function. *ACM Transactions on Recommender Systems*, 3(2):1–23, 2024.

- [99] Yihua Zhang, Chongyu Fan, Yimeng Zhang, Yuguang Yao, Jinghan Jia, Jiancheng Liu, Gaoyuan Zhang, Gaowen Liu, Ramana Kompella, Xiaoming Liu, and Sijia Liu. Unlearncanvas: A stylized image dataset to benchmark machine unlearning for diffusion models. *NeurIPS*, 2024.
- [100] Yimeng Zhang, Xin Chen, Jinghan Jia, Yihua Zhang, Chongyu Fan, Jiancheng Liu, Mingyi Hong, Ke Ding, and Sijia Liu. Defensive unlearning with adversarial training for robust concept erasure in diffusion models. *Advances in Neural Information Processing Systems*, 37:36748–36776, 2024.
- [101] Zaixi Zhang, Qi Liu, Zhenya Huang, Hao Wang, Chee-Kong Lee, and Enhong Chen. Model inversion attacks against graph neural networks. *IEEE Transactions on Knowledge and Data Engineering*, 35(9):8729–8741, 2022.
- [102] Zaixi Zhang, Qi Liu, Zhenya Huang, Hao Wang, Chengqiang Lu, Chuanren Liu, and Enhong Chen. Graphmi: Extracting private graph data from graph neural networks. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, 2021.
- [103] Zhikun Zhang, Min Chen, Michael Backes, Yun Shen, and Yang Zhang. Inference attacks against graph neural networks. In *31st USENIX Security Symposium (USENIX Security 22)*, pages 4543–4560, 2022.
- [104] Zhilu Zhang and Mert Sabuncu. Generalized cross entropy loss for training deep neural networks with noisy labels. *Advances in neural information processing systems*, 31, 2018.
- [105] Zhiwei Zhang, Fali Wang, Xiaomin Li, Zongyu Wu, Xianfeng Tang, Hui Liu, Qi He, Wenpeng Yin, and Suhang Wang. Catastrophic failure of LLM unlearning via quantization. In *The Thirteenth International Conference on Learning Representations*, 2025.