# On Strengthening and Defending Graph Reconstruction Attack with Markov Chain Approximation

Zhanke Zhou[1]  Chenyu Zhou[1]  Xuan Li[1]  Jiangchao Yao[2][3]  Quanming Yao[4]  Bo Han[1]

## Abstract

Although powerful graph neural networks (GNNs) have boosted numerous real-world applications, the potential privacy risk is still under-explored. To close this gap, we perform the first comprehensive study of *graph reconstruction attack* that aims to reconstruct the *adjacency* of nodes. We show that a range of factors in GNNs can lead to the surprising leakage of private links. Especially by taking GNNs as a Markov chain and attacking GNNs via a flexible chain approximation, we systematically explore the underneath principles of graph reconstruction attack, and propose two information theory-guided mechanisms: (1) the chain-based attack method with adaptive designs for extracting more private information; (2) the chain-based defense method that sharply reduces the attack fidelity with moderate accuracy loss. Such two objectives disclose a critical belief that to recover better in attack, you must extract more multi-aspect knowledge from the trained GNN; while to learn safer for defense, you must forget more link-sensitive information in training GNNs. Empirically, we achieve state-of-the-art results on six datasets and three common GNNs. The code is publicly available at: https://github.com/tmlr-group/MC-GRA.

## 1. Introduction

Deep learning has promoted tremendously broad research from Euclidean data like images to non-euclidean data like graphs. Specifically, graph neural networks (GNNs) (Kipf & Welling, 2016a; Gilmer et al., 2017; Zhang & Chen, 2018)

[1]Department of Computer Science, Hong Kong Baptist University [2]Cooperative Medianet Innovation Center, Shanghai Jiao Tong University [3]Shanghai AI Laboratory [4]Department of Electronic Engineering, Tsinghua Unversity. Correspondence to: Bo Han <bhanml@comp.hkbu.edu.hk>, Jiangchao Yao <Sunarker@sjtu.edu.cn>.
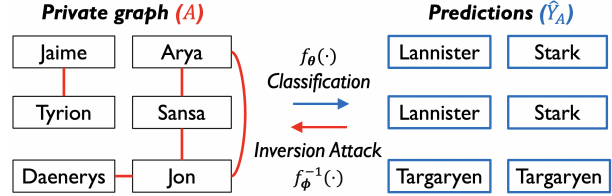
Figure 1: An illustration of Graph Reconstruction Attack. The *forward* inference of a trained model is to predict the node category $\hat{Y}_A$, *i.e.*, the family name of each character; while the *backward* inversion attack is to recover the original adjacency $A$, *i.e.*, the kinship among characters (red edges).

proposed in the recent years have drawn much attention and boosted a wide range of real-world applications, *e.g.*, social network (Fan et al., 2019), recommender systems (Wu et al., 2020a) and drug discovery (Ioannidis et al., 2020).

Nevertheless, the privacy concerns behind these applications raise with the development of the *Model Inversion Attack* (MIA) technique, which only requires a trained model and non-sensitive features to recover the sensitive information. In particular, recent progress on MIA (Fredrikson et al., 2015; Zhang et al., 2020; Struppek et al., 2022) has shown the feasible recovery of private images in high fidelity and diversity. As for the scenarios of GNNs, the similar inversion of the adjacency of the training graph is also a severe privacy threat, since links can reflect the sensitive relationship information or intellectual properties of the model's owner. We term this kind of MIA on graphs as Graph Reconstruction Attack (GRA) for simplicity and illustrate exemplars in Fig. 1. To date, only limited research has been conducted on GRA (He et al., 2021a; Zhang et al., 2021b) that is designed for ad-hoc scenarios. The general principles for strengthening and defending GRA are still unknown, which presents hidden dangers in extensive real-world applications. Thus, it is urgent to understand the vulnerability of GNNs under such attacks and explore the proper defense methods to protect GNNs and avoid privacy risks in advance.

In this work, we systematically investigate this crucial yet under-explored problem from both sides of attack and defense. Roughly, the GNNs' inference procedure can be viewed as a Markov Chain $f_{\theta} : (A, X) \rightarrow H_A \rightarrow \hat{Y}_A \leftrightarrow Y$, where the adjacent matrix $A$ and node features $X$ are taken
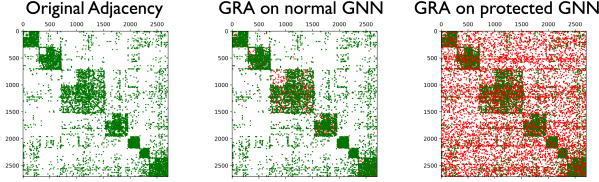
Figure 2: Recovered adjacency on Cora dataset. Green dots are correctly predicted edges while red dots are wrong ones. GRA on normal GNN leads to privacy leakage, while GRA on protected GNN cannot recover the private adjacency.
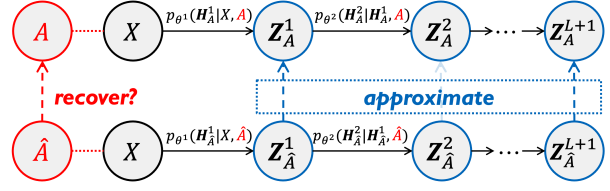


Figure 3: Modeling the GRA problem as approximating the original Markov chain (upper) by the attack chain (lower). Note that the original chain is with the original adjacency $A$, while the attack chain is with the recovered adjacency $\hat{A}$.

as the inputs to generate node embeddings $H_A$, and a linear layer with activation transforms $H_A$ into the classification outputs $\hat{Y}_A$ to predict node labels $Y$. More importantly, we reveal that every variable in $\{X, Y, H_A, \hat{Y}_A\}$ can recover adjacency to a certain extent through a simple transformation. However, different from the single variable in MIA for images, it is mysterious to understand the intriguing mechanism behind the multiple interplaying factors in GNNs, thus challenging to apply for strengthening and defending GRA.

To close the gap, we formulate the GRA problem from a novel perspective, *i.e.*, approximating the original Markov chain by the attack chain (Fig. 3). Note that such a modeling manner brings three-fold advantages: (1) adaptively supports the white-box attack that utilizes any set of prior knowledge; (2) help derive the chain-based attack and defense objectives in optimization; (3) enables analysis from the information-theoretical view. On the basis of the chain modeling, we investigate the underneath principles of the GRA problem, which are two folds. To strengthen the attack, we derive the Markov Chain-based Graph Reconstruction Attack (MC-GRA) that simulates the hidden transformation procedure of the target GNN by approximating all the known informative variables in a combinatorial manner. As for defense, we propose the Markov Chain-based Graph Privacy Bottleneck (MC-GPB), which regularizes the mutual dependency among graph representations, adjacency, and labels to alleviate privacy leakage, as shown in Fig. 2.

In short, our main contributions are summarized as follows. **(1)** To our best knowledge, we are the first to conduct a systematic study of GRA and reveal several essential and useful phenomenons (Sec. 4). **(2)** On the basis of the chain modeling, we propose a new method for the attack that boosts the attack fidelity with parameterization techniques and injected stochasticity (Sec. 5), and propose an information theory-guided principle for the defense that significantly degenerates all the attacks with only a slight accuracy loss (Sec. 6). **(3)** We provide a rigorous analysis from information-theoretical perspectives to disclose several valuable insights on what and how to strengthen and defend GRA. **(4)** Both the two proposed methods achieve state-of-the-art results on six datasets and three GNNs (Sec. 7).

## 2. Related Work

**Inversion attacks on images.** Pioneer works (Szegedy et al., 2013; Fredrikson et al., 2014; 2015; Hidano et al., 2017) introduced the Model Inversion Attack (MIA) with shallow models and justified the feasibility of MIA in recovering the monochrome images. However, they fail in attacking deep models for image classification tasks, where the reconstructed images are of low fidelity. Generative model inversion (Zhang et al., 2020) is the first to conduct MIA on convolution neural networks. Instead of directly reconstructing the private data from scratch, its inversion process is guided by a distributional prior through the generative adversarial networks (GAN) that can reveal private training data of the target model with high fidelity. Later, variational model inversion (Wang et al., 2021) further formulates the MIA as the variational inference. It generally can bring a higher attack accuracy and diversity for its equipped powerful generator StyleGAN to optimize its designed variational objective. Recent advance (Struppek et al., 2022) significantly decreases the cost of conducting MIA through relaxing the dependency between the target model and the image prior. It enables the use of a single GAN to attack a wide range of targets, requiring only minor adjustments to the attack. It shows that MIA is possible with publicly available pre-trained GANs under strong distributional shifts.

**Inversion attacks on graphs.** Early works (Duddu et al., 2020; Chanpuriya et al., 2021) attempt to reconstruct the target graph from released graph embeddings of each node that are generated by Deepwalk or GNNs. The link stealing attack (He et al., 2021a) is the first work to steal links from a GNN as the target model. It aims to conduct the MIA with three kinds of prior knowledge, including node features, partial target graph, and a shadow dataset. It considered all permutations of the three elements and proposed eight kinds of attack methods in total that are adaptive to the eight scenarios with chemical networks and social networks. Another recent work (Zhang et al., 2021b) is a learnable attack that also aims to recover the links of the original graph. With the white-box access to the target GNN model, the optimal adjacency is obtained by maximizing the classification accuracy regarding the known node labels. Please refer to Appendix C for a detailed introduction to related work.

## 3. Preliminaries

**Notations.** With adjacent matrix $A$ and node features $X$, an undirected graph is denoted as $\mathcal{G} = (A, X)$, where $A_{ij} = 1$ means there is an edge $e_{ij}$ between $v_i$ and $v_j$. For each node $v_i$, its $D$-dimension node feature is denoted as $X_{[i,:]} \in \mathbb{R}^D$, and its label $y_i \in Y$ indicates the node category. The node classification task is to predict the label $Y$ of each node via a parameterized model $f_{\boldsymbol{\theta}}(\cdot)$. $I(X; Y)$ indicates the mutual information of variables $X$ and $Y$. We summarize the frequently used notations in Table 10 of Appendix. A.1

**Graph neural networks.** Predicting node labels requires a parameterized hypothesis $f_{\boldsymbol{\theta}}$ with GNN architecture and the message propagation framework (Gilmer et al., 2017). Specifically, the forward inference of a $L$-layer GNN generates the node representations $\boldsymbol{H}_A \in \mathbb{R}^{N \times D}$ by a $L$-layer message propagation. The follow-up linear layer transforms the representations $\boldsymbol{H}_A$ to the classification probabilities $\hat{\boldsymbol{Y}}_A \in \mathbb{R}^{N \times C}$, with $C$ categories of nodes in total.

**Model inversion attack on graphs.** In this study, to catch more attention to the privacy risk of GNNs, we study the reconstruction of the graph adjacency by MIA and term it Graph Reconstruction Attack (GRA), as elaborated below.

**Definition 3.1** (Graph Reconstruction Attack). Given a set of prior knowledge $\mathcal{K}$ and a trained GNN $f_{\boldsymbol{\theta}^*}(\cdot)$, the graph reconstruction attack aims to recover the original linking relations $\hat{\boldsymbol{A}}^*$ of the training graph $\mathcal{G}_{\text{train}} = (A, X)$, namely,

$$\text{GRA:} \quad \hat{\boldsymbol{A}}^* = \arg\max_{\hat{\boldsymbol{A}}} \mathbb{P}(\hat{\boldsymbol{A}} | f_{\boldsymbol{\theta}^*}, \mathcal{K}). \quad (1)$$

Here, $\mathbb{P}(\cdot)$ is the attack method to generate $\hat{\boldsymbol{A}}$, and $\mathcal{K}$ can be any subset of $\{X, Y, \boldsymbol{H}_A, \hat{\boldsymbol{Y}}_A\}$. Note that GRA is conducted in a post-hoc manner, *i.e.*, after the training of GNN $f_{\boldsymbol{\theta}}(\cdot)$.

## 4. A Comprehensive Study of GRA

In this section, we formulate the Graph Reconstruction Attack as a Markov chain approximation problem (Sec. 4.1), quantify the privacy risk of releasing the non-sensitive features (Sec. 4.2), and investigate the training dynamics of graph representations *w.r.t.* the privacy leakage (Sec. 4.3).

### 4.1. Modeling GRA as Markov chain approximation.

To adaptively support the white-box GRA that leverages the target model and any prior knowledge; and to properly locate, present, and utilize the interplaying variables of GNN forward in a generic manner; we cast the GRA problem as approximating the original Markov chain ORI-chain by the attack chain GRA-chain, as shown in Fig. 3, namely,

$$\texttt{ORI-chain:} \boldsymbol{Z}^0 \xrightarrow[\boldsymbol{\theta}^1]{A} \boldsymbol{Z}_A^1 \xrightarrow[\boldsymbol{\theta}^2]{A} \boldsymbol{Z}_A^2 \rightarrow \cdots \xrightarrow[\boldsymbol{\theta}^{L+1}]{A} \boldsymbol{Z}_A^{L+1},$$

$$\texttt{GRA-chain:} \boldsymbol{Z}^0 \xrightarrow[\boldsymbol{\theta}^1]{\hat{\boldsymbol{A}}} \boldsymbol{Z}_{\hat{\boldsymbol{A}}}^1 \xrightarrow[\boldsymbol{\theta}^2]{\hat{\boldsymbol{A}}} \boldsymbol{Z}_{\hat{\boldsymbol{A}}}^2 \rightarrow \cdots \xrightarrow[\boldsymbol{\theta}^{L+1}]{\hat{\boldsymbol{A}}} \boldsymbol{Z}_{\hat{\boldsymbol{A}}}^{L+1},$$

$$(2)$$

Table 1: Quantitative analysis of $I(A; \boldsymbol{Z})$ with AUC metric under range $[0, 1]$. A higher AUC value means a severer privacy leakage. "—" indicates that nodes in this dataset do not have features. Besides, the **boldface** numbers mean the best results, while the underlines indicate the second-bests. The target model $f_{\boldsymbol{\theta}}$ is a two-layer GCN by default.

| MI | Cora | Citeseer | Polblogs | USA | Brazil | AIDS |
|---|---|---|---|---|---|---|
| $I(A; X)$ | <u>.781</u> | **.881** | — | — | — | .521 |
| $I(A; \boldsymbol{H}_A)$ | .766 | .760 | **.763** | **.850** | **.758** | **.584** |
| $I(A; \hat{\boldsymbol{Y}}_A)$ | .712 | .743 | <u>.772</u> | <u>.826</u> | <u>.732</u> | <u>.561</u> |
| $I(A; Y)$ | **.815** | <u>.779</u> | .705 | .728 | .613 | .536 |

Table 2: An ensemble study on the prior knowledge with AUC metric. For a generic evaluation, it is assumed that node feature $X$ is accessible (if exists), based on which we evaluate all the possible 8 combinations with 2, 3, or 4 components, where "✓" means accessible for this variable.

| $X$ | $\boldsymbol{H}_A$ | $\hat{\boldsymbol{Y}}_A$ | $Y$ | Cora | Citeseer | Polblogs | USA | Brazil | AIDS |
|---|---|---|---|---|---|---|---|---|---|
| ✓ | ✓ | | | .781 | .881 | .763 | .850 | .758 | .521 |
| ✓ | | ✓ | | .781 | .881 | .772 | .826 | .732 | .521 |
| ✓ | | | ✓ | .849 | .907 | .705 | .728 | .613 | .522 |
| ✓ | ✓ | ✓ | | .781 | .881 | .763 | .848 | .756 | .521 |
| ✓ | ✓ | | ✓ | .849 | .907 | .779 | .850 | .743 | .522 |
| ✓ | | ✓ | ✓ | .842 | .907 | .785 | .842 | .730 | .522 |
| ✓ | ✓ | ✓ | ✓ | .849 | .907 | .781 | .852 | .717 | .522 |

where $\boldsymbol{Z}^0 = X$, $\boldsymbol{Z}_A^i = \boldsymbol{H}_A^i$ for $i = 1, \cdots, L$ and $\boldsymbol{Z}_A^{L+1} = \hat{\boldsymbol{Y}}_A$. Note that GNNs' forward can be seen as a Markov chain that is discrete-time finite, non-reversible, and pairwise-independent. The probability of current state $\boldsymbol{H}_A^i$ only depends on the previous state $\boldsymbol{H}_A^{i-1}$, where the transition kernel is determined by $A$ and $\boldsymbol{\theta}^i$. Importantly, the principle of GRA to recover the adjacency $A$ by $\hat{\boldsymbol{A}}$ is to approximate latent variables $\mathcal{S}_A = \{\boldsymbol{Z}_A^i : \boldsymbol{Z}_A^i \in \mathcal{K}\}$ in ORI-chain by the corresponding $\mathcal{S}_{\hat{\boldsymbol{A}}} = \{\boldsymbol{Z}_{\hat{\boldsymbol{A}}}^i : \boldsymbol{Z}_A^i \in \mathcal{S}_A\}$ in GRA-chain.

### 4.2. What leaks privacy in ORI-chain?

Intuitively, variables in ORI-chain might contain information about ground-truth adjacency $A$ as the transition kernel is partially determined by $A$. To figure out, we quantify the direct correlation between $A$ and a single variable $\boldsymbol{Z} \in \{X, Y, \boldsymbol{H}_A, \hat{\boldsymbol{Y}}_A\}$ in ORI-chain through the informative concept of mutual information (MI) $I(A; \boldsymbol{Z})$. Following link prediction works (Zhang & Chen, 2018; Zhu et al., 2021), the AUC metric is utilized to quantify $I(A; \boldsymbol{Z})$ regarding edges in $A$ and $\hat{A}_{\boldsymbol{Z}} = \sigma(\boldsymbol{Z}\boldsymbol{Z}^\top)$, where $\sigma(\cdot)$ is the activation function. Here, the inner product transforms the informative variable $\boldsymbol{Z} \in \mathbb{R}^{N \times D}$ to the predictive adjacency $\hat{A}_{\boldsymbol{Z}} \in \mathbb{R}^{N \times N}$, where the $(i, j)$ entry in $\hat{A}_{\boldsymbol{Z}}$ indicate the existence of edge $e_{ij}$. See Appendix. E.1 for details.

**Observation 3.1.** As shown in Tab. 1, a single variable in ORI-chain can recover the original adjacency to a
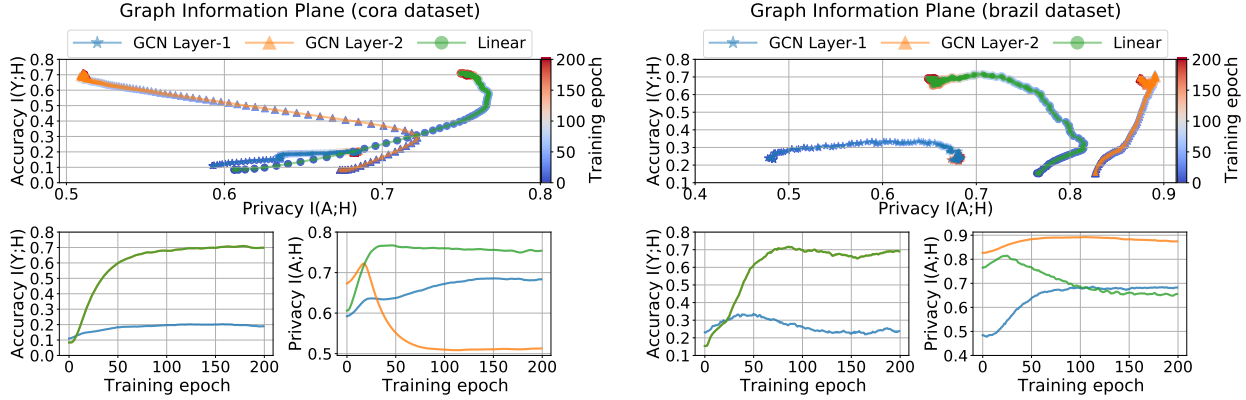
Figure 4: Graph information plane: tracking the standard training procedures of a two-layer GCN on Cora (left) and Brazil (right). The accuracy of GCN layer-2 and the linear layer is the same as $\hat{Y}_A = \texttt{Linear}(H_A^2)$ (thus with overlapped curves).

certain extent through the inner-product transformation. It is applicable to black-box attacks once obtain these variables. Besides, the model outputs $\{H_A, \hat{Y}_A\}$ generally contain more adjacency information than the original data $\{X, Y\}$.

As single variables $Z$ in ORI-chain present diverse approximation power, the stored private information might be complementary to each other in recovering adjacency. To answer, we ensemble these variables via a linear combination, namely, $\hat{A}_{esm} = 1/|\mathcal{K}| \sum_{i=1}^{|\mathcal{K}|} \hat{A}_{\mathcal{K}_i}$, where $\hat{A}_{\mathcal{K}_i} = \sigma(\mathcal{K}_i \mathcal{K}_i^\top)$.

**Observation 3.2.** As shown in Tab. 2, the straightforward and linear combination of informative terms only brings marginal improvements in recovering the adjacency. Such an observation is consistent with the chain rule of MI, *i.e.*, $\forall \mathcal{K}_i, \mathcal{K}_j \in \mathcal{K}, \ I(A; \mathcal{K}_i, \mathcal{K}_j) \geq \max\big(I(A; \mathcal{K}_i), I(A; \mathcal{K}_j)\big)$.

### 4.3. How `ORI-chain` memorizes the privacy?

For further understanding of the learning and memorization mechanisms of ORI-chain and acquiring inspiration for devising the corresponding defense approach, we track the training process by privacy $I(A; Z)$ and accuracy $I(Y; Z)$, where variable $Z \in \{H_A^1, H_A^2, \hat{Y}_A\}$ are from ORI-chain. Conceptually, we derive Graph Information Plane [1] inspired by information theory (Tishby & Zaslavsky, 2015; Shwartz-Ziv & Tishby, 2017). The anytime $Z$ in training phase is projected to the two-dimensional $\big(I(A; Z), I(Y; Z)\big)$ plane.

**Observation 3.3.** As shown in Fig. 4, the training procedure with v-shape curves contains two main phases: *fitting* and *compressing*. In the first and shorter phase, the layers increase the information about privacy. While in the second and longer phase, the layers gradually forget about privacy.

## 5. To Recover Better, You Must Extract More

To attack, one must integrate all the available prior knowledge to backward recover the adjacency. The key challenge

here is the lack of an effective way to employ all the prior knowledge and the target model in attacks. Besides, it is also hard to represent and update the recovered adjacency in a differentiable way due to the discrete nature of adjacency.

To solve, we propose the Markov Chain-based Graph Reconstruction Attack (MC-GRA) framework, as illustrated in Fig. 5(a). Here, instead of directly maximizing $I(\hat{A}; \mathcal{K})$, we choose to promote $I(Z_{\hat{A}}^i; Z_A^i = \mathcal{K}_i)$ as it provides supervision signals that can be tractably approximated. To be specific, we adopt the aforementioned chain-based modeling for extracting the knowledge stored in the target model while utilizing all the prior knowledge for optimization simultaneously. [2] The relaxation power hails from approaching the known variable $\mathcal{K}_i$ of ORI-chain by the locationally corresponding $Z_{\hat{A}}^i$ generated by GRA-chain, namely,
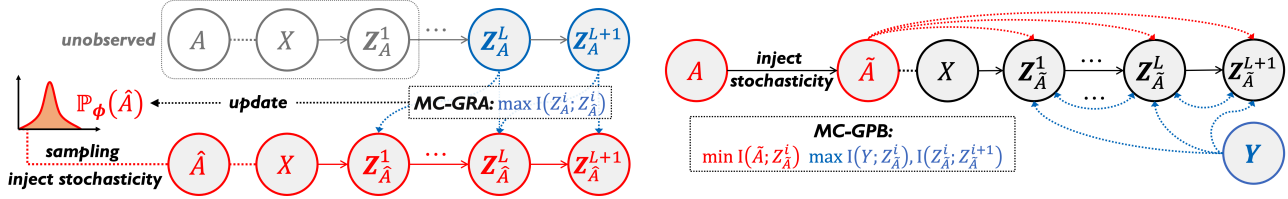
$$\text{MC-GRA: } \hat{A}^* = \arg\max_{\hat{A}} \ \underbrace{\alpha_p I(H_A; H_{\hat{A}}^i)}_{\text{propagation approximation}}$$
$$+ \underbrace{\alpha_o I(Y_A; Y_{\hat{A}}) + \alpha_s I(Y; Y_{\hat{A}})}_{\text{outputs approximation}} - \underbrace{\alpha_c H(\hat{A})}_{\text{complexity}}. \quad (3)$$

Note that MC-GRA is a maximin game: it maximizes the approximation of forward processes of the two Markov chains, while minimizing the complexity in each transition with $\hat{A}$ to avoid trivial solutions by constraining the density.

*Remark* 5.1. The adaptive power of MC-GRA comes from its leveraging any prior knowledge set. That is, the propagation approximation term in Eq. (3) for $H_A$ works once obtained, while the outputs approximation term for $\hat{Y}_A$ and $Y$. Thus, it can be utilized for all the 7 settings in Tab. 2.

**Parameterize Eq. (3) with different forms.** For approximating the original adjacency in a learnable manner, the recovered adjacency is parameterized and updated with the relaxed objective. Each time forward, an adjacency $\hat{A}$ is sampled from its parameterized distribution as $\hat{A} \sim \mathbb{P}_\phi(\hat{A})$.

---

[1] We leave the formal definition and details in Appendix. E.2.

[2] The detailed deriving is elaborated in Appendix. D.5.

(a) The attack framework MC-GRA. In forward, a recovered adjacency $\hat{A}$ is sampled from the parameterized distribution $\mathbb{P}_\phi(\hat{A})$ and injected with manual stochasticity. As for backward, the learnable parameters $\phi$ gain supervision from the MC-GRA objective Eq. (3).

(b) The defense framework MC-GPB. It solves the accuracy-privacy tradeoff by objective Eq. (4) through regularizing graph representations to make GNNs forget about private $A$ and injecting stochasticity to promote forgetting that decreases the privacy risk further.

Figure 5: Illustrations of the two proposed methods for strengthening (a) and defending (b) the GRA, respectively.

Technically, three implementations of $\mathbb{P}_\phi(\hat{A})$ with learnable weights $\phi$ are listed below with increasing complexity.

- Formulating $\hat{A}$ as the only learnable parameter and directly optimizing it, *i.e.*, $\mathbb{P}_\phi(\hat{A}) = \hat{A} \in [0,1]^{N \times N}$.

- A Gaussian distribution $\mathbb{P}_\phi = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$ with two learnable parameters $\boldsymbol{\mu}, \boldsymbol{\sigma} \in [0,1]^{N \times N}$. That is utilized to generate $\hat{A}$ as $\hat{A} = \boldsymbol{\mu} + \epsilon\boldsymbol{\sigma}$, where random noise $\epsilon \sim \mathcal{N}(0,1)$.

- A parameterized generator $f_\phi(\cdot)$ initialized with the same architecture and weights as $f_{\theta^*}(\cdot)$. It generates the probabilistic distribution by $\mathbb{P}_\phi = \sigma(\boldsymbol{H}_I \boldsymbol{H}_I^\top) \in [0,1]^{N \times N}$, where $I$ is the identity matrix and $\boldsymbol{H}_I = f_\phi(I, X)$.

**Optimize Eq. (3) with injected stochasticity.** Considering that both $\hat{A}$ and $X$ contribute to to to $\boldsymbol{Z} \in \{\boldsymbol{H}_A, \hat{\boldsymbol{Y}}_A, Y\}$, the mutual dependence among these three variables is coupled together. The spurious correlation $I(X; \hat{A}|\boldsymbol{Z})$, possibly degenerates the effectiveness of GRA (Yang et al., 2022; Miao et al., 2022). To solve, we inject stochasticity to further remove the spurious correlation among $\hat{A}$, $X$ and $\boldsymbol{Z}$, where the probability of spurious correlation naturally increases with the length of the Markov chain. Specifically, the debias power comes from the lower MI $I(\tilde{X}; \tilde{A}|\boldsymbol{Z})$, where $\tilde{X}, \tilde{A}$ are perturbed as $\tilde{X} = X \oplus X_\epsilon$, $\tilde{A} = \hat{A} \oplus A_\epsilon$. Technically, for each potential edge $e_{ij}$, its existence $a_{ij}$ is sampled from a Bernoulli distribution, *i.e.*, $a_{ij} \sim \text{Bern}(\boldsymbol{p}_{ij})$, $a_{ij} \in \{0,1\}$, and $\boldsymbol{p}_{ij} = \hat{A}_{ij} \in [0,1]$. To cooperate with the stochasticity and enable the back-propagation of gradients, the Gumbel-softmax reparameterization (Kool et al., 2019; Xie & Ermon, 2019) is applied. That is, the edge probabilities are perturbed as $\tilde{\boldsymbol{p}}_{ij} = \boldsymbol{p}_{ij} - \log(-\log(\epsilon))$, where $\epsilon \sim \text{Uniform}(0,1)$.

*Remark* 5.2. The incremental contribution of $\hat{A}$ regarding $X$ to approximate $\boldsymbol{Z}$ is $I(\boldsymbol{Z}; \hat{A}|X) = I(\boldsymbol{Z}; A, X) - I(\boldsymbol{Z}; X)$. Here, a general solution for obtaining a more informative $\hat{A}$ is to reduce $I(\boldsymbol{Z}; X)$ via perturbation and promote $I(\boldsymbol{Z}; \hat{A})$.

**Theoretical analysis about Eq. (3).** A rigorous analysis is conducted on the basis of information properties shown in Fig. 6. The non-invertible nature of GNN forwarding, which hails from the adopted non-linear operations, decreases the information entropy by layers and forms a bottleneck that

extracts informative signals from the input data. As a result, the MI of two Markov chains (Eq. (2)) is decreasing by layers, which is elaborated in the following Theorem 5.3. Based on this, we derive a tractable bound in Theorem 5.4 to estimate the attack fidelity without the ground truth $A$.

**Theorem 5.3.** *The layer-wise transformations $\boldsymbol{Z}_A^i \to \boldsymbol{Z}_A^{i+1}$ are non-invertible, e.g., $\boldsymbol{Z}_A^{i+1} = \sigma(\psi(A) \cdot \boldsymbol{Z}_A^i \cdot \boldsymbol{\theta}^i)$, where $\psi(A)$ is the graph convolution kernel, as in Eq. (2). It leads to a lower MI between the two Markov chains, i.e., $I(\boldsymbol{Z}_A^i; \boldsymbol{Z}_{\hat{A}}^i) - I(\boldsymbol{Z}_A^{i+1}; \boldsymbol{Z}_{\hat{A}}^{i+1}) \geq 0$. Proof. See Appendix.A.3.*

**Theorem 5.4** (Tractable Lower Bound of Fidelity). *The attack fidelity satisfies $I(A; \hat{A}) \geq H(\boldsymbol{H}_A) - H_b(e) - P(e)\log(|\mathcal{H}|)$, where $P(e) \triangleq P(\boldsymbol{H}_A \neq \boldsymbol{H}_{\hat{A}})$ is the probability of approximation error, $\mathcal{H}$ denotes the support of $\boldsymbol{H}_A$, and $H_b(\cdot)$ is the binary entropy. Proof. See Appendix. A.4.*

The estimated $I(A; \hat{A})$ can be a valuable reference when conducting GRA that maximizes such a MI term (see Fig. 6(b)). Besides, Theorem 5.4 also indicates that a higher approximation $I(\boldsymbol{H}_A; \boldsymbol{H}_{\hat{A}})$ with a lower error $P(e)$ can bring a higher $I(A; \hat{A})$ that indicates a higher attack fidelity. Then, we indicate the worst privacy leakage with the optimal attack fidelity as the upper bound in following Theorem 5.5.

**Theorem 5.5** (The Optimal Fidelity). *The recovering fidelity satisfies $I(A; X, Y, \boldsymbol{H}_A) - I(A; \hat{A}) \geq 0$. Solving MC-GRA sufficiently yields a solution to achieve the optimal case, i.e., $I(A; \hat{A}^*) = I(A; X, Y, \boldsymbol{H}_A)$. Proof. See Appendix. A.5.*

Theorem 5.5 indicates that MC-GRA is capable of achieving optimal recovering fidelity. Nonetheless, the remaining information of $A$, *i.e.*, $H(A|\hat{A}^*) = H(A|\mathcal{K})$, is unobservable from $\mathcal{K} = \{X, Y, \boldsymbol{H}_A\}$. Such information refers to the non-overlapping area of $A$ shown in Fig. 6(b), which cannot be recovered unless additional information is provided.

## 6. To Learn Safer, You Must Forget More

Recall in Sec. 4, graph representations naturally comprise the connectivity information, while the graph information plane shows that increasing privacy information is stored in the training phase. So, how can GNNs be *GRA-resistant*?
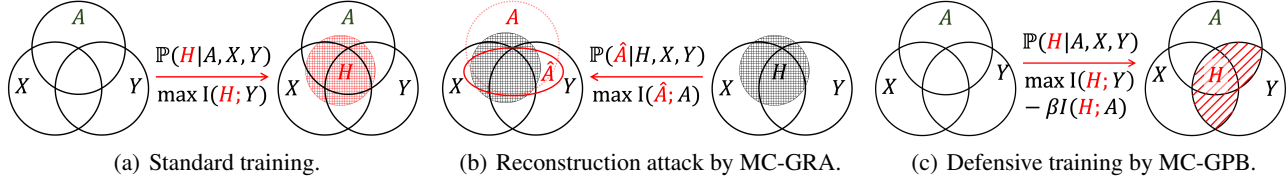
Figure 6: Illustrations of the information properties regarding the training, attacking, and defending processes.

(a) Standard training.  (b) Reconstruction attack by MC-GRA.  (c) Defensive training by MC-GPB.

For defense, one must require the GNN to *forget* the privacy information in the training process, *i.e.*, make the learned representations contain less information about adjacency. Nonetheless, it could easily degenerate the accuracy as the adjacency also essentially supports the prediction. To solve the trade-off, we proposed the *Markov Chain-based Graph Privacy Bottleneck* (MC-GPB) framework to defend against GRA (see Fig. 5(b)). Intuitively, the expected graph representations should come from a refined training process that learn the $\boldsymbol{\theta}^*$ from the original data $A, X, Y$. Inspired by the principle that *to learn, you must forget* by the information bottleneck (Tishby et al., 2000; Shwartz-Ziv & Tishby, 2017; Wu et al., 2020b) that constrains the data compression procedure $X \to Z \to Y$, we derive the defense objective as

$$\text{MC-GPB:} \boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^{L} \underbrace{-I(Y; \boldsymbol{H}_A^i)}_{\text{accuracy}} + \underbrace{\beta_p^i I(A; \boldsymbol{H}_A^i)}_{\text{privacy}}$$
$$+ \sum_{i=1}^{L-1} \underbrace{\beta_c^i I(\boldsymbol{H}_A^i; \boldsymbol{H}_A^{i+1})}_{\text{complexity}}. \quad (4)$$

Note that MC-GPB is also a maximin game: the correlation between hidden representations and labels is maximized, while that with adjacency is minimized instead. Analytically, it aims to minimize the conditional MI $I(A; \boldsymbol{H}_A^i|Y)$ through balancing accuracy $I(Y; \boldsymbol{H}_A^i)$ and privacy $I(A; \boldsymbol{H}_A^i)$. And the transformation complexity $I(\boldsymbol{H}_A^i; \boldsymbol{H}_A^{i+1})$ is constrained to relieve the smoothing effect of message propagation.

**Promote forgetting in Eq.(4) with injected stochasticity.** Making GNNs forget more about adjacency leads to lower privacy risk. For simplicity, the DropEdge (Rong et al., 2020) method is adopted, which performs random drop with probability $p$ on each observed edge of $A$. The perturbed adjacency $\tilde{A} = A \oplus A_\epsilon : A_\epsilon \perp\!\!\!\perp A, Y, \boldsymbol{Z}$, which satisfies $I(\tilde{A}; Y) \leq I(A; Y)$ and $I(\tilde{A}; \boldsymbol{Z}) \leq I(A; \boldsymbol{Z})$ (You et al., 2020; 2022). The injected stochasticity enforces the GNN model to discriminate the essential topological information $I(A; Y)$, rather than fully capturing the association between $A$ and $Y$ that can be potentially spurious (Zhao et al., 2022). As such, the redundancy $I(\tilde{A}; \boldsymbol{Z}|Y)$ is compressed to preserve privacy and maintain accuracy simultaneously.

**Promote feasibility via differentiable measurements.** Solving Eq. (3) and Eq. (4) requires tractable objectives. Given two variables, $X \in \mathbb{R}^{N \times D_x}$ and $Y \in \mathbb{R}^{N \times D_y}$, we cal-

culate the similarity $s(X, Y)$ to approximate $I(X, Y)$ considering six differentiable measurements (Kornblith et al., 2019). Technical details can be found in Appendix. E.3.

**Theoretical analysis about Eq. (4).** Regularizing the graph representations $\boldsymbol{H}_A$ with a lower $I(A; \boldsymbol{H}_A)$ indicates a lower $I(A; X, Y, \boldsymbol{H}_A)$, and thus, the optimal fidelity $I(A; \hat{\boldsymbol{A}}^*)$ is also decreased (refer to Theorem 5.5). Note that accuracy is prior to privacy in optimization with trade-offs, which corresponds to the concept of sufficient statistics.

**Proposition 6.1** (Sufficient Statistics). *Denote the sufficient statistics of $X$ as $\boldsymbol{Z}$. Namely, $\boldsymbol{Z}$ is a compression of $X$ as $\boldsymbol{Z} = f(X)$, and sufficiency satisfies $I(\boldsymbol{Z}; Y) = I(X; Y)$.*

**Theorem 6.2** (Maximum Adjacency Information). *The MI between representations $\boldsymbol{H}_A$ and adjacency $A$ satisfies that $I(A; \boldsymbol{H}_A) \leq I(A; A) = H(A)$. Proof. See Appendix. A.6.*

As such, Theorem 6.2 indicates that the graph representations might maintain the maximum information of private $A$, as $\max I(A; \boldsymbol{H}_A) = H(A)$. Thus, the only sufficient guarantee is not *safe* enough, and the representations $\boldsymbol{H}_A$ potentially stores excess adjacency information $I(A; \boldsymbol{H}_A|Y)$, as illustrated in Fig. 6(a). To reduce, we refer to the minimal sufficient statistics in Proposition 6.3, and deduce the lower bound of adjacency information in Theorem 6.4 as follows.

**Proposition 6.3** (Minimal Sufficient Statistics). *Denote sufficient statistics (Proposition 6.1) of $X$ as $\boldsymbol{Z}$, and the minimal sufficient statistics, $\boldsymbol{Z}^*$, is the optimal graph representation, namely, $\boldsymbol{Z}^* = \arg \min_{\boldsymbol{Z}: I(\boldsymbol{Z};Y)=I(X;Y)} I(\boldsymbol{Z}; X)$.*

**Theorem 6.4** (Minimum Adjacency Information). *For any sufficient graph representations $\boldsymbol{H}_A$ of adjacency $A$ w.r.t. task $Y$, its MI with $A$ satisfies that $I(A; \boldsymbol{H}_A) \geq I(A; Y)$. The minimum information $I(A; \boldsymbol{H}_A) = I(A; Y)$ can be achieved iff $I(A; \boldsymbol{H}_A|Y) = 0$. Proof. See Appendix. A.7.*

Then, the Theorem 6.5 justifies that solving MC-GPB yields an approximation to the optimal representations $\boldsymbol{H}_A^*$, as illustrated in Fig. 6(c), It satisfies sufficiency (accuracy guarantee) and contains minimal adjacency (privacy guarantee).

**Theorem 6.5.** *When degenerating $\beta_c = 0$ and $\beta^i = \beta$, MC-GPB Eq. (4) is equivalent to minimizing the Information Bottleneck Lagrangian, i.e., $\mathcal{L}(p(\boldsymbol{Z}|A)) = H(Y|\boldsymbol{Z}) + \beta I(\boldsymbol{Z}; A)$. It yields a sufficient representation $\boldsymbol{Z}$ of data $A$ for task $Y$, that is an approximation to the optimal representation $\boldsymbol{Z}^*$ in Proposition 6.3. Proof. See Appendix. A.8.*

Table 3: Results of MC-GRA with standard GNNs. Relative promotions (in %) are computed *w.r.t.* results in Tab. 2.

| $X$ | $\boldsymbol{H}_A$ | $\hat{\boldsymbol{Y}}_A$ | $Y$ | Cora | Citeseer | Polblogs | USA | Brazil | AIDS |
|---|---|---|---|---|---|---|---|---|---|
| ✓ | ✓ |  |  | .864 (10.6%↑) | .912 (3.5%↑) | .831 (8.9%↑) | .883 (3.8%↑) | .771 (1.7%↑) | .574 (10.1%↑) |
| ✓ |  | ✓ |  | .839 (7.4%↑) | .902 (2.3%↑) | .836 (8.2%↑) | .913 (10.5%↑) | .800 (9.2%↑) | .567 (8.8%↑) |
| ✓ |  |  | ✓ | .896 (5.5%↑) | .918 (1.2%↑) | .837 (18.7%↑) | .825 (13.3%↑) | .753 (22.8%↑) | .574 (9.9%↑) |
| ✓ | ✓ | ✓ |  | .866 (10.8%↑) | .921 (4.5%↑) | .839 (9.9%↑) | .878 (3.5%↑) | .776 (2.6%↑) | .572 (9.7%↑) |
| ✓ | ✓ |  | ✓ | .905 (6.5%↑) | .930 (2.5%↑) | .832 (6.8%↑) | .878 (3.5%↑) | .758 (2.0%↑) | .603 (15.5%↑) |
| ✓ |  | ✓ | ✓ | .897 (5.6%↑) | .928 (2.3%↑) | .839 (6.8%↑) | .870 (3.3%↑) | .758 (3.7%↑) | .567 (8.6%↑) |
| ✓ | ✓ | ✓ | ✓ | .904 (6.4%↑) | .931 (2.6%↑) | .853 (9.2%↑) | .870 (1.9%↑) | .760 (5.9%↑) | .588 (12.6%↑) |

Table 4: Results of GRA with MC-GPB protected GNNs. Relative reductions are computed *w.r.t.* results in Tab. 1. $I(A; \boldsymbol{H}_A), I(A; \hat{\boldsymbol{Y}}_A)$ are non-learnable GRA (He et al., 2021a) while $I(A; \boldsymbol{H}_{\hat{A}}^1)$ is the learnable GRA (Zhang et al., 2021b).

| MI | Cora | Citeseer | Polblogs | USA | Brazil | AIDS |
|---|---|---|---|---|---|---|
| $I(A; \boldsymbol{H}_A)$ | .706 (7.8%↓) | .750 (1.3%↓) | .724 (5.1%↓) | .716 (15.8%↓) | .745 (1.7%↓) | .564 (3.4%↓) |
| $I(A; \hat{\boldsymbol{Y}}_A)$ | .704 (0.1%↓) | .730 (1.7%↓) | .705 (8.7%↓) | .587 (28.9%↓) | .692 (5.5%↓) | .559 (0.4%↓) |
| $I(A; \boldsymbol{H}_{\hat{A}}^1)$ | .625 (9.9%↓) | .691 (9.8%↓) | .506 (26.3%↓) | .300 (64.5%↓) | .609 (25.1%↓) | .514 (10.6%↓) |
| Acc. | .734 (3.0%↓) | .602 (4.4%↓) | .830 (1.1%↓) | .391 (16.8%↓) | .808 (5.1%↑) | .668 (0.0%↑) |

Table 5: Results of MC-GRA with MC-GPB protected GNNs. Relative reductions are computed *w.r.t.* results in Tab. 3.

| $X$ | $\boldsymbol{H}_A$ | $\hat{\boldsymbol{Y}}_A$ | $Y$ | Cora | Citeseer | Polblogs | USA | Brazil | AIDS |
|---|---|---|---|---|---|---|---|---|---|
| ✓ | ✓ |  |  | .816 (5.5%↓) | .871 (4.4%↓) | .748 (9.9%↓) | .841 (4.7%↓) | .752 (2.4%↓) | .503 (12.3%↓) |
| ✓ |  | ✓ |  | .817 (9.7%↓) | .843 (6.5%↓) | .707 (15.4%↓) | .844 (7.5%↓) | .747 (6.6%↓) | .458 (19.2%↓) |
| ✓ |  |  | ✓ | .892 (0.4%↓) | .888 (3.2%↓) | .699 (16.4%↓) | .738 (10.5%↓) | .700 (7.0%↓) | .490 (14.6%↓) |
| ✓ | ✓ | ✓ |  | .804 (7.1%↓) | .894 (2.9%↓) | .706 (15.8%↓) | .754 (14.1%↓) | .636 (16.7%↓) | .546 (3.7%↓) |
| ✓ | ✓ |  | ✓ | .890 (1.6%↓) | .881 (5.2%↓) | .731 (12.1%↓) | .808 (5.6%↓) | .705 (6.9%↓) | .507 (15.9%↓) |
| ✓ |  | ✓ | ✓ | .858 (4.3%↓) | .903 (2.6%↓) | .791 (5.7%↓) | .768 (11.7%↓) | .656 (13.4%↓) | .511 (9.8%↓) |
| ✓ | ✓ | ✓ | ✓ | .864 (4.4%↓) | .891 (4.2%↓) | .757 (11.2%↓) | .853 (1.9%↓) | .637 (16.1%↓) | .547 (6.9%↓) |

## 7. Empirical Study

In this section, we empirically verify the two proposed methods and provide answers to the three questions. *Q1:* how effective are the proposed methods on real-world datasets with common GNNs? *Q2:* how helpful are MI constraints and injected stochasticity? *Q3:* what insights can empirical results provide to GNNs and defending GRA in practice?

**Setup.** The default target model is a two-layer GCN followed by a linear layer. We also investigate other GNN architectures, including GAT (Veličković et al., 2018) and GraphSAGE (Hamilton et al., 2017). For evaluation, we use the AUC metric as in (Zhang et al., 2021a; Zhu et al., 2021; Zhang et al., 2021b), which considers a set of thresholds. Besides, the implementation software is Pytorch (Paszke et al., 2017) while the hardware is an NVIDIA RTX 3090 GPU. Details of the six datasets are referred to Appendix B.1.

**Baselines.** Two recent works are considered as baselines here: (1) Stealing link (He et al., 2021a) that performs non-learnable GRA on the target model's outputs, which shares a similar spirit as in Sec. 4.2. (2) GraphMI (Zhang et al., 2021b) that conducts learnable GRA with prior knowledge $\mathcal{K} = \{X, Y\}$. The recovered adjacency is obtained by maximizing the classification probability with regard to labels.

### 7.1. Quantitative Results

**Attacking.** As results shown in Tab. 3, the proposed MC-GRA achieves the best results in all six datasets with various settings of prior knowledge sets. The relative promotions in AUC are gained by comparing with the linear ensemble results in Tab. 2. As can be seen, more prior knowledge with a larger $|\mathcal{K}|$ generally bring a higher attack AUC. The learnable MC-GRA brings significantly and consistently better results than the non-learnable methods, especially on the more challenging datasets, *i.e.*, Brazil and AIDS, where at most **22.8**% and **15.5**% promotion can be achieved.

**Defending.** Here, we evaluate the effectiveness of the proposed MC-GPB method in defending against GRA. First, in Tab. 4, we show that MC-GPB is able to defend all the attack methods of GRA. Especially on the Polblogs dataset, MC-GPB achieves a **13.4**% average reduction in privacy leakage at the cost of **1.1**% loss in accuracy. Besides, as in Tab. 5. we show that MC-GPB can also defend the MC-GRA, where it shows consistent and significant reductions in attack AUC. The above results verify the effectiveness of MC-GPB in defending both learnable and non-learnable GRA methods. It can potentially protect GNNs applied in real-world applications, *e.g.*, the recommendation system.

Table 6: MC-GRA with various architectures on Cora.

| $\mathcal{K}$ | GCN | | | GAT | | | GraphSAGE | | |
|---|---|---|---|---|---|---|---|---|---|
| | $L=2$ | $L=4$ | $L=6$ | $L=2$ | $L=4$ | $L=6$ | $L=2$ | $L=4$ | $L=6$ |
| $\{X, Y\}$ | .895 | .892 | .878 | .883 | .878 | .876 | .889 | .872 | .840 |
| $\{X, Y, \boldsymbol{H}_A\}$ | .904 | .900 | .884 | .897 | .885 | .874 | .892 | .8881 | .873 |
| $\{X, Y, \boldsymbol{H}_A, \hat{\boldsymbol{Y}}\}$ | .905 | .895 | .892 | .913 | .887 | .879 | .909 | .893 | .865 |
| Acc. | .792 | .661 | .248 | .637 | .651 | .630 | .614 | .443 | .145 |

Table 7: MC-GPB with various architectures on Polblogs.

| MI | GCN | | | GAT | | | GraphSAGE | | |
|---|---|---|---|---|---|---|---|---|---|
| | $L=2$ | $L=4$ | $L=6$ | $L=2$ | $L=4$ | $L=6$ | $L=2$ | $L=4$ | $L=6$ |
| $I(A; \boldsymbol{H}_A)$ | .724 | .790 | .810 | .901 | .808 | .854 | .805 | .808 | .813 |
| $I(A; \hat{\boldsymbol{Y}}_A)$ | .705 | .650 | .650 | .654 | .623 | .673 | .803 | .668 | .652 |
| $I(A; \boldsymbol{H}_{\hat{A}})$ | .506 | .577 | .532 | .542 | .656 | .536 | .599 | .769 | .468 |
| Acc. | .830 | .822 | .512 | .855 | .880 | .869 | .830 | .869 | .801 |

**Different GNN architectures.** As shown in Tab. 6 as well as Tab. 7, we show that both proposed methods are model-agnostic as they can be generalized to different kinds of GNN with various layers. Generally, a deeper model (larger $L$) can better protect privacy (lower $I(A; \boldsymbol{H}_A^L)$), which is consistent with observations in Sec. 4.3. However, it might come at the cost of severe accuracy degradation due to the well-known over-smoothing effect of GNNs in message propagation. Besides, it is found that *a more powerful model with a higher accuracy is usually more vulnerable to GRA*, which presents a higher risk of privacy leakage in practice.

## 7.2. Ablation Study

**The MI regularization.** As shown in Tab. 8, each MI component contributes to the final results. Specifically, the encoding approximation terms in MC-GRA contribute most to the attack. A potential reason is that hidden representations $\boldsymbol{H}_A$ contain more information about privacy than other variables. And thus, extracting this term brings a higher fidelity in outcomes. In addition, all three kinds of constraints contribute greatly to MC-GPB, while the contributing patterns are diverse. Thus, it is essential to have a careful balance of these three constraints with tuning hyperparameters $\beta_p^i, \beta_c^i$.

**The injected stochasticity.** As can be seen from Tab. 9, learning without injecting stochasticity generally leads to sub-optimal outcomes for both methods. That is, the manual randomness help the removal of spurious correlation for MC-GRA and boosts the forgetting about privacy for MC-GPB. In addition to MI regularization and injected stochasticity, the other ablation study can be found in Appendix. B.2.

## 7.3. Case Visualizations

**The recovered adjacency.** We show the recovered $\hat{\boldsymbol{A}}$ by various GRA methods in Fig. 8. Compared with GraphMI, MC-GRA can recover adjacency more accurately, with fewer wrong predictions and higher AUC values. As for defense, MC-GPB significantly degenerates both GRA methods, with more failure cases and much lower AUC values.

Table 8: Ablation study of two algorithms *w.r.t.* the approximation (*appr.*) and constraint (*cons.*) terms.

| variant | Cora | USA | AIDS |
|---|---|---|---|
| MC-GRA (full) | .905 | .904 | .572 |
| - w/o encoding appr. | .829 (8.3%↓) | .870 (3.7%↓) | .536 (6.2%↓) |
| - w/o decoding appr. | .854 (5.6%↓) | .849 (6.0%↓) | .490 (14.3%↓) |
| - w/o complexity cons. | .889 (1.7%↓) | .858 (5.0%↓) | .537 (11.3%↓) |
| MC-GPB (full) | .745 | .391 | .668 |
| - w/o accuracy cons. | .681 (8.6%↓) | .369 (5.6%↓) | .625 (6.4%↓) |
| - w/o privacy cons. | .707 (5.1%↓) | .249 (36.3%↓) | .480 (28.1%↓) |
| - w/o complexity cons. | .705 (5.4%↓) | .251 (35.8%↓) | .448 (32.9%↓) |

Table 9: Results of removing injecting stochasticity.

| type | case | USA | Brazil | AIDS |
|---|---|---|---|---|
| attack | $\mathcal{K}=\{X, Y\}$ | .802 (2.7%↓) | .713 (5.3%↓) | .567 (1.2%↓) |
| | $\mathcal{K}=\{X, Y, \boldsymbol{H}_A\}$ | .856 (1.3%↓) | .740 (2.3%↓) | .572 (5.1%↓) |
| | $\mathcal{K}=\{X, Y, \boldsymbol{H}_A, \hat{\boldsymbol{Y}}\}$ | .864 (0.4%↓) | .730 (3.9%↓) | .567 (3.5%↓) |
| defense | $I(A; \boldsymbol{H}_A)$ | .861 (16.2%↑) | .758 (1.7%↑) | .564 (0.0%↑) |
| | $I(A; \hat{\boldsymbol{Y}}_A)$ | .309 (47.4%↓) | .722 (4.3%↓) | .548 (2.0%↓) |
| | $I(A; \boldsymbol{H}_{\hat{A}})$ | .389 (29.7%↑) | .796 (30.7%↑) | .539 (4.9%↑) |
| | Acc. | .259 (33.8%↓) | .538 (33.4%↓) | .628 (6.0%↓) |

**A further analysis with the graph information plane.** Besides, we visualize the training procedures of MC-GPB in Fig. 7 based on the graph information plane introduced in Sec. 4.3. As shown, the privacy term is markedly reduced while that in standard training is increased, especially in the later training stage. It shows the trade-off between accuracy and privacy in training GNNs: the accuracy $I(Y; \boldsymbol{Z})$ also starts to decrease when the privacy $I(A; \boldsymbol{Z})$ is minimized to some extent. More visualizations are in Appendix. B.3.

## 8. Further Discussions

**GRA in practice.** In real-world examples regarding the threat model, the prior knowledge set $\mathcal{K}$ can be accessed by an adversary in practice. For instance, to train a GNN model for fraudulent account detection, a social network service provider uses the technology of another company. In this case, the provider will frequently send the company the model's outputs $\boldsymbol{Y}_A$ to debug and improve. Similar circumstances apply to the node embeddings $\boldsymbol{H}_A$, which are typically released. Thus, the inversion of adjacency requiring only a subset $\mathcal{K}$ of the informative variables can be a privacy threat in real-world scenarios of GNNs, which have been widely used in recommendation systems, social networks, citation networks, and drug discovery. Therefore, the user's privacy should be protected especially for personalized relationships and certain sensitive information.

**The inversion target.** Intuitively, the adjacency $A$ and the node features $X$ can be regarded as inversion targets. The key motivations to attack adjacency are the practical risks and understandability to human beings. Unlike visual images that are naturally understandable to humans, the node features are not understandable without the sufficient knowledge of human experts, while the adjacency is much easier to understand. More discussions are in Appendix. D.
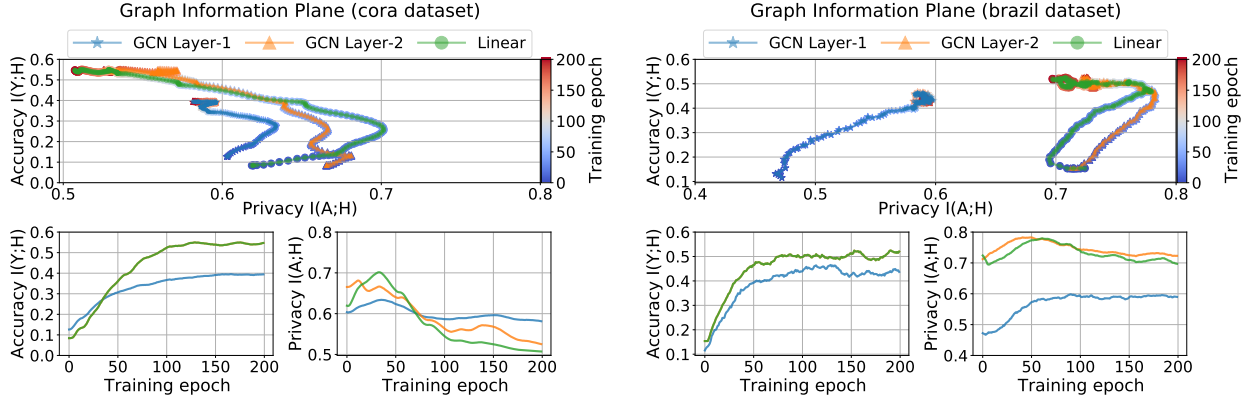
Figure 7: Graph information plane: defensive training with MC-GPB. Compared with the standard training (Fig. 4) without any constraints, MC-GPB effectively decreases the amount of privacy information contained in the graph representations.
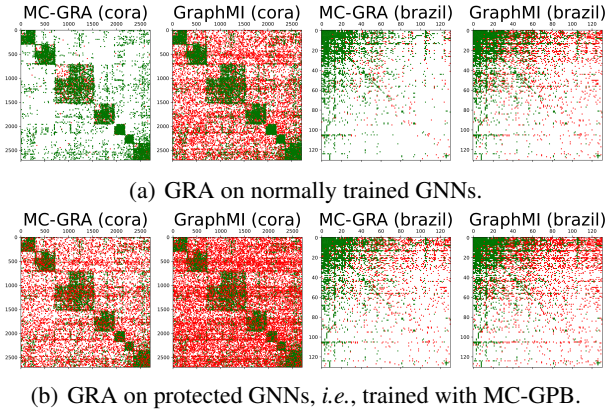


(a) GRA on normally trained GNNs.



(b) GRA on protected GNNs, *i.e.*, trained with MC-GPB.

Figure 8: Examples of recovered adjacency. Green dots are correctly predicted edges while red dots are wrong ones.

**Limitations.** This work follows the common homophily assumption that connected nodes are likely to be in the same category and possess similar features (He et al., 2021a; Zhang et al., 2021b). We leave the generalization to heterogeneous graphs as future work. Besides, our proposed method requires white-box access to the target model. The black-box scenarios with only access to the model's outputs can be more practical but also much more challenging.

**Future directions.** One general direction to enhance GRA is to extract information about adjacency from more information sources, *e.g.*, with partial edges of the target graph or an auxiliary dataset to conduct a transferring attack. GRA can be conducted on more GNN architectures and cooperated with generative models, *e.g.*, the graph auto-encoders or diffusion models. Besides, a more fine-grained study on graph properties is also intriguing, *e.g.*, density, community, number of triangles. To what extent can the above properties be recovered will shed insights into the power of GRA and the memorization effect of GNNs. Besides, applying GRA to more realistic and general settings is also promising, *e.g.*, inductive GNNs which can generalize well to unseen nodes, or the black-box scenarios where the attacker can only get

access to the outputs of the target model. As for defending against GRA in practice, a trained model might be required to completely forget about partial training data with limited budgets for updating their weights. A promising solution here is machine unlearning, especially on large-scale graphs.

**Broader impacts.** *The gun is not guilty, the person who pulled the trigger is*, said the father of the AK-47. It must be admitted that GRA (or any kind of MIA) might be misused to attack real-world targets. For this reason, it is essential to raise the awareness of such an adversary and the potential privacy risk. More importantly, investigating GRA (MIA) enables to understand the black-box deep learning models, to inspire more robust methods, to protect privacy in advance, and to make AI products safer and more trustworthy.

## 9. Conclusion

In this work, we conduct a comprehensive study of enhancing and defending Graph Reconstruction Attack. We conceptually abstract the problem as approximating the original Markov chain by the attack chain. Technically, we derive (1) the chain-based attack method with adaptive designs for extracting more private information; and (2) the chain-based defense method that sharply reduces the attack fidelity with moderate accuracy loss. Empirically, the proposed methods achieve the best results on six datasets and three GNNs.

## Acknowledgements

# References

Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., and Zhang, L. Deep learning with differential privacy. In *CCS*, 2016.

Adamic, L. A. and Glance, N. The political blogosphere and the 2004 us election: divided they blog. In *Proceedings of the 3rd international workshop on Link discovery*, 2005.

Bietti, A., Wei, C.-Y., Dudik, M., Langford, J., and Wu, S. Personalization improves privacy-accuracy tradeoffs in federated learning. In *ICML*, 2022.

Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T., Song, D., Erlingsson, U., et al. Extracting training data from large language models. In *USENIX Security*, 2021.

Chanpuriya, S., Musco, C., Sotiropoulos, K., and Tsourakakis, C. Deepwalking backwards: from embeddings back to graphs. In *ICML*, 2021.

Chen, S., Kahla, M., Jia, R., and Qi, G.-J. Knowledge-enriched distributional model inversion attacks. In *ICCV*, 2021.

Chen, Y., Yang, H., Zhang, Y., Ma, K., Liu, T., Han, B., and Cheng, J. Understanding and improving graph injection attack by promoting unnoticeability. In *ICLR*, 2022.

Cortes, C., Mohri, M., and Rostamizadeh, A. Algorithms for learning kernels based on centered alignment. *The Journal of Machine Learning Research*, 2012.

Cover, T. M. *Elements of information theory*. John Wiley & Sons, 1999.

Dai, H., Li, H., Tian, T., Huang, X., Wang, L., Zhu, J., and Song, L. Adversarial attack on graph structured data. In *ICML*, 2018.

Duddu, V., Boutet, A., and Shejwalkar, V. Quantifying privacy leakage in graph embedding. In *MobiQuitous*, 2020.

Fan, W., Ma, Y., Li, Q., He, Y., Zhao, E., Tang, J., and Yin, D. Graph neural networks for social recommendation. In *TheWebConf*, 2019.

Fang, G., Song, J., Wang, X., Shen, C., Wang, X., and Song, M. Contrastive model inversion for data-free knowledge distillation. In *IJCAI*, 2021.

Fano, R. M. Transmission of information: A statistical theory of communications. *American Journal of Physics*, 1961.

Fredrikson, M., Lantz, E., Jha, S., Lin, S., Page, D., and Ristenpart, T. Privacy in pharmacogenetics: An end-to-end case study of personalized warfarin dosing. In *USENIX Security*, 2014.

Fredrikson, M., Jha, S., and Ristenpart, T. Model inversion attacks that exploit confidence information and basic countermeasures. In *CCS*, 2015.

Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. Neural message passing for quantum chemistry. In *ICML*, 2017.

Grandvalet, Y. and Bengio, Y. Semi-supervised learning by entropy minimization. In *NIPS*, 2004.

Gretton, A., Bousquet, O., Smola, A., and Schölkopf, B. Measuring statistical dependence with hilbert-schmidt norms. In *ALT*, 2005.

Hamilton, W., Ying, Z., and Leskovec, J. Inductive representation learning on large graphs. In *NeurIPS*, 2017.

He, X., Jia, J., Backes, M., Gong, N. Z., and Zhang, Y. Stealing links from graph neural networks. In *USENIX Security*, 2021a.

He, X., Wen, R., Wu, Y., Backes, M., Shen, Y., and Zhang, Y. Node-level membership inference attacks against graph neural networks. *arXiv preprint arXiv:2102.05429*, 2021b.

Hidano, S., Murakami, T., Katsumata, S., Kiyomoto, S., and Hanaoka, G. Model inversion attacks for prediction systems: Without knowledge of non-sensitive attributes. In *PST*, 2017.

Ioannidis, V. N., Zheng, D., and Karypis, G. Few-shot link prediction via graph neural networks for covid-19 drug-repurposing. *arXiv preprint arXiv:2007.10261*, 2020.

Kahla, M., Chen, S., Just, H. A., and Jia, R. Label-only model inversion attacks via boundary repulsion. In *CVPR*, 2022.

Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2016a.

Kipf, T. N. and Welling, M. Variational graph auto-encoders. *NIPS Workshop on Bayesian Deep Learning*, 2016b.

Kool, W., Van Hoof, H., and Welling, M. Stochastic beams and where to find them: The gumbel-top-k trick for sampling sequences without replacement. In *ICML*, 2019.

Kornblith, S., Norouzi, M., Lee, H., and Hinton, G. Similarity of neural network representations revisited. In *ICML*, 2019.

LaValle, S. M., Branicky, M. S., and Lindemann, S. R. On the relationship between classical grid search and probabilistic roadmaps. *The International Journal of Robotics Research*, 2004.

Miao, S., Liu, M., and Li, P. Interpretable and generalizable graph learning via stochastic attention mechanism. In *ICML*, 2022.

Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. Automatic differentiation in PyTorch. In *ICLR*, 2017.

Peng, X., Liu, F., Zhang, J., Lan, L., Ye, J., Liu, T., and Han, B. Bilateral dependency optimization: Defending against model-inversion attacks. In *SIGKDD*, 2022.

Ribeiro, L. F., Saverese, P. H., and Figueiredo, D. R. struc2vec: Learning node representations from structural identity. In *SIGKDD*, 2017.

Riesen, K., Bunke, H., et al. Iam graph database repository for graph based pattern recognition and machine learning. In *SSPR/SPR*, 2008.

Rong, Y., Huang, W., Xu, T., and Huang, J. Dropedge: Towards deep graph convolutional networks on node classification. In *ICLR*, 2020.

Sen, P., Namata, G., Bilgic, M., Getoor, L., Galligher, B., and Eliassi-Rad, T. Collective classification in network data. *AI magazine*, 2008.

Shamir, O., Sabato, S., and Tishby, N. Learning and generalization with the information bottleneck. *Theoretical Computer Science*, 2010.

Shen, Y., He, X., Han, Y., and Zhang, Y. Model stealing attacks against inductive graph neural networks. In *IEEE SP*, 2022.

Shwartz-Ziv, R. and Tishby, N. Opening the black box of deep neural networks via information. *arXiv preprint arXiv:1703.00810*, 2017.

Struppek, L., Hintersdorf, D., Correia, A. D. A., Adler, A., and Kersting, K. Plug and play attacks: Towards robust and flexible model inversion attacks. In *ICML*, 2022.

Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.

Tishby, N. and Zaslavsky, N. Deep learning and the information bottleneck principle. In *IEEE information theory workshop*, 2015.

Tishby, N., Pereira, F. C., and Bialek, W. The information bottleneck method. *arXiv preprint physics/0004057*, 2000.

Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., and Bengio, Y. Graph attention networks. In *ICLR*, 2018.

Wang, K.-C., Fu, Y., Li, K., Khisti, A., Zemel, R., and Makhzani, A. Variational model inversion attacks. In *NeurIPS*, 2021.

Wu, S., Sun, F., Zhang, W., Xie, X., and Cui, B. Graph neural networks in recommender systems: a survey. *ACM Computing Surveys*, 2020a.

Wu, T., Ren, H., Li, P., and Leskovec, J. Graph information bottleneck. In *NeurIPS*, 2020b.

Xie, S. and Ermon, S. Reparameterizable subset sampling via continuous relaxations. In *IJCAI*, 2019.

Yang, Y.-Y., Chou, C.-N., and Chaudhuri, K. Understanding rare spurious correlations in neural network. *arXiv preprint arXiv:2202.05189*, 2022.

Yang, Z., Chang, E.-C., and Liang, Z. Adversarial neural network inversion via auxiliary knowledge alignment. *arXiv preprint arXiv:1902.08552*, 2019.

You, Y., Chen, T., Sui, Y., Chen, T., Wang, Z., and Shen, Y. Graph contrastive learning with augmentations. In *NeurIPS*, 2020.

You, Y., Chen, T., Wang, Z., and Shen, Y. Bringing your own view: Graph contrastive learning without prefabricated data augmentations. In *WSDM*, 2022.

Zeng, H., Zhou, H., Srivastava, A., Kannan, R., and Prasanna, V. Graphsaint: Graph sampling based inductive learning method. In *ICLR*, 2019.

Zhang, M. and Chen, Y. Link prediction based on graph neural networks. In *NeurIPS*, 2018.

Zhang, M., Li, P., Xia, Y., Wang, K., and Jin, L. Labeling trick: A theory of using graph neural networks for multi-node representation learning. In *NeurIPS*, 2021a.

Zhang, R., Hidano, S., and Koushanfar, F. Text revealer: Private text reconstruction via model inversion attacks against transformers. *arXiv preprint arXiv:2209.10505*, 2022a.

Zhang, Y., Jia, R., Pei, H., Wang, W., Li, B., and Song, D. The secret revealer: Generative model-inversion attacks against deep neural networks. In *CVPR*, 2020.

Zhang, Z., Liu, Q., Huang, Z., Wang, H., Lu, C., Liu, C., and Chen, E. Graphmi: Extracting private graph data from graph neural networks. In *IJCAI*, 2021b.

Zhang, Z., Chen, M., Backes, M., Shen, Y., and Zhang, Y. Inference attacks against graph neural networks. In *USENIX Security*, 2022b.

Zhao, T., Liu, G., Wang, D., Yu, W., and Jiang, M. Learning from counterfactual links for link prediction. In *ICML*, 2022.

Zhao, X., Zhang, W., Xiao, X., and Lim, B. Exploiting explanations for model inversion attacks. In *ICCV*, 2021.

Zhu, J., Yao, J., Liu, T., Yao, Q., Xu, J., and Han, B. Combating exacerbated heterogeneity for robust models in federated learning. In *ICLR*, 2023.

Zhu, Z., Zhang, Z., Xhonneux, L., and Tang, J. Neural bellman-ford networks: A general graph neural network framework for link prediction. In *NeurIPS*, 2021.

# Appendix

# A. Theoretical justification

## A.1. Notations

With adjacent matrix $A$ and node features $X$, an undirected graph is denoted as $\mathcal{G} = (A, X)$, where $A_{ij} = 1$ means there is an edge $e_{ij}$ between $v_i$ and $v_j$. For each node $v_i$, its $D$-dimension node feature is denoted as $X_{[i,:]} \in \mathbb{R}^D$, and its label $y_i \in Y = \{y_i\}_{i=1}^N$ indicates the node class. The node classification task is to predict the label $Y = \{y_i\}_{i=1}^N$ of each node via a parameterized model $f_\theta(\cdot)$, *i.e.*, $f_\theta(A, X) = \hat{Y} \leftrightarrow Y$. We summarize the frequently used notations in Table 10 as follows.

Table 10: The most frequently used notations in this work.

| notations | meanings |
|---|---|
| $\mathcal{V} = \{v_i\}_{i=1}^N$ | the set of nodes |
| $\mathcal{E} = \{e_{ij}\}_{ij=1}^M$ | the set of edges |
| $A \in \{0,1\}^{N \times N}$ | the adjacent matrix with binary elements |
| $X \in \mathbb{R}^{N \times D}$ | the node features |
| $\mathcal{G} = (A, X)$ | the input graph of a GNN |
| $Y$ | the labels of nodes |
| $\boldsymbol{H}_A$ | representation of all nodes with adjacency $A$ |
| $H(X)$ | the information entropy of random variable $X$ |
| $H(X, Y)$ | the joint entropy of variable $X$ and $Y$ |
| $I(X; Y)$ | the mutual information of $X$ and $Y$ |
| $I(X; Y \mid Z)$ | the conditional mutual information of $X$ and $Y$ when observing $Z$ |

## A.2. Preliminaries for information measures

**Definition A.1** (Informational Divergence). The informational divergence (also called relative entropy or Kullback-Leibler distance) between two probability distributions $p$ and $q$ on a finite space $\mathcal{X}$ (*i.e.*, a common alphabet) is defined as

$$D(p\|q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} = \mathbb{E}_p \left[ \log \frac{p(X)}{q(X)} \right] \tag{5}$$

*Remark* A.2. $D(p\|q)$ measures the *distance* between $p$ and $q$. However, $D(\cdot\|\cdot)$ is not a true metric, and it does not satisfy the triangular inequality. $D(p\|q)$ is non-negative and $D(p\|q) = 0$ if and only if $p = q$.

**Definition A.3** (Mutual Information). Given two discrete random variables $X$ and $Y$, the mutual information (MI) $I(X; Y)$ is the relative entropy between the joint distribution $p(x, y)$ and the product of the marginal distributions $p(x)p(y)$, namely,

$$\begin{aligned}
I(X; Y) &= D(p(x, y) \| p(x)p(y)) \\
&= \sum_{x \in X, y \in Y} p(x, y) \log \left( \frac{p(x, y)}{p(x)p(y)} \right) \\
&= \sum_{x \in X, y \in Y} p(x, y) \log \left( \frac{p(x|y)}{p(x)} \right).
\end{aligned} \tag{6}$$

*Remark* A.4. $I(X; Y)$ is symmetrical in $X$ and $Y$, *i.e.*, $I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) = I(Y; X)$.

**Proposition A.5** (Chain Rule for Entropy). $H(X_1, X_2, \cdots, X_n) = \sum_{i=1}^n H(X_i | X_1, X_2, \cdots, X_{i-1})$.

**Proposition A.6** (Chain Rule for Conditional Entropy). $H(X_1, X_2, \cdots, X_n | Y) = \sum_{i=1}^n H(X_i | X_1, X_2, \cdots, X_{i-1}, Y)$.

**Proposition A.7** (Chain Rule for Mutual Information). $I(X_1, X_2, \cdots, X_n; Y) = \sum_{i=1}^n I(X_i; Y | X_1, X_2, \cdots, X_{i-1})$.

**Corollary A.8.** $\forall A, Z_i, Z_j, I(A; Z_i, Z_j) \geq \max\left(I(A; Z_i), I(A; Z_j)\right)$.

*Proof.* As $I(A; Z_i|Z_j) \geq 0$, $I(A; Z_i, Z_j) = I(A; Z_i) + I(A; Z_i|Z_j) \geq I(A; Z_i)$. Similarly, $I(A; Z_i, Z_j) \geq I(A; Z_j)$ can be obtained. Thus, we have $I(A; Z_i, Z_j) \geq \max\big(I(A; Z_i), I(A; Z_j)\big)$. □

**Proposition A.9** (Chain Rule for Conditional Mutual Information). $I(X_1, \cdots, X_n; Y|Z) = \sum_{i=1}^n I(X_i; Y|X_1, \cdots, X_{i-1}, Z)$.

**Definition A.10** (Markov). A discrete stochastic process is called Markov if it satisfies $p(x_{i+1}|x_i, x_{i-1}, x_{i-2}, \cdots, x_1) = p(x_{i+1}|x_i)$ $\forall i$. As such, $\forall n > 1$, $p(x_1, x_2, \ldots, x_n) = p(x_1)p(x_2|x_1)p(x_3|x_2)\ldots p(x_n|x_{n-1})$.

**Definition A.11** (Causally Similarity). Suppose we have two stochastic processes $X(t)$ and $Y(t)$ defined on the ordered set $R$ with associated probability functions $p$ and $q$ and the same outcome sets $\{\overrightarrow{x}(t)\}$. We say that the two processes are *causally similar* if $p(\overrightarrow{x}(t)|\overrightarrow{x}(t-a)) = q(\overrightarrow{x}(t)|\overrightarrow{x}(t-a))$ $\forall t$ and $\forall a > 0$.

*Remark A.12. Two stochastic processes are causally similar if they are time homogenous, Markov, and share the same transition matrix. Besides, two Markov processes are also causally similar that is not necessarily time homogenous if they share the same transition matrix at the same time step.*

**Lemma A.13.** *Suppose we have two causally similar stochastic processes with probability functions at time t of $p(\overrightarrow{x_t})$ and $q(\overrightarrow{x_t})$. Then $D(p(\overrightarrow{x_t})||q(\overrightarrow{x_t})) \leq D(p(\overrightarrow{x_s})||q(\overrightarrow{x_s}))$ when $t > s$.*

**Lemma A.14.** *In a stationary Markov process, the entropy conditioned on the initial condition is non-decreasing.*

*Proof.* That is, $H(X_n|X_1) \geq H(X_n|X_1, X_2)$ as further conditioning reduces entropy. Besides, $H(X_n|X_1, X_2) = H(X_n|X_2) = H(X_{n-1}|X_1)$. Thus, $H(X_n|X_1) \geq H(X_{n-1}|X_1)$, which shows that $H(X_n|X_1)$ is non-decreasing. □

## A.3. Proof for Theorem 5.3

**Lemma A.15** (Invertible Transformations Are Invariant to MI). *The mutual information is invariant to any invertible transformations $\psi(\cdot), \phi(\cdot)$, namely, $I(X; Y) = I(\psi(X); Y) = I(X; \phi(Y)) = I(\psi(X); \phi(Y))$.*

**Lemma A.16** (Non-invertible Transformation Reduces MI). *For any non-invertible transformation $\psi(\cdot)$, it reduce the MI between $X$ and $Y$ as $I(X; Y) \geq I(\psi(X); Y) \geq I(\psi(X); \psi(Y))$.*

*Proof.* Based on Lemma A.15 and Lemma A.16, we have the following deduction. As $\boldsymbol{Z}_A^{i+1} = \sigma(\psi(A) \cdot \boldsymbol{Z}_A^i \cdot \boldsymbol{\theta}^i)$, where graph convolution kernel $\psi(A) \in \mathbb{R}^{N \times N}$ and weights $\boldsymbol{\theta}^i \in \mathbb{R}^{D \times D}$ are invertible transformations. Note the activate function $\sigma(\cdot)$ (*e.g.*, ReLU) is a non-invertible transformation that $H(X) \geq H(\sigma(X))$, and $I(X; Y) \geq I(\sigma(X); \sigma(Y))$.

$$
\begin{aligned}
I(\boldsymbol{Z}_A^{i+1} \; ; \; \boldsymbol{Z}_{\hat{\boldsymbol{A}}}^{i+1}) &= I\big(\sigma(\psi(A) \cdot \boldsymbol{Z}_A^i \cdot \boldsymbol{\theta}^i) \; ; \; \sigma(\psi(\hat{\boldsymbol{A}}) \cdot \boldsymbol{Z}_{\hat{\boldsymbol{A}}}^i \cdot \boldsymbol{\theta}^i)\big) \\
&\leq I\big(\psi(A) \cdot \boldsymbol{Z}_A^i \cdot \boldsymbol{\theta}^i \; ; \; \psi(\hat{\boldsymbol{A}}) \cdot \boldsymbol{Z}_{\hat{\boldsymbol{A}}}^i \cdot \boldsymbol{\theta}^i\big) \\
&= I\big(\psi(A) \cdot \boldsymbol{Z}_A^i \; ; \; \psi(\hat{\boldsymbol{A}}) \cdot \boldsymbol{Z}_{\hat{\boldsymbol{A}}}^i\big) = I\big(\boldsymbol{Z}_A^i \cdot \boldsymbol{\theta}^i \; ; \; \boldsymbol{Z}_{\hat{\boldsymbol{A}}}^i \cdot \boldsymbol{\theta}^i\big) \\
&= I(\boldsymbol{Z}_A^i \; ; \; \boldsymbol{Z}_{\hat{\boldsymbol{A}}}^i).
\end{aligned}
\tag{7}
$$

Thus, $\forall i \in [L-1], I(\boldsymbol{Z}_A^{i+1}; \boldsymbol{Z}_{\hat{\boldsymbol{A}}}^{i+1}) \leq I(\boldsymbol{Z}_A^i; \boldsymbol{Z}_{\hat{\boldsymbol{A}}}^i)$. The layer-wise MI of the two Markov chains is decreasing by layers. □

## A.4. Proof for Theorem 5.4

**Lemma A.17** (Finite sample bounds). *Assume all variables' empirical estimates of the mutual information are based on finite support, i.e., $K = |\hat{X}| \approx 2^{I(\hat{X}; X)}$. Then, we denote by $\hat{I}(\cdot; \cdot)$ the finite sample distribution $\hat{p}(x, y)$ for a given sample of size $n$. The generalization bounds in (Shamir et al., 2010) guarantee that*

$$
\begin{aligned}
I(\hat{X}; Y) &\leq \hat{I}(\hat{X}; Y) + O(\frac{K|\mathcal{Y}|}{\sqrt{n}}), \\
I(\hat{X}; X) &\leq \hat{I}(\hat{X}; X) + O(\frac{K}{\sqrt{n}}).
\end{aligned}
\tag{8}
$$

*Proof.* Let the random variables $X$ and $Y$ represent input and output messages with a joint probability $P(x, y)$. Let $e$ represent an occurrence of error, *i.e.*, that $X \neq \hat{X}$ with $\hat{X} = f(Y)$ being an approximate version of $X$. Fano's

inequality (Fano, 1961) (also known as the Fano converse) states that the conditional entropy

$$H(X|Y) = -\sum_{i,j} P(x_i, y_j) \log P(x_i|y_j)$$
$$\leq H_b(e) + P(e) \log(|\mathcal{X}| + 1),$$

(9)

where the probability of the communication error $P(e) \triangleq P(X \neq \hat{X}) \geq \frac{H(X|Y)-1}{\log(|\mathcal{X}|)}$, and $H_b(e)$ is the corresponding binary entropy that computed as $H_b(e) = -e \log_2(e) - (1-e) \log_2(1-e)$.

Note that the data processing inequality (Cover, 1999) indicates that any three variables $X, Y, Z$ that form a Markov chain $X \to Y \to Z$, satisfy $I(X;Y) \geq I(X;Z)$ and $I(Y;Z) \geq I(X;Z)$. As $(A, X) \xrightarrow{\theta^1} \boldsymbol{Z}_A^1$ and $(\hat{A}, X) \xrightarrow{\theta^1} \boldsymbol{Z}_{\hat{A}}^1$, we have

$$I(A; \hat{A}) \geq I(\boldsymbol{Z}_A^1 ; \boldsymbol{Z}_{\hat{A}}^1) \geq I(\boldsymbol{Z}_A^2 ; \boldsymbol{Z}_{\hat{A}}^2) \geq \cdots \geq I(\boldsymbol{Z}_A^L ; \boldsymbol{Z}_{\hat{A}}^L) = I(\boldsymbol{H}_A ; \boldsymbol{H}_{\hat{A}})$$

(10)

Then, according to Fano's inequality (Fano, 1961), the lower bound of MI $I(\boldsymbol{H}_A; \boldsymbol{H}_{\hat{A}})$ is

$$I(\boldsymbol{H}_A ; \boldsymbol{H}_{\hat{A}}) = H(\boldsymbol{H}_A) - H(\boldsymbol{H}_A \mid \boldsymbol{H}_{\hat{A}})$$
$$\geq H(\boldsymbol{H}_A) - H_b(e) - P(e) \log(|\mathcal{H}|),$$

(11)

where entropy $H(\boldsymbol{H}_A) = \mathbb{E}_{x \in \boldsymbol{H}_A}[-\log p(x)] = -\sum_{x \in \boldsymbol{H}_A} p(x) \log p(x)$, the probability of approximation error $P(e) = P(\boldsymbol{H}_A \neq \boldsymbol{H}_{\hat{A}})$, and the binary entropy $H_b(e) = -e \log_2(e) - (1-e) \log_2(1-e)$. $\mathcal{H}$ denotes the support of $\boldsymbol{H}_A$. Specifically, the approximation fidelity $I(A; \hat{A}) \geq -\sum_{x \in \boldsymbol{H}_A} p(x) \log p(x) + e \log_2(e) + (1-e) \log_2(1-e) - P(e) \log(|\mathcal{H}|)$. $\square$

### A.5. Proof for Theorem 5.5

*Proof.* To learn $\hat{A}$ given the prior knowledge $\mathcal{K} = \{X, Y, \boldsymbol{H}_A\}$, we have $H(\hat{A}) \leq H(\mathcal{K})$, and $\forall Z, I(Z; \hat{A}) \leq I(Z; \mathcal{K}) = I(Z; X, Y, \boldsymbol{H}_A)$. Thus, the recovering fidelity of $\hat{A}$ satisfies $I(A; X, Y, \boldsymbol{H}_A) - I(A; \hat{A}) \geq 0$. Then, we obtain the upper bound of the attack fidelity with the optimal recover adjacency $\hat{A}^*$, namely,

$$\hat{A}^* = \max_{\hat{A}} I(A; \hat{A}) = I(A; \mathcal{K}) = I(A; X, Y, \boldsymbol{H}_A),$$
$$s.t. \ I(A; \mathcal{K}|\hat{A}^*) = I(A; \hat{A}^*|\mathcal{K}) = 0.$$

(12)

Solving MC-GRA (Eq. (3)) that $\exists \alpha_1, \alpha_2 \in \mathbb{R}$, $\hat{A}^* = \arg\max_{\hat{A}} \sum_{i=1}^{L} \alpha_1 I(\boldsymbol{H}_A; \boldsymbol{H}_{\hat{A}}^i) + \alpha_2 I(Y; \boldsymbol{Y}_{\hat{A}})$ yields a sufficient solution to achieve the optimal fidelity, *i.e.*, $\hat{A}^* : I(A; \hat{A}^*) = I(A; X, Y, \boldsymbol{H}_A)$. However, the optimal $\hat{A}^*$ does not necessarily mean exactly recover the original $A$, as $H(A|\hat{A}^*) = H(A) - I(A; \hat{A}^*) \geq 0$. Intuitively, the perfect recovery can not be achieved due to the data compression nature of the learning process. Besides, $H(A) \geq \max_{Z \in \mathcal{K}} H(Z)$ is a usual case as the hidden dimension $D \ll N$. The remaining information, *i.e.*, $H(A|\hat{A}^*) = H(A|\mathcal{K})$, that is unobservable from $\mathcal{K} = \{X, Y, \boldsymbol{H}_A\}$, can not be recovered unless additional information is provided.

$\square$

### A.6. Proof for Theorem 6.2

*Proof.* $\forall X, Y$, we have $I(X; Y) \leq I(X; X) = H(X)$. Thus, the MI between representations $\boldsymbol{H}_A$ and adjacency $A$ satisfies that $I(A; \boldsymbol{H}_A) \leq I(A; A) = H(A)$. Which means, the upper bound of the MI, *i.e.*, the worst privacy leakage as $I(A; \boldsymbol{H}_A) \leq H(A)$, is that all the private information $H(A)$ about the adjacency is obtained for the attacker. $\square$

### A.7. Proof for Theorem 6.4

*Proof.* For any sufficient graph representations $\boldsymbol{H}_A$ of adjacency $A$ w.r.t. task $Y$ introduced in Proposition 6.1, its MI with $A$ satisfies that $I(A; Y) = I(\boldsymbol{H}_A; Y)$, as $\boldsymbol{H}_A$ can be seen as extracted from $A$. However, $I(A; \boldsymbol{H}_A) \geq I(A; Y)$ as the data processing inequality (Cover, 1999) in Markov chain $A \to \boldsymbol{H}_A \to Y$. Based on the two above conditions, the minimum information $I(A; \boldsymbol{H}_A) = I(A; Y)$ can be achieved if and only if $I(A; \boldsymbol{H}_A|Y) = 0$. That is, the optimal representations $\boldsymbol{H}_A^*$ satisfy (1) sufficient condition that $I(A; Y) = I(\boldsymbol{H}_A^*; Y)$, and (2) minimal condition that $I(A; \boldsymbol{H}_A^*) = I(A; Y)$.

Thus, $I(A; \boldsymbol{H}_A^*|Y) = I(A; \boldsymbol{H}_A^*, Y) - I(A, Y) = I(A, Y) - I(A, Y) = 0$. $\square$

16

### A.8. Proof for Theorem 6.5

*Proof.* When degenerate $\beta_c = 0$ and $\beta^i = \beta$, MC-GPB is equivalent to minimizing the Information Bottleneck Lagrangian (Shwartz-Ziv & Tishby, 2017), *i.e.*, $\mathcal{L}(p(\boldsymbol{Z}|A)) = H(Y|\boldsymbol{Z}) + \beta I(\boldsymbol{Z}; A)$, in the limit $\beta \to 0$. Specifically, $\mathcal{L}(p(\boldsymbol{Z}|A)) = H(Y|\boldsymbol{Z}) + \beta I(\boldsymbol{Z}; A) = H(Y) - I(\boldsymbol{Z}; Y) + \beta I(\boldsymbol{Z}; A) \propto -I(\boldsymbol{Z}; Y) + \beta I(\boldsymbol{Z}; A)$, where entropy $H(Y)$ is a constant. Then, we deduce the optimal case of $\min \mathcal{L}(p(\boldsymbol{Z}|A)) = \max I(\boldsymbol{Z}; Y) - \beta I(\boldsymbol{Z}; A)$ as follows.

$$
\begin{aligned}
&\max I(\boldsymbol{Z}; Y) - \beta I(\boldsymbol{Z}; A) \\
=& \max \big(I(Y; \boldsymbol{Z}, A) - I(A; Y|\boldsymbol{Z})\big) - \beta\big(I(\boldsymbol{Z}; A, Y) - I(A; Y|\boldsymbol{Z})\big) \\
=& \max I(Y; \boldsymbol{Z}, A) - (1 - \beta)I(A; Y|\boldsymbol{Z}) - \beta I(\boldsymbol{Z}; A, Y) \\
=& \max I(Y; A) - (1 - \beta)I(A; Y|\boldsymbol{Z}) - \beta I(\boldsymbol{Z}; A, Y) \\
=& \max(1 - \beta)I(A; Y) - (1 - \beta)I(A; Y|\boldsymbol{Z}) - \beta I(\boldsymbol{Z}; A|Y) \\
=&(1 - \beta)I(A; Y).
\end{aligned}
\tag{13}
$$

As the two MI terms $I(A; Y|\boldsymbol{Z}) \geq 0$ and $I(\boldsymbol{Z}; A|Y) \geq 0$, the optimal $\boldsymbol{Z}^*$ should satisfies that $I(A; Y|\boldsymbol{Z}^*) = I(\boldsymbol{Z}^*; A|Y) = 0$. As such, it yields a sufficient representation $\boldsymbol{Z}$ of data $A$ for task $Y$, that is an approximation to the minimal and sufficient representations $\boldsymbol{Z}^*$ in Proposition 6.3, *i.e.*, $\boldsymbol{Z}^* = \arg\min_{\boldsymbol{Z}:\ I(\boldsymbol{Z};Y)=I(A;Y)} I(\boldsymbol{Z}; A)$. □

## B. Full empirical study

### B.1. Datasets

Six common datasets are utilized in experiments, which are collected from four diverse domains: (1) Cora and Citeseer (Sen et al., 2008) are citation networks where nodes are documents, and edges indicate citations among them; (2) Polblogs (Adamic & Glance, 2005) is a social network of political blogs where nodes represent blogs with political leaning while edges are citations; (3) USA and Brazil (Ribeiro et al., 2017) are air-traffic networks where nodes are airports and edges denote airlines; (4) AIDS (Riesen et al., 2008) is a chemical network where each node is an atom, and each edge is a chemical bond. The data statistics are in Tab. 11.

Table 11: Dataset statistics. The hard homophily of an edge $e_{ij}$ is computed as $\mathbb{I}(y_i, y_j)$ with node labels, while the soft homophily is calculated by $\cos(x_i, x_j)$ with node features. "—" means this dataset is intrinsic without node features.

| dataset | Cora | Citeseer | Polblogs | USA | Brazil | AIDS |
|---|---|---|---|---|---|---|
| # Nodes | 2,708 | 3,327 | 1490 | 1190 | 131 | 1429 |
| # Edges | 5,278 | 4,676 | 33430 | 27164 | 2077 | 2948 |
| # Class | 7 | 6 | 2 | 4 | 4 | 14 |
| # Features | 1433 | 3703 | — | — | — | 4 |
| Soft homophily | 0.83 | 0.81 | — | — | — | 0.06 |
| Hard homophily | 0.81 | 0.74 | 0.91 | 0.70 | 0.45 | 0.51 |

### B.2. Full quantitative results

**A further comparison of attack methods.** For the attack, we further compare the proposed method (MC-GRA) with three baselines in the below table. Here, the evaluation is also with the AUC metric, where a higher value indicates a better attack performance. The **boldface** numbers represent the best results. As can be seen from the table below, the MC-GRA consistently achieves the best in all six datasets, outperforming all the baselines by a large margin.

Table 12: A further quantitative comparison of attack methods (with AUC metric).

| dataset | Cora | Citeseer | Polblogs | USA | Brazil | AIDS |
|---|---|---|---|---|---|---|
| single MI (Tab. 1) | .815 | .881 | .763 | .850 | .758 | .584 |
| ensemble (Tab. 2) | .849 | .907 | .781 | .852 | .717 | .522 |
| GraphMI | .812 | .781 | .791 | .769 | .680 | .575 |
| MC-GRA (Tab. 3) | **.904** | **.931** | **.853** | **.870** | **.760** | **.588** |

**A further comparison of defense methods.** For the defense, we compare the proposed MC-GPB with two additional defense methods, *i.e.*, adding random noise and differentiable privacy. Specifically, we inject Gaussian noise into the model prediction, termed random noise. While another baseline, termed differential privacy (Abadi et al., 2016), is achieved by adding Gaussian noise to the clipped gradients in each training iteration. The empirical results are shown in the below table, where GraphMI (Zhang et al., 2021b) is used as the attack method. As can be seen, the defending power of random noise and differential privacy comes at the price of sharply degenerating the model's accuracy. By contrast, our proposed MC-GPB significantly degenerates GRA with much lower AUC while maintaining high accuracy simultaneously.

Table 13: A further quantitative comparison of defense methods.

| dataset | Cora | | Citeseer | | Polblogs | | USA | | Brazil | | AIDS | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ACC↑ | AUC↓ | ACC↑ | AUC↓ | ACC↑ | AUC↓ | ACC↑ | AUC↓ | ACC↑ | AUC↓ | ACC↑ | AUC↓ |
| No defense | .757 | .812 | .630 | .781 | .833 | .791 | .470 | .769 | .769 | .680 | .668 | .575 |
| Random noise | .620 | .657 | .570 | .727 | .802 | .759 | .440 | .754 | .634 | .713 | .572 | .559 |
| Differential privacy | .315 | .500 | .224 | .500 | .553 | .502 | .263 | .500 | .423 | .706 | .131 | .502 |
| MC-GPB | .734 | .625 | .602 | .691 | .830 | .506 | .391 | .300 | .808 | .609 | .668 | .514 |

**Ablation study of similarity measurement.** We also conduct experiments with the influence of similarity measurement since our implementation depends on the estimation of mutual information, shown as Tab. 14. As can be seen, the MC-GRA has consistent performance across different similarity measurements, while the MC-GPB exhibits a high variance for different similarity measurements. Therein, the HSIC and CKA are generally good choices.

Table 14: Ablation study of similarity measurements (with AUC metric).

| type | case | Cora | | | | USA | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | DP | HSIC | CKA | KDE | DP | HSIC | CKA | KDE |
| attack | $\mathcal{K}=\{X,Y\}$ | .876 | .871 | .873 | .876 | .791 | .800 | .802 | .802 |
| | $\mathcal{K}=\{X,Y,\boldsymbol{H}_A\}$ | .892 | .890 | .892 | .895 | .856 | .850 | .845 | .851 |
| | $\mathcal{K}=\{X,Y,\boldsymbol{H}_A,\hat{\boldsymbol{Y}}\}$ | .898 | .898 | .904 | .896 | .846 | .852 | .818 | .840 |
| defense | $I(A;\boldsymbol{H}_A)$ | .476 | .751 | .701 | .706 | .716 | .873 | .879 | .883 |
| | $I(A;\hat{\boldsymbol{Y}}_A)$ | .508 | .688 | .705 | .704 | .587 | .542 | .872 | .873 |
| | $I(A;\boldsymbol{H}_{\hat{A}})$ | .505 | .644 | .644 | .625 | .300 | .467 | .770 | .728 |
| | Acc. | .306 | .635 | .758 | .734 | .391 | .319 | .431 | .447 |

**Ablation study of parameterization methods.** We also provide the empirical result of different parameterization methods of MC-GRA, which are mentioned in Sec. 5. Overall, the GNNs method achieves the best score out of the three methods, especially for the dataset with given node features (Cora, Citeseer, AIDS), and exceeds its counterparts by a large margin. Therein, the Gaussian parameterization generally performs better than its counterparts for graphs without node features.

Table 15: Attack with different parameterization methods (with AUC metric).

| variant | Cora | Citeseer | Polblogs | USA | Brazil | AIDS |
|---|---|---|---|---|---|---|
| Direct Matrix | .890 | .580 | .684 | .737 | .521 | .540 |
| Gaussian | .893 | .654 | .777 | .846 | .758 | .567 |
| GNNs | .891 | .889 | .803 | .776 | .731 | .662 |

**A further empirical study about the evaluation metric.** Without requiring any estimation, we utilize the AUC (area under the curve) metric with the ground-truth edges in $A$ to quantify the mutual information $I(A;\boldsymbol{Z})$. And alternatively, the metric can be replaced by AP (Average Precision), MRR (Mean Reciprocal Rank), Hit@K (the ratio of positive edges that are ranked at the K-th place or above), which are also common in the link prediction task.

The metrics mentioned above treat all edges equally indeed. An objective measurement is required to further discriminate

between high-value and low-value links. Here, the link homophily is a proper measurement. For instance, the link (Jaime, Tyrion) in Figure. 1 can be seen as a *homogeneous* link because Jaime and Tyrion have the same node labels (*i.e.*, Lannister). On the other hand, the link (Daenerys, Jon) can be seen as a *heterogeneous* link because Daenerys and Jon do not have the same node labels (as audiences in the earlier period, we do not know they are Targaryens, and they have a kinship). Formally, the homogeneous links can be denoted as $\{e_{ij} : y_i = y_j\}$ while the heterogeneous links are $\{e_{ij} : y_i \neq y_j\}$, where $y_i, y_j$ are node labels of node $i$, node $j$. In what follows, we further investigate the effectiveness of GRA on these two kinds of links.

**(1)** For the attack, the homogeneous links are much easier to recover, as shown in the table below. More importantly, it is observed that the high-value heterogeneous links are naturally protected but can still be recovered to some extent.

Table 16: A further quantitative comparison of attack methods on homogeneous or heterogeneous links (with AUC metric).

| dataset | Cora | Citeseer | Polblogs | USA | Brazil | AIDS |
|---|---|---|---|---|---|---|
| MC-GRA (Homogeneous links) | .960 | .917 | .896 | .951 | .891 | .585 |
| MC-GRA (Heterogeneous links) | .684 | .861 | .298 | .716 | .564 | .551 |
| GraphMI (Homogeneous links) | .724 | .799 | .717 | .919 | .871 | .707 |
| GraphMI (Heterogeneous links) | .569 | .675 | .391 | .666 | .728 | .437 |

**(2)** For defense, we apply the proposed MC-GPB to protect GCN against the GRA by GraphMI. In addition, we also implement a revised version, *i.e.*, MC-GPB-hetero, which only focuses on protecting the heterogeneous links of the original adjacency matrix. As results are shown in the table below, the recovery of heterogeneous links is significantly degenerated by MC-GPB and further degenerated by MC-GPB-hetero. Thus, we justify that MC-GPB and its revised version are capable of protecting the high-value heterogeneous links.

Table 17: A further quantitative comparison of attack methods on homogeneous or heterogeneous links (with AUC metric).

| dataset | Cora | Citeseer | Polblogs | USA | Brazil | AIDS |
|---|---|---|---|---|---|---|
| No defense (Heterogeneous links) | .569 | .675 | .391 | .666 | .728 | .437 |
| MC-GPB (Heterogeneous links) | .532 | .584 | .453 | .552 | .530 | .471 |
| MC-GPB-hetero (Heterogeneous links) | .493 | .515 | .210 | .494 | .416 | .423 |

**Empirical results on large-scale datasets.** Here, we conduct an empirical study on two large-scale datasets for node property prediction, *i.e.*, ENZYME (6254 nodes and 23914 edges) and OGB-Arxiv (8532 nodes and 26281 edges). Detailed statistics are shown below. Specifically, we use GraphSAINT (Zeng et al., 2019) random walk sampler to extract the subgraph for illustration. The dataset split setting of train/validate/test sets is consistent with other datasets used in this work.

Table 18: Dataset statistics of the two large-scale datasets.

| dataset | # Nodes | # Edges | # Class | # Features | Hard homophily |
|---|---|---|---|---|---|
| ENZYME | 6254 | 23914 | 3 | 18 | 0.629 |
| OGB-Arxiv | 8532 | 26281 | 40 | 128 | 0.618 |

**(1)** Then, we evaluate the performance of MC-GRA with $\mathcal{K} = \{X, Y\}$, as shown in the table below (the higher the AUC value, the better the attack performance). As can be seen, MC-GRA is also effective on these two large-scale datasets. Besides, MC-GRA outperforms the baseline GRA method by a large margin.

Table 19: A comparison of attack methods on large-scale datasets (with AUC metric).

| dataset | ENZYME | OGB-Arxiv |
|---|---|---|
| GraphMI | .494 | .828 |
| MC-GRA | .761 | .891 |

**(2)** We also conduct the experiment of our defense method MC-GPB on these two datasets, with GraphMI as the attack method. As shown in the table below (the lower, the better), MC-GPB degenerates both kinds of GRA (*i.e.*, MC-GRA and GraphMI), which empirically proves the effectiveness of our defense method on large-scale datasets.

Table 20: A comparison of attack methods on large-scale datasets our defense method MC-GPB (with AUC metric).

| dataset | ENZYME | OGB-Arxiv |
|---------|--------|-----------|
| GraphMI | .488 (1.2%↓) | .533 (35.6%↓) |
| MC-GRA | .607 (20.2%↓) | .848 (4.8%↓) |

**Attacks without node feature.** We further implement the experiments without node features $X$ in the following. Note that the usair, brazil, and polblogs datasets have no initial node feature. Therefore, calculating $I(A; X)$ in Table. 1 with these datasets is infeasible due to the lack of $X$. Here, we present the attack results of our method on Cora, Citeseer, and AIDS datasets without using the node feature.

Table 21: A further quantitative comparison of attack methods without access to the node feature (with AUC metric).

| dataset | Cora | Citeseer | AIDS |
|---------|------|----------|------|
| Dot-Product (Tab. 2) | .849 | .907 | .521 |
| GraphMI (without $X$) | .802 | .759 | .575 |
| MC-GRA ($\mathcal{K} = \{\boldsymbol{H}_A\}$) | .834 | .887 | .575 |
| MC-GRA ($\mathcal{K} = \{\hat{\boldsymbol{Y}}_A\}$) | .771 | .890 | .540 |
| MC-GRA ($\mathcal{K} = \{\boldsymbol{Y}\}$) | .864 | .853 | .525 |
| MC-GRA ($\mathcal{K} = \{\boldsymbol{H}_A, \hat{\boldsymbol{Y}}_A\}$) | .828 | .918 | .525 |
| MC-GRA ($\mathcal{K} = \{\boldsymbol{H}_A, \boldsymbol{Y}\}$) | .875 | .919 | .539 |
| MC-GRA ($\mathcal{K} = \{\hat{\boldsymbol{Y}}_A, \boldsymbol{Y}\}$) | .867 | .896 | .539 |
| MC-GRA ($\mathcal{K} = \{\boldsymbol{H}_A, \hat{\boldsymbol{Y}}_A, \boldsymbol{Y}\}$) | .883 | .914 | .580 |

As shown in the above table, without using node features $X$ as the prior knowledge, the MC-GRA is still effective in recovering adjacency with considerable AUC results. Besides, the performance of MC-GRA is still better than the baselines. In fact, node features do not always exist, *e.g.*, for the Polblogs, USA, and Brazil datasets. While on the other hand, the characteristic adjacency $A$ is indispensable for graph learning. Besides, we further justify the feasibility of our MC-GRA without node features. For the extension of GRA, a more fine-grained study on graph properties is intriguing, *e.g.*, density, community, number of triangles *w.r.t.* adjacency $A$. To what extent can the above properties be recovered will shed insights into the power of GRA and the memorization effect of GNNs.

**Why would some of the model accuracy benefit from the defense mechanism?** As shown in Tab. 4, MC-GPB can also bring improvement in classification accuracy on partial datasets. We speculate that the reason is *forgetting more might also lead to learning better in some cases*. We provide a three-fold analysis from the information-theory perspective as follows.

**(1)** For brevity, we consider a simplified objective of the graph privacy bottleneck in Eq. (4), *i.e.*, to solve $\min -I(\boldsymbol{H}; Y) + \beta \cdot I(\boldsymbol{H}; A)$ *w.r.t.* representations $\boldsymbol{H}$, graph adjacency $\boldsymbol{A}$, and node labels $\boldsymbol{Y}$. The maximin game here is to encourage the accuracy by a higher $I(\boldsymbol{H}; Y)$, and reduce the complexity by regularizing the $I(\boldsymbol{H}; A)$ with $\beta$ for the trade-off.

**(2)** In this case, the spurious correlation measured by $I(\boldsymbol{H}; A|Y)$ will also be reduced and help the inference in test time. The reason is that absorbing too much irrelevant information between $\boldsymbol{A}$ and $\boldsymbol{Y}$, which can be superficially but not causally associated, will lead to degenerated test performance for GNNs (Zhao et al., 2022). Thus, a lower $I(\boldsymbol{H}; A|Y)$ here encourages the forgetting of adjacency and might bring a better generalization power in the test-time inference of GNNs. While the optimal case, *i.e.*, $I(\boldsymbol{H}; A|Y) = 0$, is also discussed in Theorem 6.4.

**(3)** Another supporting material is that only relying on a subgraph for reasoning can also boost the test-time performance (Miao et al., 2022). The method GSAT (Miao et al., 2022) aims to extract a subgraph $G_s$ as the interpretation. It inherits the same spirit of information bottleneck in its optimization, *i.e.*, $\min -I(G_s; Y) + \beta \cdot I(G_s; G)$. The integrated subgraph sampler can explicitly remove the spurious correlation or noisy information in the entire graph $G$.

Besides, we should note that in the cases where the model does not suffer from severe spurious correlations. The defense mechanism usually induces the drop trade-off regarding the model accuracy.

### B.3. Full qualitative results

**The recovered adjacency.** Fig. 9-11 shows the recovered adjacency of each dataset, which is grouped by node label under different attack strategies. In addition, we also provide the recovered adjacency on protected GNN, which is training with our proposed MC-GPB mechanism (sub-figure (d) and (e)). As can be seen, GNN training with MC-GPB successfully resists both attacks in terms of a larger amount of wrong prediction compared to the normal GNN. For example, for the Cora dataset, MC-GRA (Fig. 9(b)) achieves a better result compared to GraphMI (Fig. 9(c)) under both normal training strategy, in terms of fewer error predictions (red dots). Whereas MC-GPB successfully defended the MC-GRA (Fig. 9(d)) and GraphMI (Fig. 9(e)), and MC-GRA still have better performance compared to GraphMI with protected GNN.



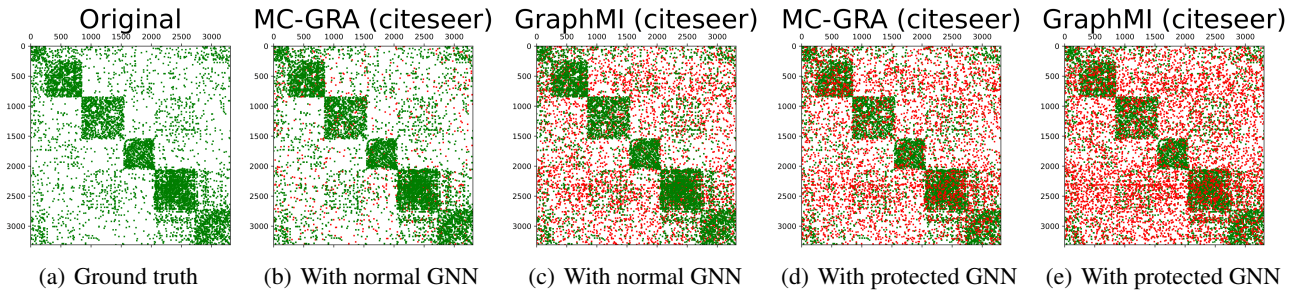Figure 9: Recovered adjacency on Cora dataset.



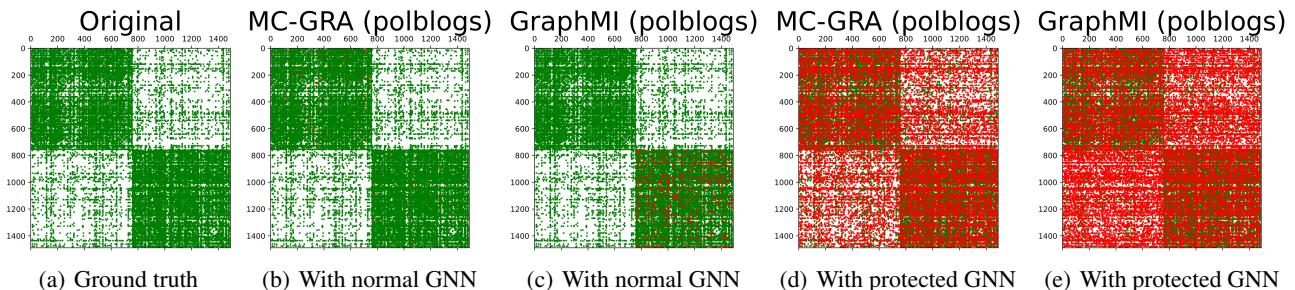Figure 10: Recovered adjacency on Citeseer dataset.



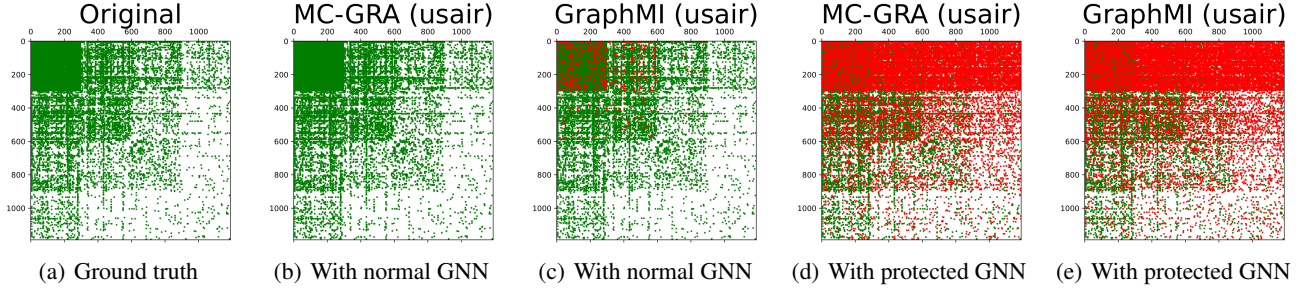Figure 11: Recovered adjacency on Polblogs dataset.

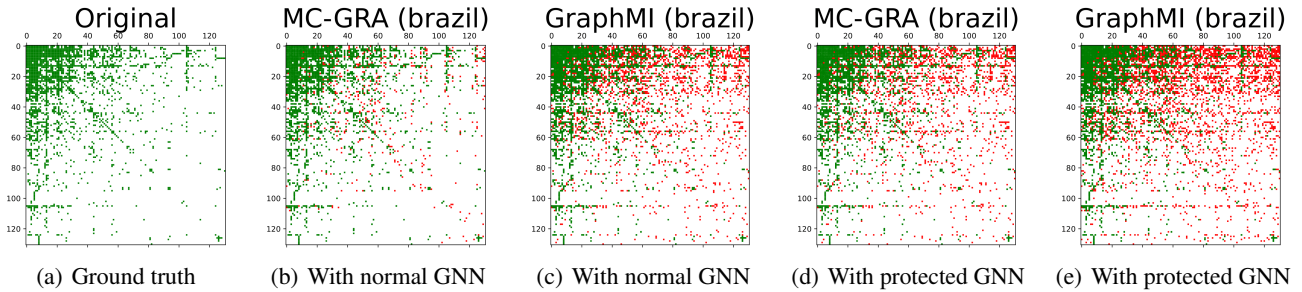Figure 12: Recovered adjacency on USA dataset.
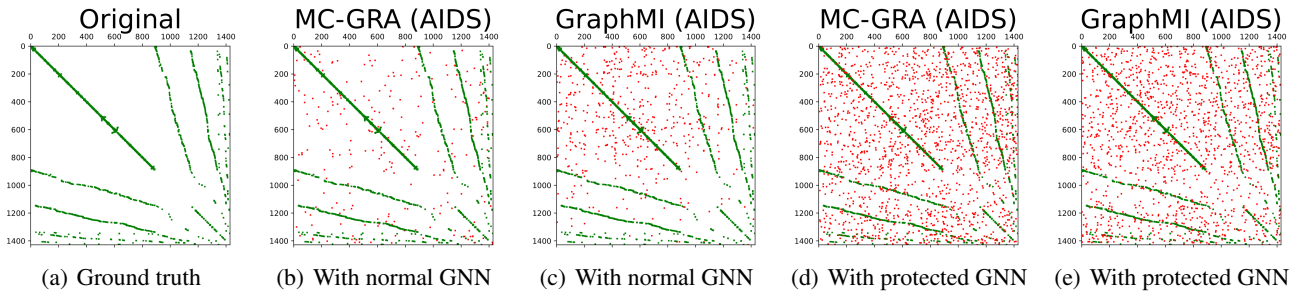


Figure 13: Recovered adjacency on Brazil dataset.



Figure 14: Recovered adjacency on AIDS dataset.

**Tracking the MI terms.** We show the learning curves of MC-GRA and MC-GPB on each dataset as follows.

For MC-GRA ( Fig. 15), most of the output and propagation loss converged to near zero, showing that the model efficiently approximates the original Markov chain. For MC-GPB, we track three constraints layer-wise and average them out to visualize the overall trend. We also record the model accuracy constraint, *i.e.*, the cross-entropy of model output, to check whether the layer and the full model have consistent patterns. Both privacy and complexity constraints, despite the fluctuation of the former in some datasets, show a downward trend throughout the training, especially for the usair dataset. The accuracy curves also have similar patterns.
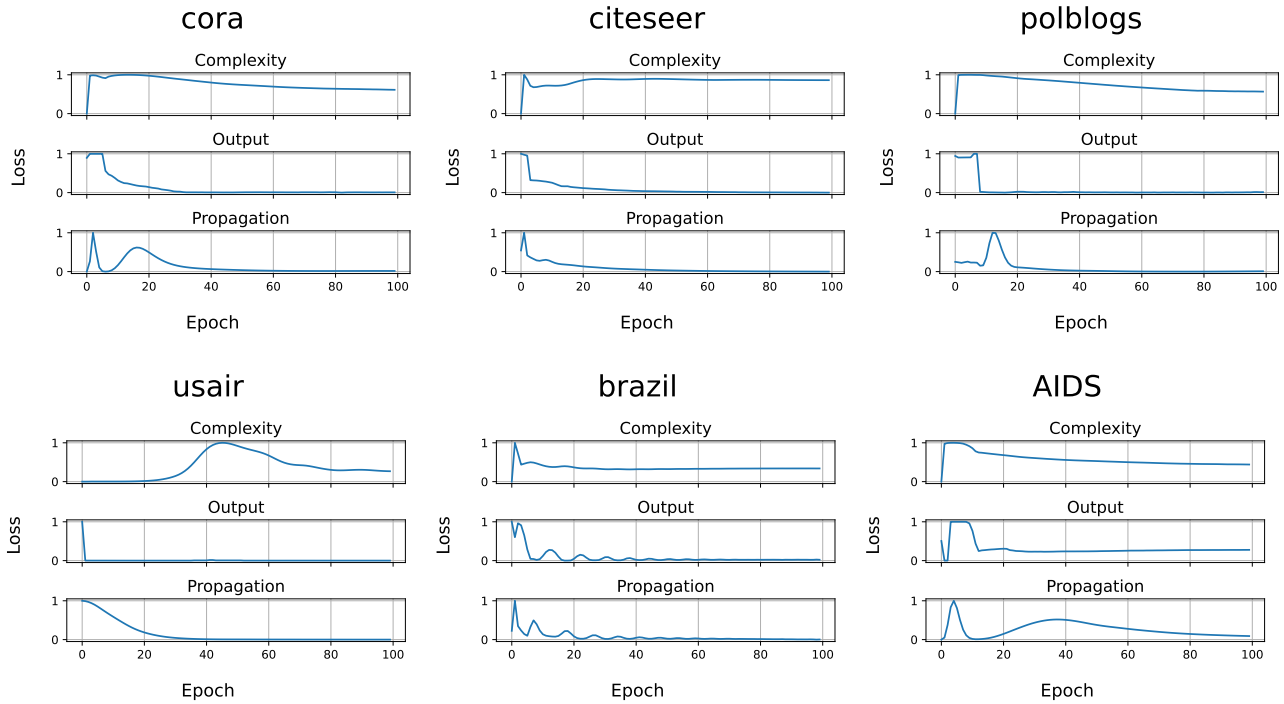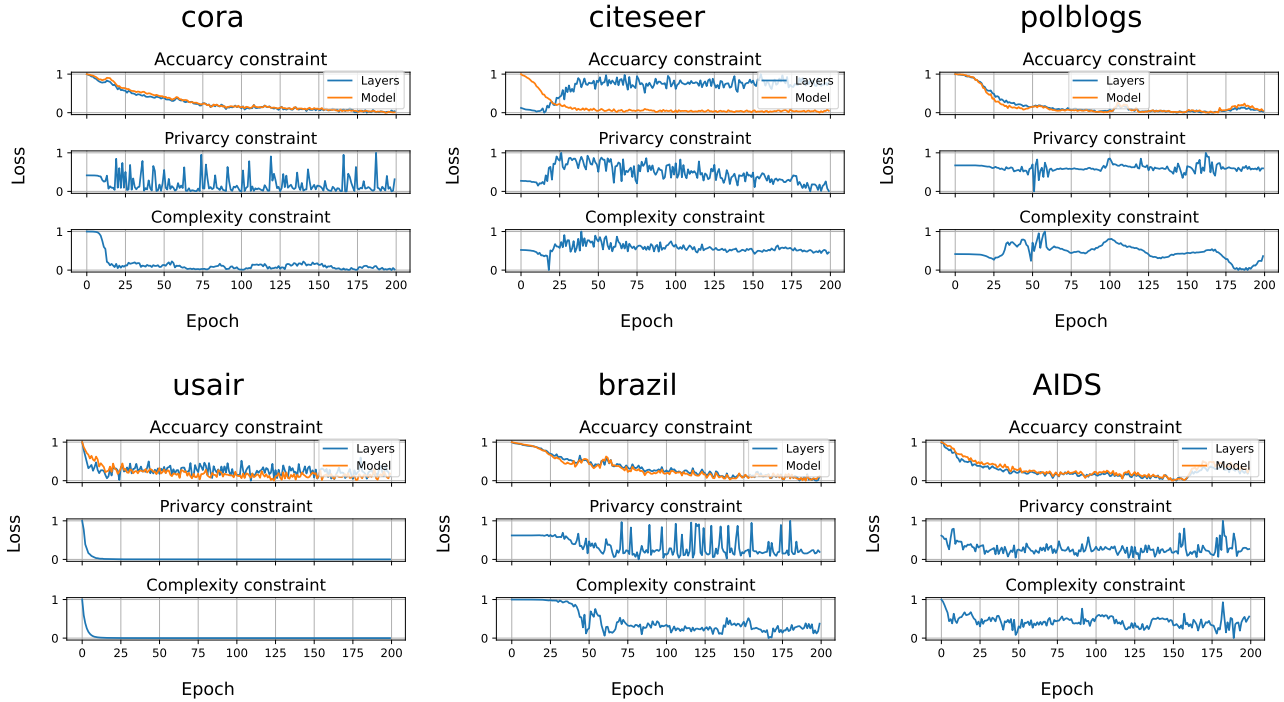
22

Figure 15: Training curves of MC-GRA on each dataset.



Figure 16: Training curves of MC-GPB on each dataset.

**A further analysis with the training dynamics.** We show the graph information planes with/without MC-GPB as follows. The model training without MC-GPB memorizes the privacy information at the beginning of training before gradually forgetting it. By applying our MC-GPB, we can enhance such forgetting procedures while preventing the model from discarding task-relevant information that might lead to a drop in accuracy.
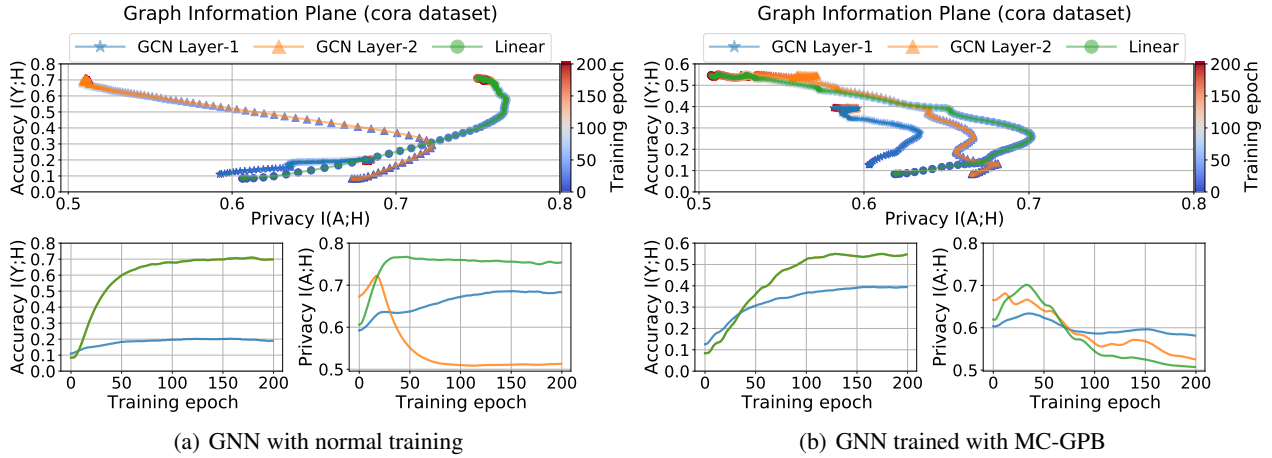
(a) GNN with normal training

(b) GNN trained with MC-GPB

Figure 17: Graph information plane on Cora dataset.



(a) GNN with normal training

(b) GNN trained with MC-GPB

Figure 18: Graph information plane on Citeseer dataset.



(a) GNN with normal training

(b) GNN trained with MC-GPB

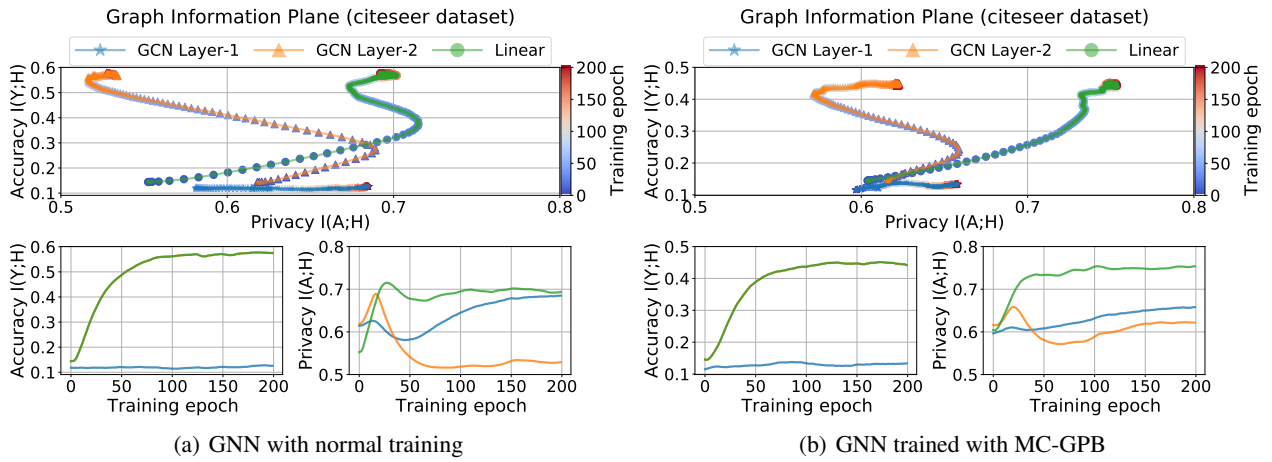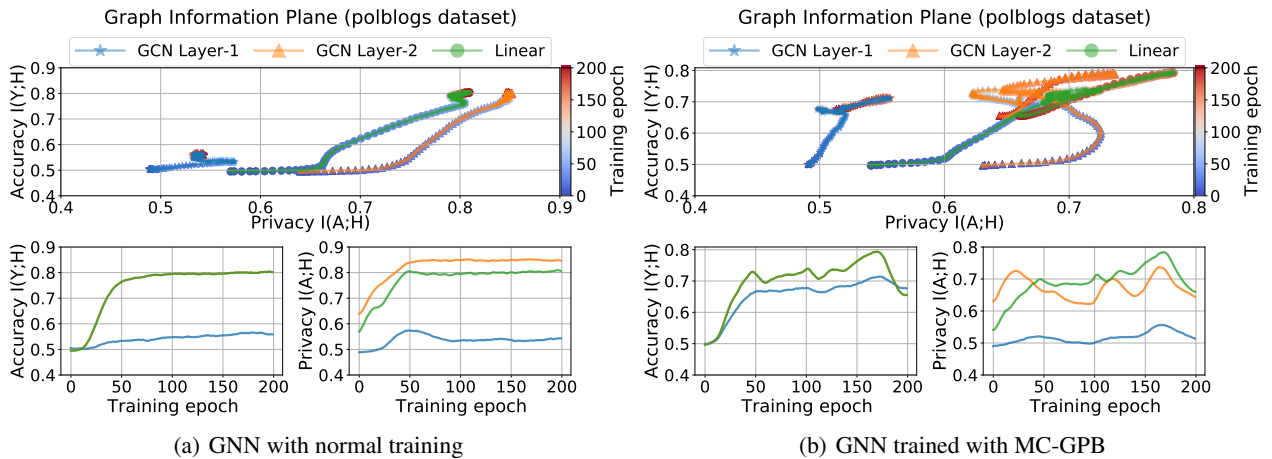Figure 19: Graph information plane on Polblogs dataset.
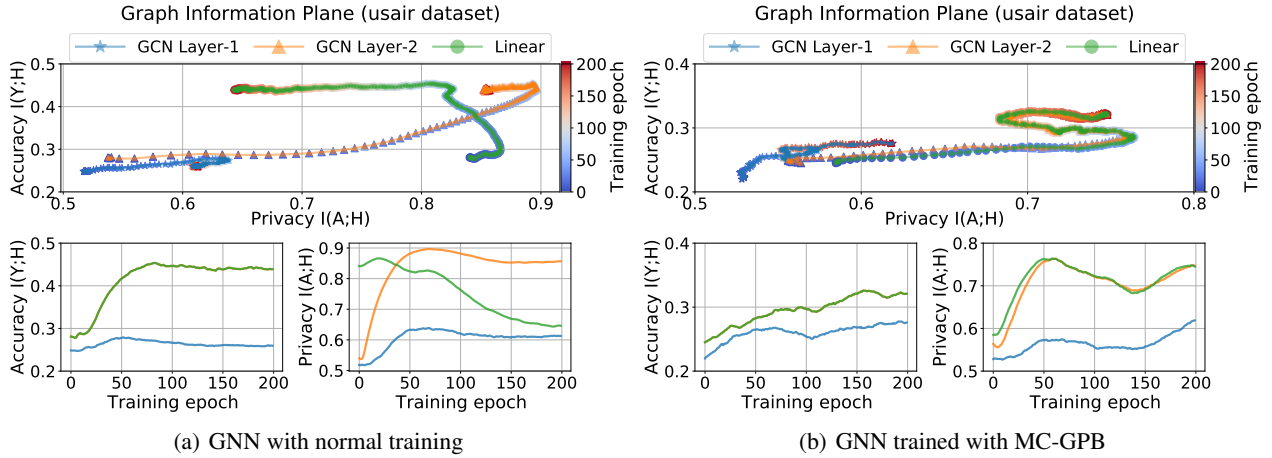
24

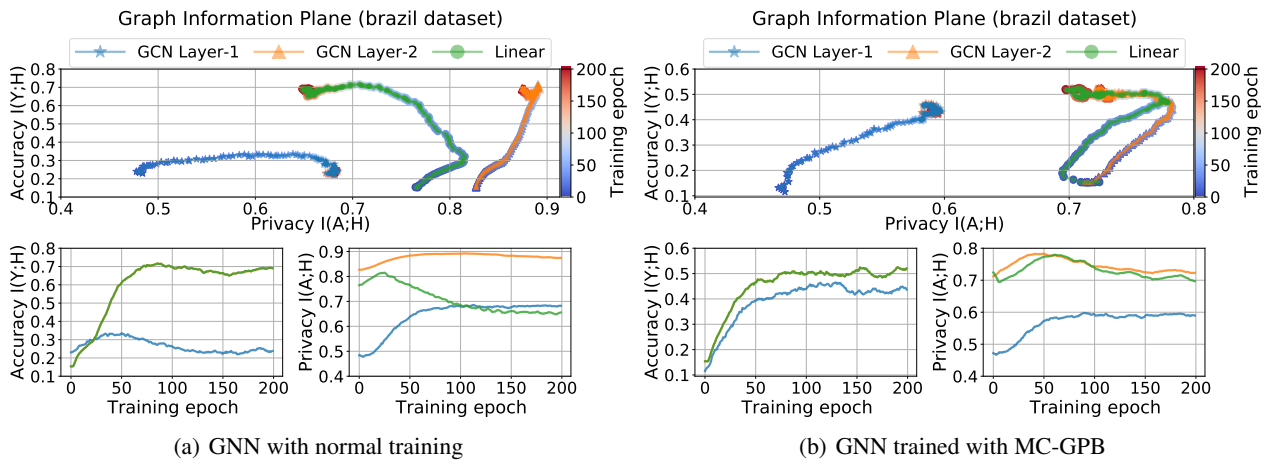Figure 20: Graph information plane on USA dataset.



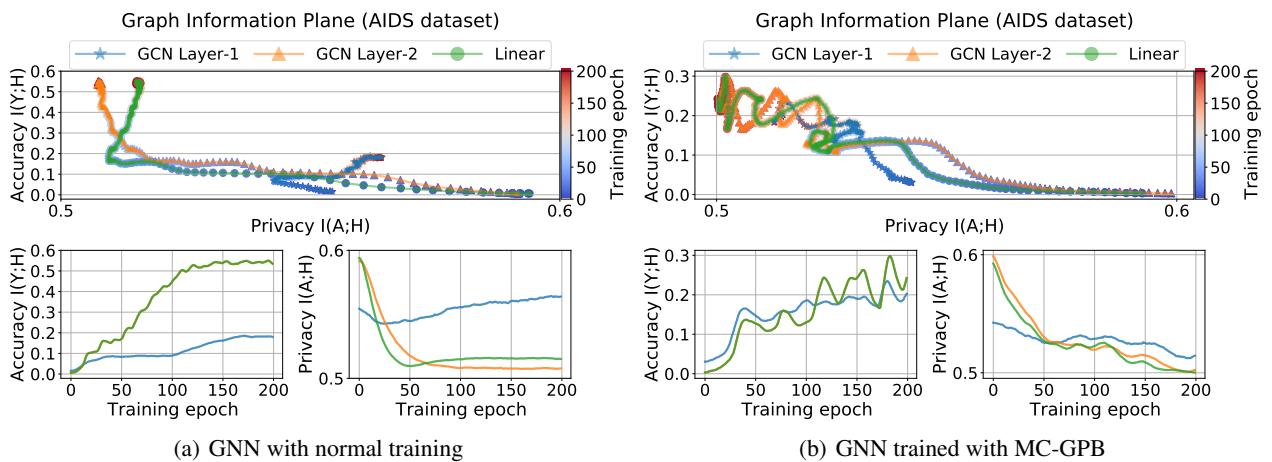Figure 21: Graph information plane on Brazil dataset.



Figure 22: Graph information plane on AIDS dataset.

# C. Related work

## C.1. Graph Neural Networks

Predicting node labels requires a parameterized hypothesis $f_{\boldsymbol{\theta}}(A, X) = \hat{Y}_A$ with GNN architecture (Kipf & Welling, 2016a; Veličković et al., 2018; Hamilton et al., 2017) and message propagation framework (Gilmer et al., 2017), where the architecture can be GCN (Kipf & Welling, 2016a), GAT (Veličković et al., 2018), or GraphSAGE (Hamilton et al., 2017). The forward inference of a $L$-layer GNN generates node representations $\boldsymbol{H}_A \in \mathbb{R}^{N \times D}$ by a $L$-layer message propagation.

Formally, let $\ell = 1 \dots L$ denote the layer index, $h_i^\ell$ is the representation of the node $i$, MESS$(\cdot)$ is a learnable mapping function to transform the input feature, AGGREGATE$(\cdot)$ captures the 1-hop information from neighborhood $\mathcal{N}(v)$ in the graph, and COMBINE$(\cdot)$ is the final combination between neighbor features and the node itself. Then, the $l$-layer operation of GNNs can be formulated as $\boldsymbol{m}_v^\ell = \text{AGGREGATE}^\ell(\{\text{MESS}(\boldsymbol{h}_u^{\ell-1}, \boldsymbol{h}_v^{\ell-1}, e_{uv}) : u \in \mathcal{N}(v)\})$, where the representation of node $v$ is $\boldsymbol{h}_v^\ell = \text{COMBINE}^\ell(\boldsymbol{h}_v^{\ell-1}, \boldsymbol{m}_v^\ell)$. After $L$-layer propagation, the final node representations $\boldsymbol{h}_e^L$ of each $e \in V$ are obtained. In addition, we summarize the detailed architectures of different GNNs in the following Table 22.

Then, the follow-up linear layer transforms $\boldsymbol{H}_A$ to classification probabilities $\hat{Y}_A \in \mathbb{R}^{N \times C}$, with $C$ categories in total. The training objective is to minimize the classification loss, *e.g.*, the cross-entropy between predictions $\hat{Y}_A$ and ground truth $Y$.

Table 22: Detailed architectures of different GNNs.

| GNN | MESS$(\cdot)$ & AGGREGATE$(\cdot)$ | COMBINE$(\cdot)$ |
|---|---|---|
| GCN | $\boldsymbol{m}_i^l = \boldsymbol{W}^l \sum_{j \in \mathcal{N}(i)} \frac{1}{\sqrt{\hat{d}_i \hat{d}_j}} \boldsymbol{h}_j^{l-1}$ | $\boldsymbol{h}_i^l = \sigma(\boldsymbol{m}_i^l + \boldsymbol{W}^l \frac{1}{\hat{d}_i} \boldsymbol{h}_i^{l-1})$ |
| GAT | $\boldsymbol{m}_i^l = \sum_{j \in \mathcal{N}(i)} \alpha_{ij} \boldsymbol{W}^l \boldsymbol{h}_j^{l-1}$ | $\boldsymbol{h}_i^l = \sigma(\boldsymbol{m}_i^l + \boldsymbol{W}^l \alpha_{ii} \boldsymbol{h}_i^{l-1})$ |
| GraphSAGE | $\boldsymbol{m}_i^l = \boldsymbol{W}^l \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} \boldsymbol{h}_j^{l-1}$ | $\boldsymbol{h}_i^l = \sigma(\boldsymbol{m}_i^l + \boldsymbol{W}^l \boldsymbol{h}_i^{l-1})$ |

## C.2. Privacy Attack on Graphs

Generally, the privacy attack on graphs can be attributed to *membership inference attack*, *model extraction attack*, and *model inversion attack*. Specifically, the membership inference attack (He et al., 2021b) aims to indicate whether a data sample is used to train the model. Besides, model extraction attack (Shen et al., 2022) is to extract information about the model parameters and reconstruct a surrogate model that behaves like the black-box model. Lastly, the model inversion attack aims to extract sensitive features of training data with only access to a trained model and non-sensitive features. We summarize the literature on the model inversion attack as follows.

## C.3. Inversion Attack on Graphs

As introduced before, most works of model inversion attack are investigated images and texts domains, leaving its effectiveness in the non-grid domain an open problem, *e.g.*, graph-structured data. While recently, several graph neural networks (GNNs) (Kipf & Welling, 2016a; Gilmer et al., 2017; Kipf & Welling, 2016b; Zhang & Chen, 2018) are proposed for graph data and boosted many real-world applications, *e.g.*, recommendation systems (Wu et al., 2020b) and drug discovery (Ioannidis et al., 2020). In graph scenarios, the target of the model inversion attack is to recover the topology of the training graph, *i.e.*, the connectivity properties *w.r.t.* each edge.

In practice, inferring links between nodes leads to a severe privacy threat when the links represent sensitive information, *e.g.*, the relationship between users in social networks. Besides, it may also compromise a model owner's intellectual property. The challenges of applying the model inversion attack to graphs are two folds. (1) The discrete nature of graph structure. It is hard to optimize in a differentiable way. Besides, the nodes and edges in a graph cannot be resized to the same shape. (2) Lack of domain knowledge as priors. Graph data are less intuitive than images, and the domain knowledge can be diverse, *e.g.*, molecules, social networks, and citation networks.

The pioneer works (Duddu et al., 2020; Chanpuriya et al., 2021) attempt to reconstruct the target graph from released graph embeddings $\boldsymbol{H} \in \mathbb{R}^{N \times D}$ of each node, that are generated by Deepwalk or GNNs. The specific attack method can be the Deepwalk Backward (Chanpuriya et al., 2021), or a decoder $f^{dec}(\boldsymbol{H}) : \mathbb{R}^{N \times D} \to \{0, 1\}^{N \times N}$ that is trained on auxiliary datasets (Duddu et al., 2020). Graph embedding attack with the auto-encoder framework is also an exciting direction as the

graph embeddings of each node can usually be accessed in practice. (Zhang et al., 2022b) systematically investigate the information leakage of graph embedding, and justify that the basic graph properties, *e.g.*, number of nodes, number of edges, and graph density, can be accurately extracted. Besides, it can determine whether a given subgraph is contained in the target graph or not. More importantly, it also shows that the graph topology can be recovered via conducting the MIA with graph embeddings.

The link stealing attack (He et al., 2021a) is the first work to *steal* links from a GNN as the target model. It aims to conduct the attack on black-box settings with three kinds of prior knowledge, including (1) node features, (2) target dataset's partial graph, and (3) a shadow dataset. This work proposed $8$ different kinds of attacks in total to be adaptive to the $2^3 = 8$ scenarios. Each proposed method for the attack was verified on chemical networks and social networks, which justified the feasibility of conducting a model inversion attack on graphs. However, it requires to be accessible to the partial graph and an auxiliary dataset. The partial graph actually contains sensitive information about the adjacency, and selecting the auxiliary dataset also requires extra information about the target graph. Thus, these methods cannot be directly utilized here. Besides, the GraphMI (Zhang et al., 2021b) is also a learnable attack that also aims to recover the links of the original graph. With the white-box access to the target model, the optimal adjacency is obtained by maximizing the classification probability *w.r.t.* given node labels $Y$, namely, $\hat{A}^* = \arg\max_{\hat{A}} \mathbb{P}(f_\theta(\hat{A}, X), Y)$. In practice, the $\hat{A}^*$ can be recursively updated by the projected gradient descent (PGD). The proposed attack method In this work is partially inspired by GraphMI, and it can be degenerated to GraphMI when the prior knowledge set is reduced to $\mathcal{K} = \{X, Y_{sub}\}$, where $Y_{sub} \subset Y$.

In addition to attacking GNNs trained for node classification tasks, a recent work (Zhang et al., 2022b) also attempted to attack the model for graph classification tasks. The shift from node-level tasks to graph-level tasks brings several unique challenges as the obtained one-dimensional embeddings $\boldsymbol{h}_\mathcal{G} \in \mathbb{R}^D$ are the compressed information of the whole graph $\mathcal{G}$. This work reconstructs the graphs with a graph auto-encoder that takes the graph embeddings as inputs and then generates the corresponding graphs. Note that the adopted graph auto-encoder is trained on an auxiliary dataset and then applied to the target dataset. It shares a similar spirit of generative attacks on degenerate that the generator (*e.g.*, a generative adversarial network) is pre-trained on public datasets.

## C.4. Model Inversion Attack on Images

Pioneer works (Fredrikson et al., 2014; 2015; Hidano et al., 2017) introduce the model inversion attack with comparably simple models, *e.g.*, linear regression, decision trees, and shallow networks. These early works justified the feasibility of model inversion attacks and succeeded in recovering the monochrome images. However, the reconstructed images are also usually of low fidelity (Szegedy et al., 2013), and they fail in attacking DNNs for image classification tasks.

*So, how can we recover the polychromatic and realistic images used for training?* Generative Model Inversion (GMI) (Zhang et al., 2020) is the first to conduct model inversion attacks on deep models, *i.e.*, the convolution neural networks. Instead of directly reconstructing the private data from scratch, its inversion process is guided by a distributional prior through the generative adversarial networks (GAN). Specifically, the used GAN is pretrained on public datasets for obtaining the generic prior knowledge of human faces via minimizing the canonical WassersteinGAN training loss, namely,

$$\mathcal{L}_{GAN}(G, D) = \mathbb{E}_x[D(x)] - \mathbb{E}_z[G(z)]. \tag{14}$$

that optimizes the generator $G$ and the discriminator $D$ in an adversarial training manner. Then, GMI introduces a diversity loss to encourage the more diverse generated images to recover more training patterns in $X^{tra}$. To be specific, with two sampled latent vectors $z_1, z_2$, the diversity loss is calculated as

$$\mathcal{L}_{div}(G) = \mathbb{E}_{z_1, z_2}\Big[\frac{||f_\theta^{feat}(G(z_1)) - f_\theta^{feat}(G(z_2))||}{||z_1 - z_2||}\Big]. \tag{15}$$

where $f_\theta^{feat}$ is the feature extractor of the target network. With these two loss functions, the full objective of GAN is as follows.

$$\min_G \max_D \mathcal{L}_{GAN}(G, D) - \lambda \mathcal{L}_{div}(G). \tag{16}$$

After training the GAN, GMI aims to find the latent vector $z$ that achieves the highest likelihood under the target network while being constrained to the data manifold learned by $G$, *i.e.*,

$$z^* = \min_z -D(G(z)) - \lambda \log[f_\theta(G(z))]. \tag{17}$$

where a lower prior loss $-D(G(z))$ require the more realistic images, while a lower identity loss $\log[f_\theta(G(z))]$ encourages the generated images to have a higher likelihood *w.r.t.* the targeted network. In summary, GMI conducts the model inversion attack in an end-to-end manner based on GANs that can reveal private training data of the target model with high fidelity, which make up for the deficiency of the early works. Besides, it also reveals the non-convex nature of model inversion attacks on deep models, where a more powerful target model can exhibit a higher privacy risk.

However, the top-one identification accuracy of face images inverted from the classifier is not that high. Is it because CNNs do not *memorize* much about private data or is it due to the *imperfect* attack algorithm? To answer, the follow-up work, Knowledge-Enriched Distributional Model Inversion (KED-MI) (Chen et al., 2021), shows that the target network maybe not be fully utilized. KED-MI further distills the useful knowledge from the target model with two designs. On the one hand, instead of generating and discriminating real or fake samples, DMI utilizes the target model to generate soft labels for supervising the GAN, *i.e.*, to minimize the loss $\mathcal{L}_{GAN} = \mathcal{L}(D) + \mathcal{L}(G)$, specifically,

$$
\begin{aligned}
\mathcal{L}(D) =& \mathcal{L}_{sup}(D) + \mathcal{L}_{unsup}(D) \\
\mathcal{L}_{sup}(D) =& -\mathbb{E}_{x \sim p_{data}(x) \sum_{k=1}^{K}} f_\theta(x) \log p_{disc}(y = k|x) \\
\mathcal{L}_{unsup}(D) =& -\mathbb{E}_{x \sim p_{data}} \log D(x) + \mathbb{E}_{z \sim noise} \log(1 - D(G(z))) \\
\mathcal{L}(G) =& ||\mathbb{E}_{x \sim p_{data}} f_\theta(x) - \mathbb{E}_{z \sim noise} f_\theta(G(z))||_2^2 + \lambda \mathcal{L}_{ent}
\end{aligned}
\tag{18}
$$

where the entropy regularization term $\mathcal{L}_{ent}$ is taken from previous work (Grandvalet & Bengio, 2004).

On the other hand, no longer recovering a sample given a label in a one-to-one manner, DMI explicitly parameterizes the distribution of private data and proceed with the model inversion attack in a new many-to-one way. Technically, the latent vectors of the generator are sampled from a learnable distribution to capture the class-wise information, while the discriminator acts as a (K+1)-classifier to differentiate the K classes of private data. Here, the latent variable $z$ is parameterized by $z = \sigma\epsilon + \mu$ by the reparameterization trick that the corresponding distribution $p_{gen}$ samples the optimal $z^*$ as follows.

$$
z^* = \min_z -\mathbb{E}_{z \sim p_{gen}} \log D(G(z)) - \lambda \mathbb{E}_{z \sim p_{gen}} \log[f_\theta(G(z))].
\tag{19}
$$

Basically, a successful attack should generate realistic and diverse samples. *So, how can we generate more diverse samples?* The Variational Model Inversion (VMI) (Wang et al., 2021) further formulates the model inversion attack as the variational inference. VMI generally can bring a higher attack accuracy and diversity for its equipped powerful generator StyleGAN to optimize its designed variational objective. Specifically, for the target class $y$, VMI approximates the target posterior with a variational distribution $q(x) \in Q_x$ from the variational family $Q_x$. The target model $f_\theta(x)$ is then denoted as $p_{\text{TAR}}(x|y)$. The variational objective is derived as follows.

$$
\begin{aligned}
q^*(x) =& \min_{q \in Q_x} D_{KL}(q(x)||p_{\text{TAR}}(x|y)) \\
=& \min_{q \in Q_x} \mathbb{E}_{q(x)} \big[ -\log p_{\text{TAR}}(x|y) + D_{KL}(q(x)||p_{\text{TAR}}(x)) \big]
\end{aligned}
\tag{20}
$$

In addition to recovering images, *can model inversion attack be applied to other extensions?* The Contrastive Model Inversion (CMI) (Fang et al., 2021) aims for data-free knowledge distillation. It recovers the training data from the target model via model inversion attacks, based on which it trains a student model. To overcome the mode collapse problem that recovered images are highly similar to each other, CMI proposes the contrastive learning objective upon the generated data to promote diversity while remaining considerable fidelity. With the similarity measurement $sim(x_1, x_2, h) = cos(h(x_1), h(x_2) = {}^{<h(x_1),h(x_2)>}/_{||h(x_1)|| \cdot ||h(x_2)||}$, where the $h(\cdot)$ projects $x_i$ to the embedding space, the contrastive loss is formulated as

$$
\mathcal{L}_{con}(X, h) = -\mathbb{E}_{x_i \in X} \left[ \log \frac{exp(sim(x_i, x_i^+, h)/\tau)}{\sum_j exp(sim(x_i, x_j^-, h)/\tau)} \right]
\tag{21}
$$

Besides, XAI-aware model inversion attack (Zhao et al., 2021) shows that the additional knowledge collected from the model inference procedure can promote the model inversion attack performances. In detail, if the model explanations *e.g.*, saliency maps of Gradients or Grad-CAM, are attainable in practice, it might do harm to privacy since these explanations can help recover private data.

In addition, (Kahla et al., 2022) conducts the first label-only model inversion attack only accessing the model's predicted labels without the confidence scores. As a machine learning model is often packed into a black-box that only generates the hard label (*i.e.*, the label of the class with the highest probability), such an attack scenario is more practical but also much more challenging to perform. Despite requiring less knowledge about the target model, this work justifies that such a black-box attack is also feasible and effective. Specially, it attempts to generate the most likelihood images for the target class, and observes that a region of high likelihood shall be located in the center of the class. Based on this observation, this work proposes to iteratively move the generated image away from the decision boundary and closer to the center.

Recent advance (Struppek et al., 2022) significantly decreases the cost of conducting a model inversion attack through relaxing the dependency between the target model and the image prior. This work enables the use of a single GAN to attack a wide range of targets, requiring only minor adjustments to the attack. Moreover, this work shows that the model inversion attack is possible even with publicly available pre-trained GANs and under strong distributional shifts.

### C.5. Model Inversion Attack on Texts

In addition to recovering training images with visual models introduced before, model inversion attack on text data with language models is also attacking more and more interests. In this domain, the input $X$ is changed from image to text (*i.e.*, sentences), and the architecture of model $f_\theta$ is also shifted from the convolutional neural network to the transformer-based language model.

A pioneer work (Carlini et al., 2021) demonstrates that large language models (*i.e.*, the GPT-2) memorize and leak individual training examples, with only black-box query access. The private information of an individual person can be accurately recovered. More importantly, this work reveals that some worst-case training examples are indeed memorized, although training examples do not have noticeably lower losses than test examples on average. Such a phenomenon is correlated with the memorization effect of DNNs and deserved further investigations.

Besides, another work, Text Revealer (Zhang et al., 2022a), firstly proposed to apply the model inversion attack to text classification with the transformer-based pretrained language models. Its attack consists of two stages: (1) collect texts from the same domain as the public dataset and extract high-frequency phrases from the public dataset as templates; (2) train a language model as the text generator on the public dataset. By minimizing the text classification loss, *i.e.*, cross-entropy, the generated text distribution becomes closer to the private dataset.

### C.6. On defending Model Inversion Attack

As for defending against the model inversion attack, a natural solution can be differential privacy (DP). Although effective in defending the membership attack (Abadi et al., 2016), the techniques of DP are proved to be ineffective with model inversion attacks (Fredrikson et al., 2014; Zhang et al., 2020).

The mutual-information-based defense (MID) (Yang et al., 2019) and Bilateral Dependency Optimization (BiDO) (Peng et al., 2022) are two representative defense methods that are specially designed for the model inversion attack. They follow a similar principle, that is, to control the mutual information among inputs $X$, hidden representations $Z$, prediction outputs $\hat{Y}$, and labels $Y$. Specifically, MID directly decreases the mutual information between $I(X; \hat{Y})$. BiDO forces the model to learn the robust representations by minimizing $I(X; Z)$ to limit redundant information that is transferred from the inputs to the latent representations, while maximizing $I(Y; Z)$ to keep the representations informative.

These two robust methods are effective in defending against model inversion attacks. The recovered images are neither correct nor realistic. However, such defense methods can do harm to the performance of the target model, as the informative signals from the input side can be overlooked under the balance of mutual information. Thus, a better trade-off between the model inversion-robustness and prediction performance is expected. In general, such an area of defending against model inversion attacks is still under-explored.

In general, the principle of conducting the model inversion attack is to utilize prior knowledge as much as possible, to extract more information from the target model, in order to generate more realistic and diverse samples. While defending against the model inversion attack, one promising solution is to store less information about input data in the weights of the model. In this way, the attacker is unable to recover the private data via querying the target model. However, it usually forms a trade-off between privacy and accuracy that such privacy-safe solutions can harm the accuracy, Thus, a better defensive approach is needed. The model inversion defense in practice where several trained models are expected to be protected without further modifications are much more challenging and essential.

# D. A Further Discussion on Graph Reconstruction Attack

## D.1. Background of the research problem

In this part, we would further clarify the background and settings of the research problem in our work, *i.e.*, Graph Reconstruction Attack. To be specific, we provide rigorous answers to the three following research questions.

- *Q1. About the black-box or white-box attack settings regarding accesses to the target model.* Here, the black-box setting indicates that the attackers can only query the target model and receive the classification results, which is the most difficult setting for the adversary. When a company employs GNN tools from another company that could be considered an adversary who possesses black-box access to the GNN model. On the contrary, the white-box setting means that the entire parameters of the target model can be obtained. Here, white-box attacks have created an increasingly serious threat to privacy due to the rise in the number of internet venues and online platforms where users can download full models.

- *Q2. About accesses to the prior knowledge.* The GNN model prediction results $\hat{Y}_A$ shared by many departments within the same corporation may be accessed by an adversary. For instance, to train a GNN model for fraudulent account detection, a social network service provider uses the technology of another business. In this situation, the supplier frequently must send the business the nodes' prediction results in order to debug or improve them. Similar circumstances apply to representations $H_A$ (*i.e.*, node embeddings), which are typically released. Furthermore, GNNs use the message passing framework, as is common knowledge, to produce representations of each node that are used in downstream tasks like node classification and link prediction. Additionally, a prior study (He et al., 2021a) also took into account the availability of an auxiliary dataset and a partial graph $A_{sub} \subseteq A$. The additional prior knowledge does successfully improve the GRA, even though it necessitates more access. As a result, the GRA can be carried out if the attacker can access the trained GNN model from malicious clients and has some leaked prior knowledge that is connected to the model.

- *Q3. About the attack target i.e., the adjacency rather than the node feature.* Intuitively, the adjacency $A$ and the node feature $X$ are all located on the input sides of the forward Markov chain, which means both of them can be the inversion target. The key motivation to attack the adjacency is two folds, *i.e.*, its practical risks and understandability to human beings. For instance, social network data is gathered with user consent in order to train GNNs for better service, such as friend classification or ad recommendation. It should be noted that user friendship data is sensitive and relational and that it should be kept secret. If the user's friendship is recovered in this scenario, the user's privacy will be compromised. A bigger security risk is presented by the fact that attackers can grasp such privacy, which is more critical. As a result, the privacy of graph adjacency data should receive greater attention and protection because it is more sensitive and intelligible than the node feature.

## D.2. Graph Reconstruction Attack in practice

Here, we provide a detailed explanation for the existence of the threat model in practice with several real-world examples.

*Q1. Why can adjacency be attacked, and should be protected in practice?* In practice, inferring links between nodes leads to a severe privacy threat when the links represent sensitive information, *e.g.*, the relationship between users in social networks. Taking social networks as an example, which require gathering interaction among individuals, GNNs can only have satisfied performance on downstream tasks like community detection or ad recommendation once the network structure is accurately characterized. However, this connection among users should be private since it is gathered with user consent and shall be kept between the service provider and users.

For instance, to train a GNN model for fraudulent account detection, a social network service provider uses the technology of another business. In this case, the supplier will frequently send the business the nodes' prediction results to debug or improve them. Similar circumstances apply to node representations, which are typically released. Thus, the model's outputs shared by many departments within the same corporation can be accessed by an adversary. The same to node features and node labels. Note that the graph reconstruction attack can be conducted with only a subset of the above informative variables, as we have empirically justified its feasibility in Section 7. The attack target here, *i.e.*, links, can reflect the model owner's sensitive relationship information or intellectual properties, which brings considerable safety risk that is orthogonal to the well-known and widely-studied adversarial attacks (Dai et al., 2018; Chen et al., 2022; Zhu et al., 2023). The inversion of adjacency is a severe privacy threat in several real-world scenarios of GNNs, which have been widely used in recommendation systems, social networks, citation networks, and drug discovery.

Thus, the users' privacy should be paid attention to and protected, especially for personal relationships and sensitive

information. The private connection among individuals shall be protected since it might become a powerful weapon for fraud syndicates to generate fake identities or threaten people with concerns about their secret relationships. Moreover, in the drug discovery scenario, the company might train its graph generation model for searching for new medicine and share its model with other companies for further development. However, the dataset for training the model contains private drug structures on the market, which is worth stealing as it involves tremendous efforts to discover and test the drug's safety before appearing on the market.

*Q2. Why investigating the GRA is meaningful and practical?* One cannot ignore the importance of privacy leaking of existing GNNs, regardless of the possibility of the attack. The model inversion attack also receives noticeable attention in the visual domain (Fredrikson et al., 2014; Hidano et al., 2017; Chen et al., 2021; Wang et al., 2021; Zhao et al., 2021; Kahla et al., 2022) and the natural language processing domain (Carlini et al., 2021; Zhang et al., 2022a). For the graph domain, previous works (He et al., 2021a; Zhang et al., 2021b) have justified the practical feasibility and possible impact on real life.

The GRA, or the general MIA, is a well-defined problem that attacks growing interests. We have justified that the adjacency matrix contains rich private information and is prone to attack. The purpose of this work is to illustrate the flaw of the current GNNs training process and provide a direction to protect the model. One cannot wait until a mistake is made to fix it.

### D.3. Issues of existing attack or defense methods

In this part, we further elaborate the challenges of the studied GRA problems and issues of existing methods (introduced in Appendix. C), which are summarized in the following three folds.

- Directly applying existing methods (Fredrikson et al., 2015) to graphs can be easily sub-optimal. The attribute is that they are originally designed for grid data like images, overlooking the inherently topological and semantical properties of graphs. Besides, another modality gap is the lacking of a distributional prior (*e.g.*, a public face dataset) stored in a pre-trained generative adversarial network that is used to guide the inversion process of graphs. Thus, it hinders several generative attack methods (Zhang et al., 2020; Struppek et al., 2022) to be applicable to graphs.

- Considering the inductive nature that graphs can be collected from diverse domains, the fetched prior knowledge set $\mathcal{K}$ is of vital importance. However, $|\mathcal{K}|$ can be 2, 3, or 4 (with 7 combinations in total), while trivially treating each case with a specific method is quite non-general. Thus, how the adaptively utilized each $\mathcal{K}_i \in \mathcal{K}$ in the form of *one* generic objective of combination optimization, is the main challenge to solving the non-convex problem here.

- As for the defense, the differential privacy techniques are proven to be of little help to defend against MIA (Zhang et al., 2020; 2021b), although it can be helpful to defend against the membership attack. On the other hand, the improvement of privacy guarantee might come at the cost of degenerating the empirical performance (Bietti et al., 2022). Thus, an effective defense method customized for GNNs that nicely balancing of accuracy and privacy is expected.

**Technical contribution.** To better clarify the technical contribution of our work, we provide a brief summary with regard to existing works as follows.

- For the attack, we propose the Markov Chain-based Graph Reconstruction Attack (MC-GRA) that boosts the attack fidelity with parameterization techniques and injected stochasticity. Unlike existing GRA methods designed for ad-hoc scenarios, our proposed MC-GRA aims to locate, present, and utilize the interplaying variables of GNN forward in a generic manner. It can adaptively support the white-box GRA that leverages the target model and any prior knowledge. That is, *to recover better, you must extract more.*

- For the defense, existing works have justified that differential privacy (DP) is ineffective with GRA (or general MIA). And currently, there is a lack of an effective way to defend GRA. In this work, we propose the Markov Chain-based Graph Privacy Bottleneck (MC-GPB), an information theory-guided principle that significantly degenerates GRA with only a slight accuracy loss. The MC-GPB requires the GNN to forget the privacy information in the training process, *i.e.*, to make the learned representations contain less information about adjacency. That is, *to learn safer, you must forget more.*

- To the best of our knowledge, we are the first to conduct a systematic study of GRA from both sides of attack and defense. By taking GNNs as a Markov chain and attacking GNNs via a flexible chain approximation, we systematically explore the underlying principles of GRA and reveal several essential phenomena. In addition, we also provide a rigorous analysis from information-theoretical perspectives to disclose several valuable insights on how to strengthen and defend GRA.

### D.4. The information-theoretic principles of GRA

Basically, the objective of the attack is to recover the adjacency, as was stated earlier. On the other hand, the defense consists of acquiring a dependable and thoroughly trained model that is resistant to assault. In order to launch an attack, one must first collect and combine all of the relevant prior knowledge and then do a backward recovery concerning the adjacency; As a sort of defense, rather than that, it is necessary to mandate that the GNN forget all of the information on the adjacency during its training phase.

Specifically, the correlation between each $\mathcal{K}_i \in \mathcal{K}$ and its counterpart in the forward process with the recovered adjacency $\hat{A}$ should be encouraged to enhance the GRA. For example, $\hat{Y}_A \in \mathcal{K}$ should be approximated by $\hat{Y}_{\hat{A}}$, *i.e.*, a higher $I(\hat{Y}_A; \hat{Y}_{\hat{A}})$, that is essential to obtain a high fidelity $I(A; \hat{A})$. On the contrary, constraining the correlation between intermediate variables and the original adjacency, *e.g.*, a lower $I(A; \boldsymbol{H}_A^i)$ for $i = 1, 2, \cdots, L$, is a natural solution to defend the GRA. As such, even these variables are leaked, the attacker is scarcely possible to recover $A$.

### D.5. Deriving the MC-GRA objective

For approximating $A$ by $\hat{A}$, the basic objective of attacking and its relaxed form to optimize $\hat{A}$ are derived as follows.

**The basic attack objective.** Intuitively, given a prior knowledge set $\mathcal{K} \subseteq \{X, Y, \boldsymbol{H}_A, \hat{Y}_A\}$, the optimal recovered adjacency $\hat{A}^*$ can be obtained by directly maximizing its correlation with each term $\mathcal{K}_i \in \mathcal{K}$, *i.e.*, solving the Basic-GRA,

$$\hat{A}^* = \arg\max_{\hat{A}} I(\hat{A}; \mathcal{K}) \triangleq \sum_{\mathcal{K}_i \in \mathcal{K}} \alpha_i I(\hat{A}; \mathcal{K}_i). \tag{22}$$

where the hyper-parameters $\{\alpha_i\}_{i=1}^{|\mathcal{K}|}$ balance the MI terms $\{I(\hat{A}; \mathcal{K}_i)\}_{i=1}^{|\mathcal{K}|}$. Intuitively, maximizing $I(\hat{A}; \mathcal{K})$ enables to extract information in $\mathcal{K}$ and store it in $\hat{A}$ to approximate $A$, namely,

$$\max_{\hat{A}} H(\hat{A}) \approx H(\mathcal{K}) \;\Rightarrow\; \max_{\hat{A}} I(A; \hat{A}) \approx I(A; \mathcal{K}). \tag{23}$$

The Basic-GRA can be applied to *black-box* settings, however, it can also be sub-optimal as locations of $\hat{A}$ and $\mathcal{K}_i$ are distant: $\hat{A}$ is in the front-end while $\mathcal{K}_i$ is in the back-end. Which means, the information unrelated to adjacency induced by the `ORI-chain`, *i.e.*, $H(\mathcal{K}_i|A)$, will be also stored in $H(\hat{A})$. That is, $\max_{\hat{A}} I(A; \hat{A}) \approx I(A; \mathcal{K})$ might come at the cost of $H(\hat{A}) \approx H(\mathcal{K})$. Besides, the knowledge stored in the target model $f_{\boldsymbol{\theta}^*}(\cdot)$ is entirely not utilized. Thus, the recovered $\hat{A}^*$ is not good enough, and a refined objective is required.

**The relaxed attack objective.** For extracting the target model and relaxing the optimization simultaneously, we replace $\hat{A}$ in Eq. (22) by the latent variable $\boldsymbol{Z}_{\hat{A}}^j$ generated by `GRA-chain` $f_{\boldsymbol{\theta}^*}(\hat{A}, X)$. We promote $I(\boldsymbol{Z}_A^j; \boldsymbol{Z}_{\hat{A}}^j)$ instead of $I(\boldsymbol{Z}_A^j; \hat{A})$, as it provides supervision signals that can be tractably approximated. Here, $\boldsymbol{Z}_{\hat{A}}^j$ shares the same location as $\mathcal{K}_i$ in the chain ($\boldsymbol{Z}_A^j = \mathcal{K}_i$). The derived new objective is

$$\text{MC-GRA:} \quad \hat{A}^* = \arg\max_{\hat{A}} \quad \underbrace{\alpha_p I(\boldsymbol{H}_A; \boldsymbol{H}_{\hat{A}}^i)}_{\text{propagation approximation}} + \underbrace{\alpha_o I(\boldsymbol{Y}_A; \boldsymbol{Y}_{\hat{A}}) + \alpha_s I(Y; \boldsymbol{Y}_{\hat{A}})}_{\text{outputs approximation}} - \underbrace{\alpha_c H(\hat{A})}_{\text{complexity}}. \tag{24}$$

Note that MC-GRA is a maximin game: it maximizes the approximation of encoding and decoding processes of the two Markov chains, while minimizing the complexity to avoid trivial solutions by constraining the graph density. Compared with Basic-GRA (Eq. (22)), it approximates latent variables $\mathcal{S}_A$ in `ORI-chain` by $\mathcal{S}_{\hat{A}}$ in `GRA-chain`, *i.e.*,

$$\max_{\hat{A}} I(\mathcal{S}_A; \mathcal{S}_{\hat{A}}) \;\Rightarrow\; \max_{\hat{A}} I(A; \hat{A}), \;\; s.t. \;\; \mathcal{S}_A = \{\boldsymbol{Z}_A^i : \boldsymbol{Z}_A^i \in \mathcal{K}\}, \; \mathcal{S}_{\hat{A}} = \{\boldsymbol{Z}_{\hat{A}}^i : \boldsymbol{Z}_A^i \in \mathcal{S}_A\}. \tag{25}$$

## E. Implementation Details

In this section, we provide a detailed introduction to the technical designs and implementation details of the proposed methods in our work. Specifically, Appendix. E.1 and Appendix. E.2 are technical details of the empirical study in Sec. 4, Appendix. E.3 and Appendix. E.4 are elaboration for the methodology in Sec. 5 and Sec. 6.

### E.1. Quantifying privacy leakage

As introduced in Sec. 4, studying the direct correlation between ground-truth adjacency $A$ and single variable $\boldsymbol{Z} \in \{X, Y, \boldsymbol{H}_A, \hat{\boldsymbol{Y}}_A\}$ can provide insight into the studied GRA problem. The informative concept of mutual information (MI) that $I(X; Y) = D_{KL}[p(x, y) \| p(x)p(y)] = H(X) - H(X|Y)$ is a measure of the symmetric correlation between the two variables, which suits our needs perfectly. To avoid the cumbersome calculation of MI, a surrogate estimation *w.r.t.* the existence of edges in $A$, *i.e.*, the AUC (area under the curve) metric is utilized (Zhang & Chen, 2018; Zhu et al., 2021) to quantify $I(A; \boldsymbol{Z})$ regarding edges in $A$ and $\hat{A}_{\boldsymbol{Z}} = \sigma(\boldsymbol{Z}\boldsymbol{Z}^\top)$ can be an efficient solution here.

Specifically, suppose $A \in \{0, 1\}^{N \times N}$ and $\boldsymbol{Z} \in \mathbb{R}^{N \times D}$, where $N$ is the number of nodes and $D$ is the hidden size. We define $I(A; \boldsymbol{Z}) \triangleq I(A; \hat{A}_{\boldsymbol{Z}})$, where $\hat{A}_{\boldsymbol{Z}} = \sigma(\boldsymbol{Z}\boldsymbol{Z}^\top) \in \mathbb{R}^{N \times N}$ indicates the predictive existence of each edge. Note that the dot product $\boldsymbol{Z}\boldsymbol{Z}^\top$ followed by the activate function $\sigma(\cdot)$ is in direct proportion to the cosine similarity, which is commonly utilized in investigating the distribution of representations. Here, a higher MI $I(A; \boldsymbol{Z})$ indicates a lower expectation of distance $\mathbb{E}_{(i,j) \sim A} d(z_i, z_j)$, where measurement $d(\cdot, \cdot)$ can be cosine, euclidean, etc. Alternatively, the activation function $\sigma(\cdot)$ can be ReLU, Sigmoid, etc. For quantifying $I(A; \hat{A}_{\boldsymbol{Z}})$, the AUC (area under the curve) is utilized as the metric (Zhang & Chen, 2018; Zhu et al., 2021) regarding edges in factual adjacency $A$ and recovered adjacency $\hat{A}_{\boldsymbol{Z}}$.

### E.2. Graph information plane

Recall that in Sec. 4.3, we track the aforementioned MI terms in the training process for further study. Here, the training dynamics of representations H in each layer are projected to the two dimensional $\big(I(A; \boldsymbol{H}), I(Y; \boldsymbol{H})\big)$ plane. The $I(X; \boldsymbol{H})$ is not considered as node features do not always exist, while the characteristic adjacency $A$ and labels $Y$ are indispensable for supervised graph learning. Thus, we derive the Graph Information Plane as defined in the following Def. E.1.

**Definition E.1** (Graph Information Plane). For any node representations $\boldsymbol{H}$ of a graph, it can be seen as encoded from the adjacency $A$ and decoded into the prediction objective $Y$. The sample complexity of the graph learning model is determined by the encoder MI $I(A; \boldsymbol{H}_A)$, while the generalization error is indicated by the decoder MI $I(Y; \boldsymbol{H})$. And technically, the AUC is utilized for computing $I(A; \boldsymbol{H})$ *w.r.t.* edges in $A$, while the Accuracy is used to measure $I(Y; \boldsymbol{H})$. Here, the node representations $\boldsymbol{H}$ can be uniquely mapped to the plane with coordinates $\big(I(A; \boldsymbol{H}), I(Y; \boldsymbol{H})\big)$.

### E.3. Differentiable similarity estimations

Solving Eq. (3) and Eq. (4) requires to derive tractable objectives. Given two variables, $X \in \mathbb{R}^{N \times D_x}$ and $Y \in \mathbb{R}^{N \times D_y}$, we estimate their similarity $s(X, Y)$ with six following differentiable measurements (Kornblith et al., 2019).

- *Dot Product-Based Similarity (DP).* That is, dot products measure the similarity between samples' features as $\mathrm{DP}(X, Y) = \mathrm{tr}(XX^\top YY^\top) = \|Y^\top X\|_F^2$. Note the $ij$-th element of $XX^\top$ are dot products of feature $x_i$ and $x_j$.

- *Hilbert-Schmidt Independence Criterion (HSIC).* HSIC first takes a nonlinear feature transformation of each variable, and then measures the norm of cross-covariance between these features. The empirical estimation (Gretton et al., 2005) is $\mathrm{HSIC}(K, L) = \frac{1}{(n-1)^2} \mathrm{tr}(KHLH)$. Specifically, $K_{ij} = k(x_i, x_j)$ and $L_{ij} = l(y_i, y_j)$, $k$ and $l$ are kernels, and $H$ is the centering matrix of size $N$.

- *Centered Kernel Alignment (CKA).* CKA (Cortes et al., 2012) is devised based on HSIC. It cooperates with HSIC with normalization to become invariant to isotropic scaling, *i.e.*, $\forall \alpha, \beta \in \mathbb{R}^+, s(X, Y) = s(\alpha X, \beta Y)$. In calculation, $\mathrm{CKA}(K, L) = {\mathrm{HSIC}(K,L)}\big/{\sqrt{\mathrm{HSIC}(K,K)\mathrm{HSIC}(L,L)}}$.

- *Kernel Density Estimation (KDE).* KDE estimates the margin and joint PDF (Probability Density Function) of the data as $K_{\mathbf{H}}(\mathbf{x}) = (2\pi)^{-d/2} |\mathbf{H}|^{-1/2} e^{-\frac{1}{2}\mathbf{x}^{\mathbf{T}} \mathbf{H}^{-1} \mathbf{x}}$. KDE generates the $p_X$ and $p_Y$ along with kernels $k_X$ and $k_Y$, which are then used to compute the joint PDF $p_{XY}$ by simply taking the dot product, based on which the MI $I(X; Y)$ is then computed.

- *The Kullback-Leibler Divergence (KL).* Given two probability distributions $X$ and $Y$, the KL divergence is computed as $\mathrm{KL}(X, Y) = \sum_{z \in \mathcal{X}} X(z) \log({X(z)}/{Y(z)})$.

- *Mean Squared Error (MSE).* MSE calculates the point-wise distance of two distributions to indicate their similarity as $\mathrm{MSE}(X, Y) = \frac{1}{n} \sum_{i=1}^n (X_i - Y_i)^2$.

### E.4. Full algorithm

Recall that the two proposed methods, MC-GRA and MC-GPB, that are illustrated in Fig. 5(a) and Fig. 5(b). Here, regarding these two methods, we respectively elaborate the full algorithms in Alg. 1 and Alg. 2 as follows.

---

**Algorithm 1** Markov Chain-based Graph Reconstruction Attack.

---

**Require:** Target model $f_{\boldsymbol{\theta}^*}$, prior knowledge set $\mathcal{K}$, similarity measurement $s(\cdot)$
 1: initialize the parameterized distribution $\mathbb{P}_{\boldsymbol{\phi}}(\hat{\boldsymbol{A}})$ with parameters $\boldsymbol{\phi}$
 2: collect $\mathcal{S}_A = \{\boldsymbol{Z}_A^i : \boldsymbol{Z}_A^i \in \mathcal{K}\}$
 3: **for** $i = 1 \ldots n$ **do**
 4:    sample an adjacency from the distribution $\hat{\boldsymbol{A}} \sim \mathbb{P}_{\boldsymbol{\phi}}(\hat{\boldsymbol{A}})$
 5:    inject stochasticity as $\tilde{X} = X \oplus X_\epsilon$, $\tilde{\boldsymbol{A}} = \tilde{\boldsymbol{A}} \oplus A_\epsilon$
 6:    obtain $\mathcal{S}_{\hat{\boldsymbol{A}}} = \{\boldsymbol{Z}_{\hat{\boldsymbol{A}}}^i : \boldsymbol{Z}_A^i \in \mathcal{S}_A\}$ by the forward of GRA-chain as $f_{\boldsymbol{\theta}^*}(\tilde{\boldsymbol{A}}, \tilde{X})$
 7:    update $\boldsymbol{\phi}$ by maximizing the MC-GRA objective in Eq. (3) with $\mathcal{S}_{\hat{\boldsymbol{A}}}$ and $\mathcal{S}_A$
 8: **end for**
 9: **return** The optimal recovered adjacency $\hat{\boldsymbol{A}}^*$ that $\hat{\boldsymbol{A}}^* \sim \mathbb{P}_{\boldsymbol{\phi}^*}(\hat{\boldsymbol{A}})$

---

**Algorithm 2** Markov Chain-based Defensive Training Against Graph Reconstruction Attack.

---

**Require:** Graph data $A, X, Y$ and similarity measurement $s(\cdot)$
 1: initialize parameters $\boldsymbol{\theta}$ of the GNN $f_{\boldsymbol{\theta}}$
 2: **for** $i = 1 \ldots n$ **do**
 3:    inject stochasticity as $\tilde{A} = A \oplus A_\epsilon$ by randomly dropping edges
 4:    obtain the hidden representations in each layer and outputs by forwarding $f_{\boldsymbol{\theta}}(\tilde{A}, X)$
 5:    update $\boldsymbol{\theta}$ by minimizing the MC-GPB objective in Eq. (4)
 6: **end for**
 7: **return** The trained model $f_{\boldsymbol{\theta}^*}$

---

### E.5. Reproduction

The source code is publicly available at: `https://github.com/tmlr-group/MC-GRA`. In addition, we summarize the search space of hyperparameters and optimal cases as follows. The optimal hyperparameters are obtained by random search or grid search (LaValle et al., 2004).

Table 23: The search space of hyper-parameters in MC-GPB.

| component | name | type | range |
|---|---|---|---|
| Privacy constraint | $\beta_p^1$ (GNN layer-1) | float | $(0, 10)$ |
| | $\beta_p^2$ (GNN layer-2) | float | $(0, 10)$ |
| | $\beta_p^3$ (linear layer) | float | $(0, 10)$ |
| Complexity constraint | $\beta_c^1$ (GNN layer-1) | float | $(0, 10)$ |
| | $\beta_c^2$ (GNN layer-2) | float | $(0, 10)$ |
| Similarity measurement | metric $s(\cdot, \cdot)$ | category | DP, HSIC, CKA, KDE, KL, MSE |

Table 24: The optimal hyper-parameters for MC-GPB.

| component | name | Cora | Citeseer | Polblogs | USA | Brazil | AIDS |
|---|---|---|---|---|---|---|---|
| Privacy constraint | $\beta_p^1$ (GNN layer-1) | 1.3 | 0.09 | 3.00 | 6.60 | 1.90 | 2.40 |
| | $\beta_p^2$ (GNN layer-2) | 1.3 | 0.006 | 2.00 | 1.00 | 2.50 | 3.90 |
| | $\beta_p^3$ (linear layer) | 1.7 | 0.01 | 2.00 | 0.50 | 1.00 | 1.30 |
| Complexity constraint | $\beta_c^1$ (GNN layer-1) | 1.4 | 5e-10 | 1.00 | 1.30 | 1.20 | 1.30 |
| | $\beta_c^2$ (GNN layer-2) | 1.5 | 1e-10 | 1.00 | 3.80 | 1.20 | 1.30 |
| Similarity measurement | metric $s(\cdot, \cdot)$ | KDE | KDE | KDE | DP | KDE | KDE |

Table 25: The search space of hyper-parameters in MC-GRA.

| component | name | type | range |
|---|---|---|---|
| Complexity constraint | $\alpha_c$ (information entropy of $\hat{A}$) | float | $(10^{-4}, 10^4)$ |
| Propagation approximation | $\alpha_p$ (between $H_{\hat{A}}$ and $H_A$) | float | $(10^{-4}, 10^4)$ |
| Output approximation | $\alpha_o$ (between $\hat{Y}_{\hat{A}}$ and $\hat{Y}_A$) | float | $(10^{-4}, 10^4)$ |
| Label approximation | $\alpha_s$ (between $\hat{Y}_{\hat{A}}$ and $Y$) | float | $(10^{-4}, 10^4)$ |
| Similarity measurement | metric $s(\cdot, \cdot)$ | category | DP, HSIC, CKA, KDE, KL, MSE |

Table 26: The optimal hyper-parameters for MC-GRA.

| dataset | prior $\mathcal{K}$ | $\alpha_c$ | $\alpha_p$ | $\alpha_o$ | $\alpha_s$ | $s(\cdot, \cdot)$ |
|---|---|---|---|---|---|---|
| cora | $\{X, H_A\}$ | 10000 | 1000 | 0 | 0 | MSE |
| | $\{X, Y_A\}$ | 100 | 0 | 0.1 | 0 | KL |
| | $\{X, Y\}$ | 0.01 | 0 | 0 | 1 | CKA |
| | $\{X, H_A, Y_A\}$ | 0.1 | 100 | 100 | 0 | MSE |
| | $\{X, H_A, Y\}$ | 0.0001 | 10 | 0 | 1 | MSE |
| | $\{X, Y_A, Y\}$ | 10 | 0 | 0.01 | 1 | MSE |
| | $\{X, H_A, Y_A, Y\}$ | 10 | 10 | 1000 | 1 | MSE |
| citeseer | $\{X, H_A\}$ | 0.01 | 10 | 0 | 0 | KL |
| | $\{X, Y_A\}$ | 0 | 0 | 1 | 0 | KL |
| | $\{X, Y\}$ | 100 | 0 | 0 | 1 | MSE |
| | $\{X, H_A, Y_A\}$ | 10 | 100 | 0.001 | 0 | MSE |
| | $\{X, H_A, Y\}$ | 0.0001 | 100 | 0 | 1 | KL |
| | $\{X, Y_A, Y\}$ | 1 | 0 | 10000 | 1 | KL |
| | $\{X, H_A, Y_A, Y\}$ | 0.001 | 1000 | 0.001 | 1 | KL |
| polblogs | $\{X, H_A\}$ | 0 | 1000 | 0 | 0 | KL |
| | $\{X, Y_A\}$ | 100 | 0 | 1 | 0 | DP |
| | $\{X, Y\}$ | 1000 | 0 | 0 | 1 | MSE |
| | $\{X, H_A, Y_A\}$ | 0.1 | 1000 | 0 | 0 | MSE |
| | $\{X, H_A, Y\}$ | 100 | 0.001 | 0 | 1 | HSIC |
| | $\{X, Y_A, Y\}$ | 100 | 0 | 0 | 1 | CKA |
| | $\{X, H_A, Y_A, Y\}$ | 10000 | 0.001 | 1000 | 1 | HSIC |
| usair | $\{X, H_A\}$ | 100 | 100 | 0 | 0 | MSE |
| | $\{X, Y_A\}$ | 0.01 | 0 | 0.001 | 0 | MSE |
| | $\{X, Y\}$ | 0.0001 | 0 | 0 | 1 | DP |
| | $\{X, H_A, Y_A\}$ | 0.0001 | 0.01 | 0.0001 | 0 | HSIC |
| | $\{X, H_A, Y\}$ | 1000 | 0.001 | 0 | 1 | HSIC |
| | $\{X, Y_A, Y\}$ | 1000 | 0 | 0.01 | 1 | CKA |
| | $\{X, H_A, Y_A, Y\}$ | 10 | 0.01 | 0.01 | 1 | DP |
| brazil | $\{X, H_A\}$ | 100 | 0.01 | 0 | 0 | KL |
| | $\{X, Y_A\}$ | 1 | 0 | 100 | 0 | MSE |
| | $\{X, Y\}$ | 0.001 | 0 | 0 | 1 | KDE |
| | $\{X, H_A, Y_A\}$ | 0.0001 | 100 | 1000 | 0 | MSE |
| | $\{X, H_A, Y\}$ | 0.0001 | 100 | 0 | 1 | HSIC |
| | $\{X, Y_A, Y\}$ | 100 | 0 | 0.1 | 1 | DP |
| | $\{X, H_A, Y_A, Y\}$ | 0.1 | 0.1 | 100 | 1 | KL |
| AIDS | $\{X, H_A\}$ | 10000 | 1000 | 0 | 0 | KL |
| | $\{X, Y_A\}$ | 1 | 0 | 0.01 | 0 | MSE |
| | $\{X, Y\}$ | 0.0001 | 0 | 0 | 1 | CKA |
| | $\{X, H_A, Y_A\}$ | 0.001 | 1 | 100 | 0 | MSE |
| | $\{X, H_A, Y\}$ | 0.0001 | 0.1 | 0 | 1 | MSE |
| | $\{X, Y_A, Y\}$ | 0 | 0 | 0 | 1 | MSE |
| | $\{X, H_A, Y_A, Y\}$ | 0 | 1000 | 0.0001 | 1 | KL |