
Self-Supervised Representation Learning via Latent Graph Prediction

Yaochen Xie^{*1} Zhao Xu^{*1} Shuiwang Ji¹

Abstract

Self-supervised learning (SSL) of graph neural networks is emerging as a promising way of leveraging unlabeled data. Currently, most methods are based on contrastive learning adapted from the image domain, which requires view generation and a sufficient number of negative samples. In contrast, existing predictive models do not require negative sampling, but lack theoretical guidance on the design of pretext training tasks. In this work, we propose the *LaGraph*, a theoretically grounded predictive SSL framework based on latent graph prediction. Learning objectives of *LaGraph* are derived as self-supervised upper bounds to objectives for predicting unobserved latent graphs. In addition to its improved performance, *LaGraph* provides explanations for recent successes of predictive models that include invariance-based objectives. We provide theoretical analysis comparing *LaGraph* to related methods in different domains. Our experimental results demonstrate the superiority of *LaGraph* in performance and the robustness to the decreasing training sample size on both graph-level and node-level tasks.

1. Introduction

Self-supervised learning (SSL) methods seek to use supervisions provided by data itself and design effective pretext learning tasks. These methods allow deep models to learn from a massive amount of unlabeled data and have achieved promising successes in natural language processing (Devlin et al., 2019; Wu et al., 2019; Wang et al., 2019) and image tasks (Batson & Royer, 2019; Xie et al., 2020; He et al., 2020; Chen et al., 2020). To use unlabeled graph data, earlier studies (Perozzi et al., 2014a; Grover & Leskovec, 2016) adapt sequence-based SSL methods (Mikolov et al.,

2013b;a) to learn node representations. Inspired by the recent success of SSL in the image domain, a variety of SSL methods based on graph neural networks (GNNs) have been proposed in different learning paradigms. In particular, recent studies (Veličković et al., 2019; Zhu et al., 2020; Thakoor et al., 2021; Hassani & Khasahmadi, 2020; You et al., 2020) construct SSL tasks as unsupervised approaches to learn representations from graph data at either node-level or graph-level; Hu et al. (2020) propose SSL strategies to pre-train GNNs for downstream tasks; and other studies (Jin et al., 2020; Kim & Oh, 2021) employ SSL as auxiliary tasks to boost the performance of main learning tasks.

Common taxonomies in recent survey works (Xie et al., 2022; Liu et al., 2021b) consider two categories of SSL methods to train GNNs; namely, contrastive methods and predictive methods. Contrastive methods employ pair-wise discrimination as their pretext learning tasks. It performs transformations or augmentations to obtain multiple views from a graph and trains GNNs to discriminate between jointly sampled view pairs and independently sampled view pairs. In contrast, predictive methods (Hamilton et al., 2017; Hwang et al., 2020; Rong et al., 2020) train GNNs to predict certain labels obtained from the input graph, such as node reconstruction, connectivity reconstruction, graph statistical properties, and domain knowledge-based targets.

Adapted from the image domain, current state-of-the-art SSL methods for graphs are mostly contrastive. As a drawback, they usually depend on a large training sample size to include a sufficient number of negative samples. With limited computing resources, contrastive methods may not be applicable to large-scale graphs without suffering from performance loss. To address the drawback, BGRL (Thakoor et al., 2021) adapts BYOL (Grill et al., 2020) to the graph domain. BGRL still obtains different views from each given graph, but it eliminates the requirement of negative samples by replacing contrastive objectives with the prediction of offline embedding. BGRL has achieved competitive performance to the contrastive methods. However, unlike contrastive methods grounded by mutual information estimation and maximization, BYOL and BGRL lack theoretical guidance and require implementation measures to prevent collapsing to trivial representations, such as stop gradient, EMA, and normalization layers.

^{*}Equal contribution ¹Department of Computer Science & Engineering, Texas A&M University, College Station, USA. Correspondence to: Yaochen Xie <ethanyxc@tamu.edu>, Shuiwang Ji <sji@tamu.edu>.

In this work, we propose *LaGraph*, a predictive SSL framework for representation learning of graph data, based on self-supervised latent graph prediction. In particular, we describe the notion of the latent graph and introduce the latent graph prediction as a pretext learning task. We adapt the supervised objective of latent graph prediction into a self-supervised setting by deriving its self-supervised upper bounds, according to which we present the learning framework of *LaGraph*. We provide further justifications of *LaGraph* by comparing it with theoretically sound methods in different domains. Our experimental results demonstrate the effectiveness of *LaGraph* on both graph-level and node-level representation learning, where a remarkable performance boost is achieved on a majority of datasets with higher stability to smaller batch sizes or training on subsets of nodes. Our code is available under the DIG library¹ (Liu et al., 2021a).

Relations with Prior Work: Both *LaGraph* and some existing contrastive methods (You et al., 2020; Zhu et al., 2020; Hu et al., 2020) apply node masking. While those contrastive methods use node masking as an augmentation to obtain different views for contrast, *LaGraph* employs it for the computation of the invariance term in its predictive objective. In addition, the objective of BGRL has a similar formulation to the invariance regularization term in our objective. The objectives of *LaGraph* and BGRL are from different grounding and have essential differences in their computing and effects. While the objective of BGRL is designed and engineered as a variant of contrastive methods, the *LaGraph* objectives are derived as a whole from the latent graph prediction. Our derived theorems associated with *LaGraph* objectives can explain the success of BGRL to some extent and provide guidance on better adopting objectives related to the invariance regularization on graphs.

2. Methods

2.1. Notations and Problem Formulation

We consider an undirected graph $G = (V, E)$ with a set of attributed nodes V and a set of edges E . We formulate the graph data as a tuple of matrices (\mathbf{A}, \mathbf{X}) , where $\mathbf{A} \in \mathbb{R}^{|V| \times |V|}$ denotes the adjacency matrix and $\mathbf{X} \in \mathbb{R}^{|V| \times d}$ denotes the node features of dimension d . We employ a graph encoder \mathcal{E} based on graph neural networks (GNNs) to encode each node or graph into a corresponding representation. Namely, we compute the node-level representations or node embedding by $\mathbf{H} = \mathcal{E}(\mathbf{A}, \mathbf{X}) \in \mathbb{R}^{|V| \times q}$ and the graph-level representation or graph embedding by $\mathbf{z} = \mathcal{R}(\mathbf{H}) \in \mathbb{R}^{1 \times q}$, where q denotes the embedding dimension and $\mathcal{R} : \mathbb{R}^{|V| \times q} \rightarrow \mathbb{R}^{1 \times q}$ is a readout function.

¹<https://github.com/divelab/DIG>.

Self-supervised representation learning is employed to train the graph encoder \mathcal{E} on a set of K graphs $\{G_i\}_{i=1}^K$ without labels from downstream tasks. In particular, we seek to design effective pre-text learning tasks, whose labels are obtained by task designation or from given data, to train the graph encoder \mathcal{E} and produce informative representations for downstream tasks. Depending on the pre-text learning tasks, the encoder \mathcal{E} is usually trained together with some prediction head \mathcal{D} for predictive SSL or a discriminator for contrastive SSL.

2.2. Latent Graph Prediction

Our method considers latent graph prediction as a pretext task to train graph neural networks. In this subsection, we introduce the general notion of latent data, followed by its specific definition for graph data, and the construction of the learning task. For any observed data instance \mathbf{x} , we assume that there exists a corresponding latent data $\mathbf{x}_{\mathcal{I}}$, determining the semantic of \mathbf{x} , such that the latent data $\mathbf{x}_{\mathcal{I}}$ is generated from a prior $p(\mathbf{x}_{\mathcal{I}})$ and the observed data instance is further generated from a certain distribution conditioned on the latent data, *i.e.*, $p(\mathbf{x}|\mathbf{x}_{\mathcal{I}})$. The most common case for the pair of observed data and latent data is the noisy data and its clean version.

When it comes to graph data, we consider the case that an observed graph data $G = (\mathbf{A}, \mathbf{X})$ is (noisily) generated from its latent graph $G_{\ell} = (\mathbf{A}, \mathbf{F})$ with the same node set and edge set, where node feature matrices \mathbf{X} and \mathbf{F} for the two graphs have the same dimensionality. We make two assumptions about the graphs without loss of generality. First, we assume that the observed feature vector \mathbf{x}_v of each node v in an observed graph is independently generated from a certain distribution conditioned on the corresponding latent graph. In other words, how \mathbf{x}_v is generated from the latent feature \mathbf{f}_v is not affected by the generation of other observed feature vectors. Second, we assume that the conditional distribution of the observed graph is centered at the latent graph, *i.e.*, $\mathbb{E}[\mathbf{X}|G_{\ell}] = \mathbf{F}$. The above assumptions are natural when we have little knowledge about the generation process and are commonly used in other types of data such as the non-structural and zero-mean noise in images. In cases where the generation processes of different nodes are related or the distribution is not centered at \mathbf{F} , we can still consider the related or biased components into the latent feature and therefore have the assumptions satisfied.

As the latent data usually determine the semantic meaning of observed data, we believe the prediction of the latent graph can provide informative supervision for the learning of both graph-level and node-level representations. We are hence interested in constructing the learning task of latent graph prediction. To perform latent graph prediction, it is straightforward to employ a graph neural network $f :$

$\{0, 1\}^{|\mathcal{V}| \times |\mathcal{V}|} \times \mathbb{R}^{|\mathcal{V}| \times d} \rightarrow \mathbb{R}^{|\mathcal{V}| \times d}$ that takes an observed graph $G = (\mathbf{A}, \mathbf{X})$ as inputs and predicts the feature matrix of its latent graph $G_{\mathcal{I}} = (\mathbf{A}, \mathbf{F})$. When the ground truth of the latent feature matrix \mathbf{F} is known, the learning objective can be designed as

$$f^* = \arg \min_f \mathbb{E} \|f(\mathbf{A}, \mathbf{X}) - \mathbf{F}\|^2. \quad (1)$$

Intuitively, the latent graph prediction can be considered as a generalized task from noisy data reconstruction that predicts the signal from the noisy data with the objective $\arg \min_f \mathbb{E} \|f(\mathbf{x}) - \mathbf{s}\|^2$, where the mapping from the signal to the noisy data $p(\mathbf{x}|\mathbf{s})$ can usually be explicitly modeled and samples of signal (ground truth) can usually be captured. In the data reconstruction case, pairs of (\mathbf{x}, \mathbf{s}) can be therefore directly captured or synthetically generated given a certain noise model $p(\mathbf{x}|\mathbf{s})$. However, when the task is generalized to latent graph prediction, there is a key challenge preventing us from directly applying the prediction task. That is, whereas there are natural supervisions for noisy data reconstruction, the latent graph is not observed and we are unable to explicitly model the mapping from latent graphs to observed graphs, *i.e.*, the conditional distribution $p(G|G_{\mathcal{I}})$.

2.3. Self-Supervised Upper Bounds for Latent Graph Prediction

As discussed in the previous subsection, unlike typical noisy data reconstruction tasks, the latent graph is not observed and $p(G|G_{\mathcal{I}})$ cannot be modeled explicitly. This makes it difficult to construct a direct learning task for latent graph prediction using the objective in Equation (1). We therefore seek to optimize an alternative objective that approximately optimizes the objective in Equation (1) without requiring the distribution $p(G|G_{\mathcal{I}})$, nor features \mathbf{F} of the latent graph. We now introduce the proposed self-supervised objective for latent graph prediction.

We derive our self-supervised objective without involving \mathbf{F} by constructing an upper bound of the objective in Equation (1). Specifically, we let $J \subset \{0, \dots, |\mathcal{V}| - 1\}$ be an arbitrary subset of node indices, J^c denote the complement of set J , and $\mathbf{X}_{J^c} := \mathbb{1}_{J^c} \odot \mathbf{X} + \mathbb{1}_J \odot \mathbf{M}$ be the feature matrix with features of nodes in V_J masked, where \odot denotes element-wise multiplication, $\mathbf{M} \in \mathbb{R}^{|\mathcal{V}| \times d}$ denotes a matrix consisting of independent random noise or zeros as masking values, and $\mathbb{1}_J \in \mathbb{R}^{|\mathcal{V}| \times d}$ denotes an indicator matrix such that $\mathbb{1}_J[i, :] = \mathbf{1}, \forall i \in J$ and $\mathbb{1}_J[i, :] = \mathbf{0}, \forall i \notin J$. We describe the self-supervised upper bound in Theorem 2.1, whose proof is provided in Appendix A.

Theorem 2.1. *Consider a graph $G = (\mathbf{A}, \mathbf{X})$ and its latent graph $G_{\mathcal{I}} = (\mathbf{A}, \mathbf{F})$. We let the variance of any elements in \mathbf{X} be bounded by σ^2 and J be a subset of nodes V in the graph G . For any graph neural network $f : \{0, 1\}^{|\mathcal{V}| \times |\mathcal{V}|} \times$*

$\mathbb{R}^{|\mathcal{V}| \times d} \rightarrow \mathbb{R}^{|\mathcal{V}| \times d}$, we have the following inequality

$$\begin{aligned} & \mathbb{E}_{\mathbf{A}, \mathbf{X}, \mathbf{F}} \left[\|f(\mathbf{A}, \mathbf{X}) - \mathbf{F}\|^2 + \|\mathbf{X} - \mathbf{F}\|^2 \right] \\ & \leq \mathbb{E}_{\mathbf{A}, \mathbf{X}} \|f(\mathbf{A}, \mathbf{X}) - \mathbf{X}\|^2 + \\ & 2\sigma|\mathcal{V}| \mathbb{E}_J \left[\frac{\mathbb{E}_{\mathbf{A}, \mathbf{X}} \|f_J(\mathbf{A}, \mathbf{X}) - f_J(\mathbf{A}, \mathbf{X}_{J^c})\|^2}{|J|} \right]^{1/2}. \end{aligned} \quad (2)$$

Intuitively, the first component in the upper bound derived in Theorem 2.1 measures the reconstruction error on the feature matrix \mathbf{X} of the given observed graph G , enforcing the intermediate representations to be informative. The second component controls how much information is accessible from the input feature of a node v_i when reconstructing the feature of v_i , by encouraging the output of a node to be invariant to the missing of its features in the input graph. We then call the first component a reconstruction term and the second component an invariance regularization term. Note that the invariance regularization is only computed on masked nodes in contrast to the BGRL objective, based on different theoretical grounding and leading to a different effect. A more detailed discussion is provided in Section 3.

In tasks of self-supervised representation learning, we are more interested in graph-level or node-level representations than predicted latent graphs. In these cases, we expect the representations also hold the invariance property held by the final outputs. We, therefore, seek to apply the invariance regularization to the representations, since a regularization applied to the output does not necessarily control the information accessibility of representations produced intermediately in the graph neural network. To do so, we separately consider the encoder \mathcal{E} and decoder \mathcal{D} in the graph neural network f . We introduce certain assumptions to the decoder network \mathcal{D} and the readout function \mathcal{R} , and derive two additional upper bounds for node-level and graph-level representation learning, respectively in the following corollaries. Proofs of the corollaries are provided in Appendix B.

Corollary 2.2. *Let $G = (\mathbf{A}, \mathbf{X})$ be a given graph, $G_{\mathcal{I}} = (\mathbf{A}, \mathbf{F})$ be its latent graph, \mathcal{E} and \mathcal{D} be a graph encoder and a prediction head (decoder) consisting of fully-connected layers. If the prediction head \mathcal{D} is ℓ -Lipschitz continuous with respect to l_2 -norm, we further have the following inequality,*

$$\begin{aligned} & \mathbb{E} \left[\|\mathcal{D}(\mathbf{H}) - \mathbf{F}\|^2 + \|\mathbf{X} - \mathbf{F}\|^2 \right] \leq \mathbb{E} \|\mathcal{D}(\mathbf{H}) - \mathbf{X}\|^2 \\ & + 2\sigma|\mathcal{V}| \ell \mathbb{E}_J \left[\frac{\mathbb{E} \|\mathbf{H}_J - \mathbf{H}'_J\|^2}{|J|} \right]^{1/2}, \end{aligned} \quad (3)$$

where $\mathbf{H} = \mathcal{E}(\mathbf{A}, \mathbf{X})$ and $\mathbf{H}' = \mathcal{E}(\mathbf{A}, \mathbf{X}_{J^c})$ denote the node embedding of the given graph and the masked graph,

respectively, and $\mathbf{H}_J := \mathbf{H}[J, :]$ selects rows with indices in J .

Corollary 2.3. *Let $G = (\mathbf{A}, \mathbf{X})$ be a given graph, $G_{\mathcal{T}} = (\mathbf{A}, \mathbf{F})$ be its hidden latent graph, \mathcal{E} be a graph encoder, \mathcal{R} be a readout function satisfying k -Bilipschitz continuity with respect to l_2 -norm, and \mathcal{D} be a prediction head (decoder). If the prediction head \mathcal{D} is ℓ -Lipschitz continuous with respect to l_2 -norm, we have the following inequality,*

$$\mathbb{E}[\|\mathcal{D}(\mathbf{H}) - \mathbf{F}\|^2 + \|\mathbf{X} - \mathbf{F}\|^2] \leq \mathbb{E}\|\mathcal{D}(\mathbf{H}) - \mathbf{X}\|^2 + 2\sigma|V|k\ell\mathbb{E}_J\left[\frac{\mathbb{E}\|z - z'\|^2}{|J|}\right]^{1/2}, \quad (4)$$

where $z = \mathcal{R}(\mathbf{H})$ and $z' = \mathcal{R}(\mathbf{H}')$ denote the graph-level representations of the given graph and the masked graph, respectively.

We note that the assumptions and restrictions are natural or practically satisfiable. The assumption that the variance of each element in \mathbf{X} is bounded by σ holds when node features are from $\{0, 1\}^d$ or when feature normalization is applied. The ℓ -Lipschitz continuous property is common for neural networks. And the k -Bilipschitz continuity can be satisfied by applying an injective readout function such as global sum pooling, which is commonly used in graph-level tasks.

2.4. The LaGraph Framework

We design our self-supervised learning framework according to upper bounds derived in Corollary 2.2 and Corollary 2.3. To train encoder \mathcal{E} together with decoder \mathcal{D} under self-supervision, we input to the encoder both the given graph (\mathbf{A}, \mathbf{X}) and its variation $(\mathbf{A}, \mathbf{X}_{J^c})$ with a random subset J of node indices for nodes to be masked and obtain node-level representations $\mathbf{H} = \mathcal{E}(\mathbf{A}, \mathbf{X})$ and $\mathbf{H}' = \mathcal{E}(\mathbf{A}, \mathbf{X}_{J^c})$ for the two graphs respectively. The self-supervised losses are computed on input node features, reconstructed node features, and representations, as demonstrated in Figure 1.

In particular, we consider a mini-batch of N graphs $\{(\mathbf{A}_i, \mathbf{X}_i)\}_{i=1}^N$ and their corresponding masked variation $\{(\mathbf{A}_i, \mathbf{X}_{(i, J_i^c)})\}_{i=1}^N$ where J_i denotes the node indices subset for the i -th graph. The self-supervised loss for node-level representation learning follows Corollary 2.2 and is computed as

$$L_{node}(\mathcal{E}, \mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \|\mathcal{D}(\mathbf{A}_i, \mathbf{H}_i) - \mathbf{X}_i\|^2 / |V_i| + \alpha \left[\frac{\sum_i \|\mathbf{1}_{J_i} \odot \mathbf{H}_i - \mathbf{1}_{J_i} \odot \mathbf{H}'_i\|^2}{\sum_i |J_i|} \right]^{1/2}, \quad (5)$$

where α is a hyper-parameter corresponding to the multiplier $2\sigma\ell$ in Corollary 2.2. To fulfill the conditions in Corollary 2.2, we employ fully-connected layers instead of graph convolutional layers in the decoder \mathcal{D} .

Similarly, using the same notations above, the self-supervised loss for graph-level representation learning follows Corollary 2.3 and is computed as

$$L_{graph}(\mathcal{E}, \mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \|\mathcal{D}(\mathbf{A}_i, \mathbf{H}_i) - \mathbf{X}_i\|^2 / |V_i| + \alpha' \left[\frac{\sum_i \|\mathbf{z}_i - \mathbf{z}'_i\|^2}{\sum_i |J_i|} \right]^{1/2}, \quad (6)$$

where $\mathbf{z}_i = \mathcal{R}(\mathbf{H}_i)$ and $\mathbf{z}'_i = \mathcal{R}(\mathbf{H}'_i)$ denote the graph-level representations obtained by applying readout function \mathcal{R} to the node-level representations, respectively, and α' is a hyper-parameter corresponding to the multiplier $2\sigma k\ell$ in Corollary 2.3. To fulfill the conditions in Corollary 2.3, we employ global sum pooling as the readout function \mathcal{R} , where as the decoder \mathcal{D} here can consist of either fully-connected layers or graph convolutional layers.

The pseudo-code for node-level and graph-level objective computations are provided in Algorithm 1 and Algorithm 2, respectively.

Algorithm 1 LaGraph node-level objective

Inputs: A mini-batch of graphs $\{G_1, \dots, G_N\}$, the encoder \mathcal{E} , the prediction head \mathcal{D} , and the hyper-parameter α .
 $\triangleright G_i = (\mathbf{A}_i, \mathbf{X}_i)$

for i in $1, \dots, N$ **do**

 Generate random $J_i \in \{0, 1\}^{|V_i| \times 1}$, $\mathbf{M} \in \mathbb{R}^{|V_i| \times d}$

$\mathbf{X}_{i, J_i^c} \leftarrow \mathbf{1}_{J_i^c} \odot \mathbf{X}_i + \mathbf{1}_{J_i} \odot \mathbf{M}$ \triangleright Randomly mask nodes

$\mathbf{H}_i \leftarrow \mathcal{E}(\mathbf{A}_i, \mathbf{X}_i)$ \triangleright Compute node representations

$\mathbf{H}'_i \leftarrow \mathcal{E}(\mathbf{A}_i, \mathbf{X}_{i, J_i^c})$

$\mathbf{X}_{rec, i} \leftarrow \mathcal{D}(\mathbf{A}_i, \mathbf{H}_i)$ \triangleright Reconstructed node attributes

$\ell_{rec, i} \leftarrow \|\mathbf{X}_{rec, i} - \mathbf{X}_i\|^2 / |V_i|$

$\ell_{inv, i} \leftarrow \|\mathbf{1}_{J_i} \odot \mathbf{H}_i - \mathbf{1}_{J_i} \odot \mathbf{H}'_i\|^2$

end for

$$L(\mathcal{E}, \mathcal{D}; \{G_1, \dots, G_N\}) = \frac{1}{N} \sum_i \ell_{rec, i} + \alpha (\sum_i \ell_{inv, i} / \sum_i |J_i|)^{1/2}$$

3. Theoretical Analysis and Relations with Prior Work

In this section, we further theoretically justify and motivate *LaGraph* by providing comparisons and connections between our method and existing related methods, including denoising autoencoders (Vincent et al., 2010; Wang et al., 2017), information bottleneck principle (Tishby et al., 1999),

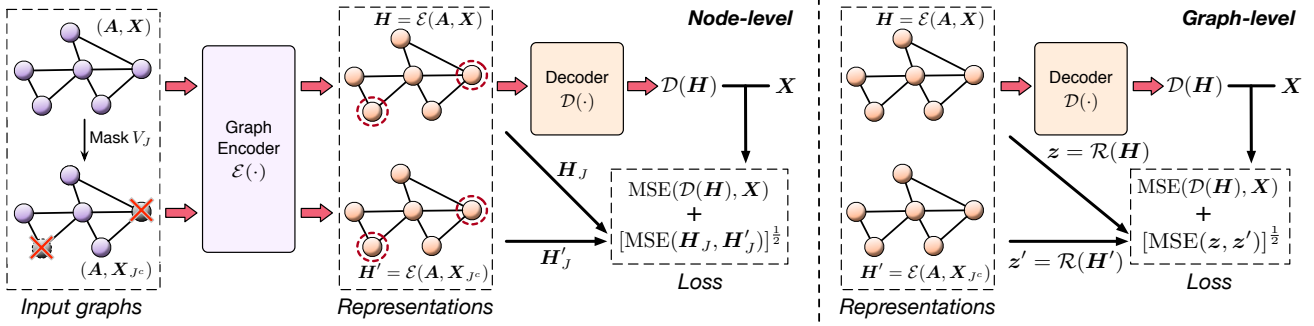


Figure 1. Overview of the *LaGraph* framework. Given a training graph, we randomly mask a small portion $V_J \in V$ of its nodes and input both the original graph and masked graph to the encoder \mathcal{E} . Crossed nodes in the figure have all their attributes masked but topology preserved. The final loss consists of a reconstruction loss on node features and an invariance loss between representations of the original graph and the masked graph. We omit the encoding part of the graph-level framework as frameworks for the two levels mainly differ in whether the invariance term is computed on representations of masked nodes or graph-level representations obtained by \mathcal{R} .

Algorithm 2 *LaGraph* graph-level objective

Inputs: A mini-batch of graphs $\{G_1, \dots, G_N\}$, the encoder \mathcal{E} , the prediction head \mathcal{D} , the readout function \mathcal{R} , and the hyper-parameter α .
 $\triangleright G_i = (A_i, X_i)$

for i in $1, \dots, N$ **do**

Generate random $J_i \in \{0, 1\}^{|V_i| \times 1}$, $M \in \mathbb{R}^{|V_i| \times d}$

$X_{i, J_i^c} \leftarrow \mathbb{1}_{J_i^c} \odot X + \mathbb{1}_{J_i} \odot M$ \triangleright Randomly mask nodes

$H_i \leftarrow \mathcal{E}(A_i, X_i)$ \triangleright Compute node embeddings

$H'_i \leftarrow \mathcal{E}(A_i, X_{i, J_i^c})$

$z_i \leftarrow \mathcal{R}(H_i)$ \triangleright Readout graph representations

$z'_i \leftarrow \mathcal{R}(H'_i)$

$X_{rec} \leftarrow \mathcal{D}(A_i, H_i)$ \triangleright Reconstructed node attributes

$\ell_{rec, i} \leftarrow \|X_{rec, i} - X_i\|^2 / |V_i|$

$\ell_{inv, i} \leftarrow \|z_i - z'_i\|^2$

end for

$$L(\mathcal{E}, \mathcal{D}; \{G_1, \dots, G_N\}) = \frac{1}{N} \sum_i \ell_{rec, i} + \alpha \left(\frac{\sum_i \ell_{inv, i}}{\sum_i |J_i|} \right)^{1/2}$$

and contrastive methods based on local-global mutual information maximization (Veličković et al., 2019; Sun et al., 2019; Hassani & Khasahmadi, 2020). We also discuss the relation and difference to BGRL (Thakoor et al., 2021) and Barlow-Twin (Zbontar et al., 2021).

3.1. Denoising Autoencoders

Denoising autoencoders employ an encoder-decoder network architecture and perform self-supervised training by masking or corrupting a portion of dimensions of the given data and reconstructing the masked or corrupted value given their context. Such an approach has been also applied for self-supervised image denoising (Batson & Royer, 2019), known as blind-spot denoising. Similar to our method, the denoising autoencoder can be also viewed as an approxima-

tion of the latent graph prediction. Using the same notation in Section 2, we formulate the connection between latent graph prediction and the graph denoising autoencoder in the following theorem.

Theorem 3.1. *Let J be a uniformly sampled subset of node indices of the given graph (A, X) , \mathcal{F} be the class of all graph neural networks, and \mathcal{F}^* be the class of graph neural networks such that $f_J^*(A, X)$ does not depend on X_J , for any J and $f^* \in \mathcal{F}^*$. Given any graph neural network $f \in \mathcal{F}$, there exist $f^* \in \mathcal{F}^*$ and $f' \in \mathcal{F}$ such that*

$$\mathbb{E}_{A, X, F} \left[\|f(A, X) - F\|^2 + \|X - F\|^2 \right] \quad (7)$$

$$= \mathbb{E}_{A, X} \|f(A, X) - X\|^2 + \mathbb{E}_{A, X, F} [2\langle f(A, X) - F, X - F \rangle] \quad (8)$$

$$\approx \mathbb{E}_{A, X} \|f^*(A, X) - X\|^2 \quad (9)$$

$$= |V| \mathbb{E}_J \mathbb{E}_{A, X} \|f'_J(A, X_{J^c}) - X_J\|^2 / |J|. \quad (10)$$

Equation (7) is proved in the proof of Theorem 1. It can be verified that the second term, *i.e.*, the expectation of the inner product, in Equation (7) reduces to zero when the neural network f satisfies that $f_J(A, X)$ does not depend on X_J , for any J , according to Batson and Royer (2019). The objective can be therefore approximated by Equation (8) with the neural network f^* satisfying such a property. To let any graph neural network f satisfy the property, one can apply masks to a portion of nodes indexed by J so that their original value is inaccessible by f when predicting $f_J(A, X)$. Therefore, the latent graph prediction objective under supervision can be further approximated by Equation (9), which describes the objective of a graph denoising autoencoder.

A substantial difference between our method and the denoising autoencoder lies in how to handle the inner product term in Equation (7). In particular, the denoising autoencoder

forces the term to be zero by assuming certain properties of the graph neural network, whereas our method derives an upper bound, *i.e.*, the invariance term, for the inner product. Theoretically, the graph denoising autoencoder is equivalent to our framework with an infinite weight scalar for the invariance term. As a drawback, when $f_J(\mathbf{A}, \mathbf{X})$ does not depend on \mathbf{X}_J , the learned representations can be less informative as representations of nodes in V_J do not include the information of \mathbf{X}_J , for any J , leading to performance loss. Our proposed upper bounds allow an encoder to access a certain level of information of the masked nodes, whose representations can be as good as ones from supervised learning. In fact, our method can be viewed as an autoencoder with an invariance regularization.

3.2. The Information Bottleneck Principle

The information bottleneck principle (Tishby et al., 1999) is a technique for data compression and signal processing in the field of information theory, and has been widely applied in deep learning problems (Tishby & Zaslavsky, 2015; Saxe et al., 2018). Let \mathbf{X} be a random variable to be compressed, $\tilde{\mathbf{X}}$ be an observed relevant variable, and \mathbf{Z} denote the compressed representation of \mathbf{X} . The information bottleneck principle seeks to optimize the following problem

$$\mathbf{Z}^* = \arg \min_{\mathbf{Z}} I(\mathbf{Z}; \tilde{\mathbf{X}}) - \beta I(\mathbf{Z}; \mathbf{X}), \quad (11)$$

where $I(\cdot; \cdot)$ denotes the mutual information and $\beta > 1$ is a Lagrange multiplier. The work Barlow Twin (Zbontar et al., 2021) has discussed a connection between the information bottleneck principle and self-supervised learning. In particular, to apply information bottleneck to SSL, one usually obtain $\tilde{\mathbf{X}}$ by performing augmentations or distortions on the given data \mathbf{X} . And Equation (11) can be rewritten into

$$\mathbf{Z}^* = \arg \min_{\mathbf{Z}} [H(\mathbf{Z}) - H(\mathbf{Z}|\tilde{\mathbf{X}})] \quad (12)$$

$$- \beta [H(\mathbf{Z}) - H(\mathbf{Z}|\mathbf{X})] \quad (13)$$

$$= \arg \min_{\mathbf{Z}} H(\mathbf{Z}|\mathbf{X}) - \lambda H(\mathbf{Z}), \quad (14)$$

where $\lambda = \frac{\beta-1}{\beta} > 0$ is a weight scalar. Intuitively, the conditional entropy $H(\mathbf{Z}|\mathbf{X})$ is to be minimized, indicating that the distortion should add no additional information to the representation \mathbf{Z} . In other words, the representation \mathbf{Z} should be as invariant as possible to distortions applied to \mathbf{X} . In addition, the entropy $H(\mathbf{Z})$ is to be maximized, indicating that the representation \mathbf{Z} itself should be as informative as possible.

The two terms in objectives of *LaGraph* correspond to the terms in Equation (14). In particular, the invariance term corresponding to $H(\mathbf{Z}|\mathbf{X})$ and the reconstruction term aims to ensure informative representations, *i.e.*, to maximize $H(\mathbf{Z})$. Objectives in existing SSL methods such as BYOL (Grill

et al., 2020), its variation BGRL (Thakoor et al., 2021) in graph domain, and Barlow Twin (Zbontar et al., 2021) also include invariance terms corresponding to $H(\mathbf{Z}|\mathbf{X})$. To encourage informative representations, Barlow Twin further includes a redundancy reduction term to minimize the cross-correlation between different dimensions of the representation, as a proxy of the maximization of $H(\mathbf{Z})$. In addition, the InfoNCE (NT-XENT) loss employed in some contrastive learning methods (You et al., 2020; Zhu et al., 2020) induces a similar effect, according to Zbontar et al. (2021). Both Equation (14) and the derivation of *LaGraph* objectives indicate the importance of the invariance term in SSL objectives. In addition, compared to the redundancy reduction term in Barlow Twin and the noise contrast in InfoNCE, *LaGraph* objectives can directly guarantee the learning of informative representations measured by the reconstruction capability.

3.3. Contrastive Learning by Maximizing Local-Global Mutual Information

Motivated by Deep InfoMax (Hjelm et al., 2019), recent graph self-supervised learning methods (Veličković et al., 2019; Sun et al., 2019; Hassani & Khasahmadi, 2020) constructs their learning tasks by maximizing the mutual information between local (node-level) representations and a global (graph-level) summary of the graph. Practically, as a k -layer encoder \mathcal{E} has the receptive field of at most k -hop neighborhood, the goal becomes the maximization of the mutual information between local representations and their k -hop neighborhood, formulated as

$$\mathcal{E}^* = \arg \max_{\mathcal{E}} \sum_{i=1}^{|V|} I(\mathbf{X}_i^{(k)}; \mathcal{E}_i(\mathbf{A}, \mathbf{X})), \quad (15)$$

where I denotes the mutual information, $\mathbf{X}_i^{(k)}$ is the k -hop neighborhood of node i , \mathcal{E} is a graph encoder with k GNN layers, and $\mathcal{E}_i(\mathbf{A}, \mathbf{X})$ denotes the local representation of node i . The learning objective is motivated by the goal that the local representations should contain as much the global information of the entire graph (or the k -hop neighborhood) as possible.

As for *LaGraph*, the reconstruction term encourages representations to contain sufficient information to reconstruct the input features while the invariance term limits the information accessibility from a local node when reconstructing its features. The two terms in the objective jointly promote node representations to learn limited local information and as much contextual information from the neighborhood as possible for reconstruction. It hence has a similar effect to the local-global mutual information maximization.

3.4. Other Invariance-Based Objectives

Recent self-supervised learning objectives such as BGRL, Barlow-Twin, and the consistency regularization (Wei et al., 2021) have similar invariance terms as one in the LaGraph objective. Specifically, BGRL minimizes the difference between representations of two augmented views. In spite of the similarity, the invariance terms in LaGraph and other objectives have different grounding and effects.

Regarding how the objectives are computed, the invariance term in the LaGraph objective for node-level representation learning is computed only on masked nodes, in contrast to BGRL and Barlow-Twins objectives where invariance of all nodes are computed. It is worth noting that the proposed objective is an upper bound to the latent graph prediction only if the invariance is computed on the masked nodes, according to the derivation in the proof of Theorem 1. Intuitively, during the computation of a node representation, the invariance term in LaGraph enforces the encoder to capture less information from the node itself and more contextual information. Computing the invariance regularization term on unmasked nodes could lead to a contradicted effect, i.e., discouraging encoders to capture information from contextual nodes, as it lets the representation remain consistent when its masked neighbor nodes are changed. We believe the derivation and the intuition of the proposed objective can provide insights on adopting the invariance regularization into graph self-supervised learning studies.

4. Experiments

We conduct experiments on both node-level and graph-level self-supervised representation learning tasks with datasets used in two most recent state-of-the-art methods for SSL (You et al., 2020; Thakoor et al., 2021). For graph-level tasks, we follow GraphCL (You et al., 2020) to perform evaluations on eight graph classification datasets (Wale & Karypis, 2006; Borgwardt et al., 2005; Dobson & Doig, 2003; Debnath et al., 1991; Yanardag & Vishwanathan, 2015) from TUDataset (Morris et al., 2020). For node-level tasks, as the citation network datasets (McCallum et al., 2000; Giles et al., 1998; Sen et al., 2008) are recognized to be saturated and unreliable for GNN evaluation (Shchur et al., 2018; Thakoor et al., 2021), we follow Thakoor et al. (2021) to include four transductive node classification datasets from Shchur et al. (2018), including Amazon Computers, Amazon Photos from the Amazon Co-purchase Graph (McAuley et al., 2015), Coauthor CS, and Coauthor Physics from the Microsoft Academic Graph (Sinha et al., 2015). We further include three larger-scale inductive datasets, PPI, Reddit, and Flickr, for node-level classification used in SUBG-CON (Jiao et al., 2020).

We follow You et al. (2020) and Zhu et al. (2020) for the

standard linear evaluation protocols at graph-level and node-level, respectively. In particular, for both levels, we first train the graph encoder on unlabeled graph datasets with the corresponding self-supervised objective. We then compute and freeze the corresponding representations and train a linear classification model on top of the fixed representations with their corresponding labels. Linear SVM and the regularized logistic regression are employed as linear classifiers for graph-level datasets and node-level datasets, according to You et al. (2020) and Zhu et al. (2020), respectively. For inductive node-level datasets, the self-supervised training is only performed on graphs in the training datasets whereas the test graphs are unavailable during the self-supervised training.

4.1. Comparisons with Baselines

We perform experiments on both graph-level and node-level datasets to demonstrate the effectiveness of *LaGraph*. We construct our model and losses according to Section 2.4. Detailed model configurations, training settings, and dataset statistics are provided in Appendix C.

Graph-level Datasets. We evaluate the performance of *LaGraph* in terms of the linear classification accuracy and compare it with three kernel-based methods including graphlet kernel (GL) (Shervashidze et al., 2009), Weisfeiler-Lehman kernel (WL) (Shervashidze et al., 2011), and deep graph kernel (DGK) (Yanardag & Vishwanathan, 2015), together with five unsupervised methods including Node2Vec (Grover & Leskovec, 2016), Sub2Vec (Adhikari et al., 2018), Graph2Vec (Narayanan et al., 2017), GAE and VGAE (Kipf & Welling, 2016). We further compare the results with recent SOTA SSL methods based on contrastive learning, including InfoGraph (Sun et al., 2019), MV-GRL (Hassani & Khasahmadi, 2020), and GraphCL (You et al., 2020). Results in Table 1 show that *LaGraph* outperforms the current SOTA methods on a majority of datasets and is on par with the best performance on the rest of datasets. Additional results adopting *LaGraph* as a pre-training strategy under the semi-supervised learning setting are provided in Appendix D.

Node-level Datasets. We perform node-level experiments on both transductive and inductive learning tasks. Transductive self-supervised learning of node representation allows utilization of all data at hand to pre-train GNNs for downstream tasks. Although labels of nodes are not visible during pre-training, patterns and information present in all nodes are observed. In contrast to transductive learning, inductive self-supervised learning only allows using a portion of data to pre-train GNNs, while holding out a certain amount of data for downstream tasks. Our inductive tasks include two cases. First, the PPI dataset consists of 24 graphs, and the

Table 1. Performance on graph-level classification tasks, scores are averaged over 5 runs. Bold and underlined numbers highlight the top-2 performance. OOM indicates running out-of-memory on a 56GB Nvidia A6000 GPU.

	NCII	PROTEINS	DD	MUTAG	COLLAB	RDT-B	RDT-M5K	IMDB-B
GL	–	–	–	81.7±2.1	–	77.3±0.2	41.0±0.2	65.9±1.0
WL	80.0±0.5	72.9±0.6	–	80.7±3.0	–	68.8±0.4	46.1±0.2	72.3±3.4
DGK	80.3±0.5	73.3±0.8	–	87.4±2.7	–	78.0±0.4	41.3±0.2	67.0±0.6
Node2Vec	54.9±1.6	57.5±3.6	75.1±0.5	72.6±10.2	55.7±0.2	73.8±0.5	34.1±0.4	50.0±0.8
Sub2Vec	52.8±1.5	53.0±5.6	73.6±1.5	61.1±15.8	62.1±1.4	71.5±0.4	36.7±0.4	55.3±1.5
Graph2Vec	73.2±1.8	73.3±2.1	76.2±0.1	83.2±9.3	59.9±0.0	75.8±1.0	47.9±0.3	71.1±0.5
GAE	73.3±0.6	74.1±0.5	77.9±0.5	84.0±0.6	56.3±0.1	74.8±0.2	37.6±1.6	52.1±0.2
VGAE	73.7±0.3	74.0±0.5	77.6±0.4	84.4±0.6	56.3±0.0	74.8±0.2	39.1±1.6	52.1±0.2
InfoGraph	76.2±1.1	74.4±0.3	72.9±1.8	89.0±1.1	70.7±1.1	82.5±1.4	53.5±1.0	73.0±0.9
GraphCL	77.9±0.4	74.4±0.5	78.6±0.4	86.8±1.3	71.4±1.2	89.5±0.8	56.0±0.3	71.1±0.4
MVGRL	75.1±0.5	71.5±0.3	OOM	89.7±1.1	OOM	84.5±0.6	OOM	74.2±0.7
LaGraph	79.9±0.5	75.2±0.4	78.1±0.4	90.2±1.1	77.6±0.2	90.4±0.8	56.4±0.4	73.7±0.9

training and testing nodes are split by graphs. In this case, the inductive task is considered across multiple graphs. In other words, node representations are learned from training graphs, and the encoder is evaluated on testing graphs. Second, Flickr and Reddit each consist of only one graph, the training and testing nodes are from the same graph. During self-supervised training, all test nodes are masked-out. During evaluation, all training nodes are masked-out, i.e., test nodes are unseen nodes of the graph during train. For both cases of inductive learning, data used during the self-supervised training stage and data used during evaluation stage are distinct, but the feature dimensionality should be the same for data used in both stages.

For the evaluation of transductive learning, we compare the performance of *LaGraph* in terms of linear classification accuracy with DeepWalk (Perozzi et al., 2014b), GAE, VGAE, and six contrastive learning methods including Deep Graph InfoMax (DGI) (Veličković et al., 2019), GMI (Peng et al., 2020), MVGRL (Hassani & Khasahmadi, 2020), GRACE (Zhu et al., 2020), GCA (Zhu et al., 2021), and BGRL (Thakoor et al., 2021), where BGRL is the current state-of-the-art SSL method for node-level representation learning. We further include the results of directly performing linear classification on raw node features (raw features) and by supervised training for references. To be consistent with Thakoor et al. (2021), we have ensured that the GPU memory consumption of *LaGraph* is under 16GB for the four transductive datasets. We then perform additional experiments on the larger-scale inductive datasets (Zitnik & Leskovec, 2017; Zeng et al., 2020; Hamilton et al., 2017) and compare our results in terms of micro-averaged F1-score with DeepWalk, unsupervised GraphSAGE (Hamilton et al., 2017), DGI, GMI, SUBG-CON (Jiao et al., 2020) and BGRL. Results for both transductive datasets and inductive datasets shown in Table 2. As there is no official BGRL implementation available at the time our experiments are conducted, results with * are obtained from an unoffi-

cial public implementation². Results suggest competitive performance of *LaGraph* compared to the existing SOTA methods. Moreover, *LaGraph* consumes even less memory than BGRL, which requires twice the memory for its GNN encoders for the EMA parameter update.

Experiment Environment Details. We train graph-level datasets on a 11GB GeForce RTX 2080 Ti GPU, and node-level datasets on a 56GB Nvidia RTX A6000 GPU. Our experiments are implemented with PyTorch 1.7.0 and PyTorch Geometric 1.7.0. All neural networks employ batch normalization (Ioffe & Szegedy, 2015), and are optimized with Adam optimizer (Kingma & Ba, 2014). We initialize GNNs with Xavier initialization (Glorot & Bengio, 2010).

4.2. Ablation Study

We further conduct three ablation studies to explore model robustness to smaller batch sizes on graph-level data and to the training with sub-graphs on large-scale node-level datasets. An additional ablation study on the effect of optimizing different objectives is provided in Appendix E.

Robustness to Batch Sizes. Different from contrastive learning methods, *LaGraph* does not require negative samples to perform noise contrast or pair-wise discrimination. Therefore, an advantage of *LaGraph* is that the performance is robust to the batch size as it does not depend on large batch sizes with sufficient negative samples. To verify the statement, we perform an ablation study on how model performance changes when decreasing the batch size from 128 to 8 for graph-level datasets. We include corresponding results of GraphCL which uses InfoNCE for references and show the comparisons in Figure 2. The results indicate while contrastive methods based on InfoNCE suffer from significant performance loss with a small batch size, *LaGraph* are more robust to the batch size.

²https://github.com/namkyeong/bgrrl_pytorch.

Table 2. Performance on node-level datasets, 20 runs averaged. Results of SSL methods with the best performance are highlighted in bold numbers. *Left*: Mean classification accuracy on transductive datasets, with baseline results from Thakoor et al. (2021). *Right*: Micro-averaged F1 scores on larger-scale inductive datasets, with baseline results from Thakoor et al. (2021) and Jiao et al. (2020).

Transductive	Am.Comp.	Am.Pht.	Co.CS	Co.Phy	Inductive	PPI	Flickr	Reddit
Raw features	73.8±0.0	78.5±0.0	90.4±0.0	93.6±0.0	Raw feat.	42.5±0.3	20.3±0.2	58.5±0.1
DeepWalk	85.7±0.1	89.4±0.1	84.6±0.2	91.8±0.2	GAE	75.7±0.0	50.7±0.2	OOM
GAE	87.7±0.3	92.7±0.3	92.4±0.2	95.3±0.1	VGAE	75.8±0.0	50.4±0.2	OOM
VGAE	88.1±0.3	92.8±0.3	92.5±0.2	95.3±0.1	Super-GCN	51.5±0.6	48.7±0.3	93.3±0.1
Supervised	86.5±0.5	92.4±0.2	93.0±0.3	95.7±0.2	Super-GAT	97.3±0.2	OOM	OOM
DGI	84.0±0.5	91.6±0.2	92.2±0.6	94.5±0.5	GraphSAGE	46.5±0.7	36.5±1.0	90.8±1.1
GMI	82.2±0.3	90.7±0.2	OOM	OOM	DGI	63.8±0.2	42.9±0.1	94.0±0.1
MVGRL	87.5±0.1	91.7±0.1	92.1±0.1	95.3±0.0	GMI	65.0±0.0	44.5±0.2	95.0±0.0
GRACE	87.5±0.2	92.2±0.2	92.9±0.0	95.3±0.0	SUBG-CON	66.9±0.2	48.8±0.1	95.2±0.0
GCA	88.9±0.2	92.5±0.2	93.1±0.0	95.7±0.0	BGRL-GCN	69.6±0.2	50.0±0.3*	OOM*
BGRL	89.7±0.3	92.9±0.3	93.2±0.2	95.6±0.1	BGRL-GAT	70.5±0.1	44.2±0.1*	OOM*
LaGraph	88.0±0.3	93.5±0.4	93.3±0.2	95.8±0.1	LaGraph	74.6±0.0	51.3±0.1	95.2±0.0

Table 3. Model performance when trained on a subset of nodes.

	# nodes sampled	100	1,000	2,500	5,000	10,000	all
	% nodes sampled	0.22%	2.24%	5.60%	11.20%	22.41%	100.00%
Flickr	F1-score - <i>LaGraph</i>	6.07	51.12	51.12	51.27	51.29	51.26
	Memory - <i>LaGraph</i>	1389MB	1465MB	1553MB	1725MB	2065MB	4211MB
	F1-score - GraphCL	45.27	45.27	45.27	45.38	45.45	45.48
	Memory - GraphCL	1647MB	2599MB	4137MB	6741MB	11905MB	47939MB
	% nodes sampled	0.07%	0.65%	1.63%	3.25%	6.50%	100.00%
Reddit	F1-score - <i>LaGraph</i>	5.76	95.05	95.06	95.08	95.09	95.22
	Memory - <i>LaGraph</i>	1403MB	1475MB	1585MB	1783MB	2161MB	16933MB
	F1-score - GraphCL	93.24	93.24	93.25	93.31	93.32	OOM
	Memory - GraphCL	4199MB	6117MB	6687MB	9297MB	14495MB	OOM

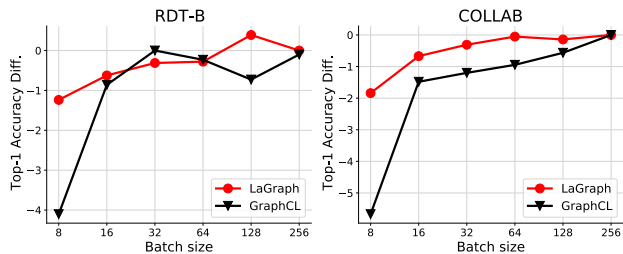


Figure 2. Model robustness to small batch sizes on RDT-B and COLLAB. Shown are relative changes in accuracy over different batch sizes compared to the batch size of 256.

Training on Sub-graphs for Large-scale Datasets.

Training graph encoders on all nodes for some large-scale graphs can be heavily expensive in computation. We hence conduct an ablation study on how training graph encoders on a portion of sampled nodes instead of the entire graph affects the effectiveness of training. Results in Table 3 suggest that the model performance remains stable when decreasing the number of nodes until the number becomes extremely small. The collapse is due to the very sparse connectivity and *LaGraph* fails to reconstruct a node from its neighbor nodes as there are no neighbors at all. In contrast, though GraphCL does not collapse at extremely small subsets, it

suffers more from performance loss above 1,000 nodes and consumes significantly more GPU memory.

5. Conclusions and Future Directions

We introduced *LaGraph*, a state-of-the-art predictive SSL framework whose objectives are based on self-supervised latent graph prediction. We provided theoretical analysis and discussed the relationship between *LaGraph* and theories in different related domains. Experimental results demonstrate the strong effectiveness of the proposed framework and the stability to the training scale for both graph-level and node-level tasks. Currently, our framework mainly considers the latent graph regarding its node features. Further investigation into a latent graph prediction framework that includes richer information such as edge features and latent connectivity into self-supervision can potentially bring additional improvement to the performance. We discuss more future directions in Appendix F.

Acknowledgments

This work was supported in part by National Science Foundation grant IIS-2006861.

References

- Adhikari, B., Zhang, Y., Ramakrishnan, N., and Prakash, B. A. Sub2Vec: Feature learning for subgraphs. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pp. 170–182. Springer, 2018.
- Batson, J. and Royer, L. Noise2Self: Blind denoising by self-supervision. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pp. 524–533, 2019.
- Borgwardt, K. M., Ong, C. S., Schönauer, S., Vishwanathan, S. V. N., Smola, A. J., and Kriegel, H.-P. Protein function prediction via graph kernels. *Bioinformatics*, 21:i47–i56, 06 2005. ISSN 1367-4803.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- Debnath, A. K., Lopez de Compadre, R. L., Debnath, G., Shusterman, A. J., and Hansch, C. Structure-activity relationship of mutagenic aromatic and heteroaromatic nitro compounds correlation with molecular orbital energies and hydrophobicity. *Journal of Medicinal Chemistry*, 34(2):786–797, 02 1991.
- Devlin, J., Chang, M., Lee, K., and Toutanova, K. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 4171–4186, 2019.
- Dobson, P. D. and Doig, A. J. Distinguishing enzyme structures from non-enzymes without alignments. *Journal of Molecular Biology*, 330(4):771–783, 2003. ISSN 0022-2836.
- Giles, C. L., Bollacker, K. D., and Lawrence, S. Cite-seer: An automatic citation indexing system. In *Proceedings of the Third ACM Conference on Digital Libraries*, pp. 89–98. Association for Computing Machinery, 1998. ISBN 0897919653.
- Glorot, X. and Bengio, Y. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9, pp. 249–256, 2010.
- Grill, J.-B., Strub, F., Altché, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., Piot, B., kavukcuoglu, k., Munos, R., and Valko, M. Bootstrap your own latent: A new approach to self-supervised learning. In *Advances in Neural Information Processing Systems*, volume 33, pp. 21271–21284, 2020.
- Grover, A. and Leskovec, J. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 855–864, 2016.
- Hamilton, W., Ying, Z., and Leskovec, J. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, pp. 1024–1034, 2017.
- Hassani, K. and Khasahmadi, A. H. Contrastive multi-view representation learning on graphs. In *International Conference on Machine Learning*, pp. 4116–4126. PMLR, 2020.
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9729–9738, 2020.
- Hjelm, R. D., Fedorov, A., Lavoie-Marchildon, S., Grewal, K., Bachman, P., Trischler, A., and Bengio, Y. Learning deep representations by mutual information estimation and maximization. In *International Conference on Learning Representations*, 2019.
- Hu, W., Liu, B., Gomes, J., Zitnik, M., Liang, P., Pande, V., and Leskovec, J. Strategies for pre-training graph neural networks. In *International Conference on Learning Representations*, 2020.
- Hwang, D., Park, J., Kwon, S., Kim, K., Ha, J.-W., and Kim, H. J. Self-supervised auxiliary learning with meta-paths for heterogeneous graphs. In *Advances in Neural Information Processing Systems*, pp. 10294–10305, 2020.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456. PMLR, 2015.
- Jiao, Y., Xiong, Y., Zhang, J., Zhang, Y., Zhang, T., and Zhu, Y. Sub-graph contrast for scalable self-supervised graph representation learning. In *IEEE International Conference on Data Mining*, pp. 222–231, 2020.
- Jin, W., Derr, T., Liu, H., Wang, Y., Wang, S., Liu, Z., and Tang, J. Self-supervised learning on graphs: Deep insights and new direction. *arXiv preprint arXiv:2006.10141*, 2020.
- Kim, D. and Oh, A. How to find your friendly neighborhood: Graph attention design with self-supervision. In *International Conference on Learning Representations*, 2021.

- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kipf, T. N. and Welling, M. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308*, 2016.
- Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017.
- Laine, S., Karras, T., Lehtinen, J., and Aila, T. High-quality self-supervised deep image denoising. *Advances in Neural Information Processing Systems*, 32:6970–6980, 2019.
- Liu, M., Luo, Y., Wang, L., Xie, Y., Yuan, H., Gui, S., Yu, H., Xu, Z., Zhang, J., Liu, Y., Yan, K., Liu, H., Fu, C., Oztekin, B. M., Zhang, X., and Ji, S. DIG: A turnkey library for diving into graph deep learning research. *Journal of Machine Learning Research*, 22(240):1–9, 2021a. URL <http://jmlr.org/papers/v22/21-0343.html>.
- Liu, X., Zhang, F., Hou, Z., Mian, L., Wang, Z., Zhang, J., and Tang, J. Self-supervised learning: Generative or contrastive. *IEEE Transactions on Knowledge and Data Engineering*, 2021b.
- McAuley, J., Targett, C., Shi, Q., and Van Den Hengel, A. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 43–52, 2015.
- McCallum, A. K., Nigam, K., Rennie, J., and Seymore, K. Automating the construction of internet portals with machine learning. *Information Retrieval*, 3(2):127–163, 2000.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. Efficient estimation of word representations in vector space. In *International Conference on Learning Representations*, 2013a.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G., and Dean, J. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, 2013b.
- Morris, C., Kriege, N. M., Bause, F., Kersting, K., Mutzel, P., and Neumann, M. TUDataset: A collection of benchmark datasets for learning with graphs. In *ICML 2020 Workshop on Graph Representation Learning and Beyond (GRL+ 2020)*, 2020.
- Narayanan, A., Chandramohan, M., Venkatesan, R., Chen, L., Liu, Y., and Jaiswal, S. graph2vec: Learning distributed representations of graphs. *arXiv preprint arXiv:1707.05005*, 2017.
- Peng, Z., Huang, W., Luo, M., Zheng, Q., Rong, Y., Xu, T., and Huang, J. Graph representation learning via graphical mutual information maximization. In *Proceedings of The Web Conference 2020*, pp. 259–270, 2020.
- Perozzi, B., Al-Rfou, R., and Skiena, S. Deepwalk: On-line learning of social representations. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 701–710, 2014a.
- Perozzi, B., Al-Rfou, R., and Skiena, S. DeepWalk: On-line learning of social representations. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 701–710, 2014b.
- Rong, Y., Bian, Y., Xu, T., Xie, W., Wei, Y., Huang, W., and Huang, J. Self-supervised graph transformer on large-scale molecular data. *Advances in Neural Information Processing Systems*, pp. 12559–12571, 2020.
- Saxe, A. M., Bansal, Y., Dapello, J., Advani, M., Kolchinsky, A., Tracey, B. D., and Cox, D. D. On the information bottleneck theory of deep learning. In *International Conference on Learning Representations*, 2018.
- Sen, P., Namata, G., Bilgic, M., Getoor, L., Galligher, B., and Eliassi-Rad, T. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.
- Shchur, O., Mumme, M., Bojchevski, A., and Günnemann, S. Pitfalls of graph neural network evaluation. *Relational Representation Learning Workshop, NeurIPS*, 2018.
- Shervashidze, N., Vishwanathan, S., Petri, T., Mehlhorn, K., and Borgwardt, K. Efficient graphlet kernels for large graph comparison. In *Artificial intelligence and statistics*, pp. 488–495. PMLR, 2009.
- Shervashidze, N., Schweitzer, P., Van Leeuwen, E. J., Mehlhorn, K., and Borgwardt, K. M. Weisfeiler-Lehman graph kernels. *Journal of Machine Learning Research*, 12(9), 2011.
- Sinha, A., Shen, Z., Song, Y., Ma, H., Eide, D., Hsu, B.-J., and Wang, K. An overview of microsoft academic service (mas) and applications. In *Proceedings of the 24th International Conference on World Wide Web*, pp. 243–246, 2015.
- Sun, F.-Y., Hoffman, J., Verma, V., and Tang, J. Infograph: Unsupervised and semi-supervised graph-level representation learning via mutual information maximization. In *International Conference on Learning Representations*, 2019.

- Thakoor, S., Tallec, C., Azar, M. G., Munos, R., Veličković, P., and Valko, M. Bootstrapped representation learning on graphs. In *ICLR 2021 Workshop on Geometrical and Topological Representation Learning*, 2021.
- Tishby, N. and Zaslavsky, N. Deep learning and the information bottleneck principle. In *2015 IEEE Information Theory Workshop (ITW)*, pp. 1–5. IEEE, 2015.
- Tishby, N., Pereira, F. C., and Bialek, W. The information bottleneck method. In *The 37th annual Allerton Conference on Communication, Control, and Computing*, pp. 368–377, 1999.
- Ulyanov, D., Vedaldi, A., and Lempitsky, V. Deep image prior. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 9446–9454, 2018.
- Veličković, P., Fedus, W., Hamilton, W. L., Liò, P., Bengio, Y., and Hjelm, D. Deep graph infomax. In *International Conference on Learning Representations*, 2019.
- Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., Manzagol, P.-A., and Bottou, L. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of Machine Learning Research*, 11(12), 2010.
- Wale, N. and Karypis, G. Comparison of descriptor spaces for chemical compound retrieval and classification. In *Sixth International Conference on Data Mining*, pp. 678–689, 2006.
- Wang, C., Pan, S., Long, G., Zhu, X., and Jiang, J. MGAE: Marginalized graph autoencoder for graph clustering. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, pp. 889–898, 2017.
- Wang, H., Wang, X., Xiong, W., Yu, M., Guo, X., Chang, S., and Wang, W. Y. Self-supervised learning for contextualized extractive summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2221–2227, 2019.
- Wei, C., Shen, K., Chen, Y., and Ma, T. Theoretical analysis of self-training with deep networks on unlabeled data. In *International Conference on Learning Representations*, 2021.
- Wu, J., Wang, X., and Wang, W. Y. Self-supervised dialogue learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 3857–3867, 2019.
- Xie, J., Kelley, S., and Szymanski, B. K. Overlapping community detection in networks: The state-of-the-art and comparative study. *ACM Computing Surveys (CSUR)*, 45(4):1–35, 2013.
- Xie, Y., Wang, Z., and Ji, S. Noise2Same: Optimizing a self-supervised bound for image denoising. In *Proceedings of the 34th Annual Conference on Neural Information Processing Systems*, pp. 20320–20330, 2020.
- Xie, Y., Xu, Z., Zhang, J., Wang, Z., and Ji, S. Self-supervised learning of graph neural networks: A unified review. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- Xu, K., Hu, W., Leskovec, J., and Jegelka, S. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
- Yanardag, P. and Vishwanathan, S. Deep graph kernels. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1365–1374, 2015.
- You, Y., Chen, T., Sui, Y., Chen, T., Wang, Z., and Shen, Y. Graph contrastive learning with augmentations. In *Advances in Neural Information Processing Systems*, volume 33, pp. 5812–5823, 2020.
- Zbontar, J., Jing, L., Misra, I., LeCun, Y., and Deny, S. Barlow twins: Self-supervised learning via redundancy reduction. In *International Conference on Machine Learning*, pp. 12310–12320. PMLR, 2021.
- Zeng, H., Zhou, H., Srivastava, A., Kannan, R., and Prasanna, V. GraphSAINT: Graph sampling based inductive learning method. In *International Conference on Learning Representations*, 2020.
- Zhu, Y., Xu, Y., Yu, F., Liu, Q., Wu, S., and Wang, L. Deep graph contrastive representation learning. In *ICML Workshop on Graph Representation Learning and Beyond*, 2020.
- Zhu, Y., Xu, Y., Yu, F., Liu, Q., Wu, S., and Wang, L. Graph contrastive learning with adaptive augmentation. In *Proceedings of the Web Conference 2021*, pp. 2069–2080, 2021.
- Zitnik, M. and Leskovec, J. Predicting multicellular function through multi-layer tissue networks. *Bioinformatics*, 33(14):i190–i198, 2017.

A. Proof of Theorem 1

Proof. We first derive the relationship between the supervised objective of latent graph prediction $\mathbb{E} \|f(\mathbf{A}, \mathbf{X}) - \mathbf{F}\|^2$ and the self-supervised reconstruction loss $\mathbb{E} \|f(\mathbf{A}, \mathbf{X}) - \mathbf{X}\|^2$ in the following equations,

$$\mathbb{E} \|f(\mathbf{A}, \mathbf{X}) - \mathbf{X}\|^2 = \mathbb{E} \|(f(\mathbf{A}, \mathbf{X}) - \mathbf{F}) - (\mathbf{X} - \mathbf{F})\|^2 \quad (16)$$

$$= \mathbb{E} [\|f(\mathbf{A}, \mathbf{X}) - \mathbf{F}\|^2 + \|\mathbf{X} - \mathbf{F}\|^2 - 2\langle f(\mathbf{A}, \mathbf{X}) - \mathbf{F}, \mathbf{X} - \mathbf{F} \rangle], \quad (17)$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product along all dimensions $\{0, \dots, d|V| - 1\}$. The expectation $\mathbb{E} \|f(\mathbf{A}, \mathbf{X})\|$ in the above equation is not relevant to the neural network f . It hence can be considered as a constant during the optimization of f . To derive an upper bound to $\mathbb{E} \|f(\mathbf{A}, \mathbf{X}) - \mathbf{F}\|^2$, we only need to derive an upper bound of its equivalent $\mathbb{E} \|f(\mathbf{A}, \mathbf{X}) - \mathbf{X}\|^2 + 2\mathbb{E} \langle f(\mathbf{A}, \mathbf{X}) - \mathbf{F}, \mathbf{X} - \mathbf{F} \rangle$. As \mathbf{F} is unobserved, our goal is to derive an upper bound to eliminate the need of \mathbf{F} for the inner product term $\langle f(\mathbf{A}, \mathbf{X}) - \mathbf{F}, \mathbf{X} - \mathbf{F} \rangle$. To do so, we apply the definition of latent graph $\mathbb{E}[\mathbf{X}|\mathbf{A}, \mathbf{F}] = \mathbf{F}$ and rewrite the inner product into the following form.

$$\mathbb{E} \langle f(\mathbf{A}, \mathbf{X}) - \mathbf{F}, \mathbf{X} - \mathbf{F} \rangle = \mathbb{E}_{\mathbf{A}, \mathbf{F}} \mathbb{E}_{\mathbf{X}} \left[\sum_i (f_i(\mathbf{A}, \mathbf{X}) - \mathbf{F}_i)(\mathbf{X}_i - \mathbf{F}_i) | \mathbf{A}, \mathbf{F} \right] \quad (18)$$

$$= \sum_i \mathbb{E}_{\mathbf{A}, \mathbf{F}} \left[\mathbb{E}[(f_i(\mathbf{A}, \mathbf{X}) - \mathbf{F}_i)(\mathbf{X}_i - \mathbf{F}_i) | \mathbf{A}, \mathbf{F}] - \mathbb{E}[f_i(\mathbf{A}, \mathbf{X}) - \mathbf{F}_i | \mathbf{F}] \mathbb{E}[\mathbf{X}_i - \mathbf{F}_i | \mathbf{A}, \mathbf{F}] \right] \quad (19)$$

$$= \sum_i \mathbb{E}_{\mathbf{A}, \mathbf{F}} [\text{Cov}(f_i(\mathbf{A}, \mathbf{X}) - \mathbf{F}_i, \mathbf{X}_i - \mathbf{F}_i | \mathbf{A}, \mathbf{F})] \quad (20)$$

$$= \sum_i \mathbb{E}_{\mathbf{A}, \mathbf{F}} [\text{Cov}(f_i(\mathbf{A}, \mathbf{X}), \mathbf{X}_i | \mathbf{A}, \mathbf{F})], \quad (21)$$

where i sums over all dimensions $\{0, \dots, d|V| - 1\}$, f_i and \mathbf{X}_i denotes the i -th element of the flattened matrices. Note that we employ $\mathbb{E}[\mathbf{X} - \mathbf{F} | \mathbf{A}, \mathbf{F}] = \mathbf{0}$ to let Equation (17) hold, according to the definition of latent graphs. Letting J be a uniformly sampled subset of all node indices $\{0, \dots, |V| - 1\}$, the right hand side of the above equation satisfies

$$\text{RHS} = \mathbb{E}_J \frac{|V|}{|J|} \sum_{j \in J} \sum_{k=0}^{d-1} \mathbb{E}_{\mathbf{A}, \mathbf{F}} [\text{Cov}(f_{jd+k}(\mathbf{A}, \mathbf{X}), \mathbf{X}_{jd+k} | \mathbf{A}, \mathbf{F})], \quad (22)$$

where $f_{jd+k} \in \mathbb{R}$ and $\mathbf{X}_{jd+k} \in \mathbb{R}$ denote the $(jd+k)$ -th element of corresponding matrices, *i.e.*, the k -th element of the node v_j , whereas $\mathbf{X}_J \in \mathbb{R}^{|J| \times d}$ denotes the feature matrix of nodes in V_J . Given the bounded variance $\text{Var}(\mathbf{X}_i) \leq \sigma^2, \forall i$, we bound the above term as

$$\text{RHS} = \mathbb{E}_J \frac{|V|}{|J|} \sum_{j \in J, k} \mathbb{E}_{\mathbf{A}, \mathbf{F}} [\text{Cov}(f_{jd+k}(\mathbf{A}, \mathbf{X}) - f_{jd+k}(\mathbf{A}, \mathbf{X}_{J^c}), \mathbf{X}_{jd+k} | \mathbf{A}, \mathbf{F})] \quad (23)$$

$$\leq \mathbb{E}_J \frac{|V|}{|J|} \sum_{j \in J, k} \mathbb{E}_{\mathbf{A}, \mathbf{F}} [\text{Var}(f_{jd+k}(\mathbf{A}, \mathbf{X}) - f_{jd+k}(\mathbf{A}, \mathbf{X}_{J^c}) | \mathbf{A}, \mathbf{F}) \cdot \text{Var}(\mathbf{X}_{jd+k})]^{1/2} \quad (24)$$

$$\leq |V| \mathbb{E}_J \left(\frac{1}{|J|} \sum_{j \in J, k} \mathbb{E}_{\mathbf{A}, \mathbf{F}} [\text{Var}(f_{jd+k}(\mathbf{A}, \mathbf{X}) - f_{jd+k}(\mathbf{A}, \mathbf{X}_{J^c}) | \mathbf{A}, \mathbf{F}) \cdot \sigma^2] \right)^{1/2} \quad (25)$$

$$\leq \sigma |V| \mathbb{E}_J \left(\frac{1}{|J|} \sum_{j \in J, k} \mathbb{E}_{\mathbf{A}, \mathbf{F}} \left[\mathbb{E} \left[[f_{jd+k}(\mathbf{A}, \mathbf{X}) - f_{jd+k}(\mathbf{A}, \mathbf{X}_{J^c})]^2 | \mathbf{A}, \mathbf{F} \right] \right] \right)^{1/2} \quad (26)$$

$$= \sigma |V| \mathbb{E}_J \left(\frac{1}{|J|} \sum_{j \in J, k} \mathbb{E} \left[f_{jd+k}(\mathbf{A}, \mathbf{X}) - f_{jd+k}(\mathbf{A}, \mathbf{X}_{J^c}) \right]^2 \right)^{1/2} \quad (27)$$

$$= \sigma |V| \mathbb{E}_J \left(\frac{1}{|J|} \mathbb{E} \|f_J(\mathbf{A}, \mathbf{X}) - f_J(\mathbf{A}, \mathbf{X}_{J^c})\|^2 \right)^{1/2}. \quad (28)$$

Above inequalities and equations are derived based on the fact that $f_J(\mathbf{A}, \mathbf{X}_{J^c})$ does not correlate to \mathbf{X}_{j+d+k} as $j \notin J^c$ for Equation (21), the Cauchy-Schwarz inequality for Inequality (22), and $(\mathbb{E}X)^2 \leq \mathbb{E}X^2$ for Inequality (23). We complete the proof of Theorem 1 by combining Equation (15) and Inequality (26),

$$\mathbb{E}[\|f(\mathbf{A}, \mathbf{X}) - \mathbf{F}\|^2 + \|\mathbf{X} - \mathbf{F}\|^2] = \mathbb{E}\|f(\mathbf{A}, \mathbf{X}) - \mathbf{X}\|^2 + 2\langle f(\mathbf{A}, \mathbf{X}) - \mathbf{F}, \mathbf{X} - \mathbf{F} \rangle \quad (29)$$

$$\leq \mathbb{E}\|f(\mathbf{A}, \mathbf{X}) - \mathbf{X}\|^2 + 2\sigma|V|\mathbb{E}_J \left(\frac{1}{|J|} \mathbb{E}\|f_J(\mathbf{A}, \mathbf{X}) - f_J(\mathbf{A}, \mathbf{X}_{J^c})\|^2 \right)^{1/2}. \quad (30)$$

□

B. Proof of Corollary 1 and 2

Proof. We first prove Corollary 1. Consider an ℓ -Lipschitz continuous prediction head with respect to l_2 -norm consists of fully connected layers. We have

$$\|f_J(\mathbf{A}, \mathbf{X}) - f_J(\mathbf{A}, \mathbf{X}_{J^c})\|_2 = \|\mathcal{D}(\mathbf{H}_J) - \mathcal{D}(\mathbf{H}'_J)\|_2 \leq \ell \|\mathbf{H}_J - \mathbf{H}'_J\|_2. \quad (31)$$

We therefore have the following inequality

$$\mathbb{E}\|f_J(\mathbf{A}, \mathbf{X}) - f_J(\mathbf{A}, \mathbf{X}_{J^c})\|_2^2 \leq \mathbb{E}\left[\ell^2 \|\mathbf{H}_J - \mathbf{H}'_J\|_2^2\right]. \quad (32)$$

We apply the above inequality to Theorem 1 and obtain the following inequality

$$\begin{aligned} & \mathbb{E}[\|f(\mathbf{A}, \mathbf{X}) - \mathbf{F}\|^2 + \|\mathbf{X} - \mathbf{F}\|^2] \\ & \leq \mathbb{E}\|f(\mathbf{A}, \mathbf{X}) - \mathbf{X}\|^2 + 2\sigma|V|\mathbb{E}_J \left(\frac{1}{|J|} \mathbb{E}\|f_J(\mathbf{A}, \mathbf{X}) - f_J(\mathbf{A}, \mathbf{X}_{J^c})\|^2 \right)^{1/2} \end{aligned} \quad (33)$$

$$\leq \mathbb{E}\|f(\mathbf{A}, \mathbf{X}) - \mathbf{X}\|^2 + 2\sigma|V|\mathbb{E}_J \left(\frac{1}{|J|} \mathbb{E}\left[\ell^2 \|\mathbf{H}_J - \mathbf{H}'_J\|_2^2\right] \right)^{1/2} \quad (34)$$

$$= \mathbb{E}\|f(\mathbf{A}, \mathbf{X}) - \mathbf{X}\|^2 + 2\sigma|V|\ell\mathbb{E}_J \left(\mathbb{E}\|\mathbf{H}_J - \mathbf{H}'_J\|_2^2 / |J| \right)^{1/2}, \quad (35)$$

which completes the proof of Corollary 1.

Similarly, for Corollary 2, we have

$$\|f_J(\mathbf{A}, \mathbf{X}) - f_J(\mathbf{A}, \mathbf{X}_{J^c})\|_2 = \|\mathcal{D}(\mathbf{H}_J) - \mathcal{D}(\mathbf{H}'_J)\|_2 \leq \ell \|\mathbf{H}_J - \mathbf{H}'_J\|_2. \quad (36)$$

Given an ℓ_r -Bilipschitz continuous readout function \mathcal{R} , the following inequalities hold,

$$\frac{1}{\ell_r} \|\mathbf{H}_J - \mathbf{H}'_J\|_2 \leq \|\mathcal{R}(\mathbf{H}_J) - \mathcal{R}(\mathbf{H}'_J)\|_2 \leq \ell_r \|\mathbf{H}_J - \mathbf{H}'_J\|_2. \quad (37)$$

We therefore have

$$\begin{aligned} & \mathbb{E}[\|f(\mathbf{A}, \mathbf{X}) - \mathbf{F}\|^2 + \|\mathbf{X} - \mathbf{F}\|^2] \\ & \leq \mathbb{E}\|f(\mathbf{A}, \mathbf{X}) - \mathbf{X}\|^2 + 2\sigma|V|\ell\mathbb{E}_J \left(\mathbb{E}\|\mathbf{H}_J - \mathbf{H}'_J\|_2^2 / |J| \right)^{1/2} \end{aligned} \quad (38)$$

$$\leq \mathbb{E}\|f(\mathbf{A}, \mathbf{X}) - \mathbf{X}\|^2 + 2\sigma|V|\ell\ell_r\mathbb{E}_J \left(\mathbb{E}\|\mathcal{R}(\mathbf{H}_J) - \mathcal{R}(\mathbf{H}'_J)\|_2^2 / |J| \right)^{1/2} \quad (39)$$

$$= \mathbb{E}\|f(\mathbf{A}, \mathbf{X}) - \mathbf{X}\|^2 + 2\sigma|V|k\ell\mathbb{E}_J \left(\mathbb{E}\|\mathbf{z} - \mathbf{z}'\|_2^2 / |J| \right)^{1/2}, \quad (40)$$

which completes the proof of Corollary 2. □

C. Experiment Settings and Model Configurations

Dataset Statistics. Statistics including number of graphs, averaged number of nodes, averaged number of edges, and node attribute dimensions are summarized in Table 4.

Table 4. Summary and statistics of common graph datasets for self-supervised learning.

Datasets	Evaluation task	# graphs	Avg. nodes	Avg. edges	# features
NCI1		4110	29.87	32.30	37
PROTEINS		1178	39.06	72.82	3
DD		188	284.32	715.66	89
MUTAG	Graph-level	1113	17.93	19.79	7
COLLAB	classification	5000	74.49	2457.78	1
RDT-B		2000	429.63	497.75	1
RDT-M5K		4999	508.52	594.87	1
IMDB-B		1000	19.77	96.53	1
Amazon Computer	Transductive	1	13,752	245,861	767
Amazon Photos	Node-level	1	7,650	119,081	745
Coauthor CS	classification	1	81,894	81,894	6,806
Coauthor Physics		1	247,962	247,962	8,415
PPI	Inductive	24	2,373	34,133	50
Flickr	Node-level	1	89,250	899,756	500
Reddit	classification	1	232,965	11,606,919	602

Models for Graph-Level Datasets. We employ a 3-layer GIN (Xu et al., 2019) as the graph encoder \mathcal{E} , and a 2-layer MLP as the decoder \mathcal{D} . Following GraphCL (You et al., 2020), we use a hidden dimension of size 32 and concatenate the embedding at each encoding layer to obtain the final representation. To fulfill the conditions in Corollary 2, we apply global sum pooling as the readout function \mathcal{R} . The obtained graph representation is then taken by a SVM classifier with a 10-fold evaluation. For graph datasets that do not come with node attributes, we apply the one-hot vector of the degree for each node as the node attributes so that the node degrees are reconstructed. Certain thresholds for max degrees are applied to reduce computational cost and avoid over sparse node features. The neural network is trained using the loss described in Equation (6). We mask all attributes of the sampled nodes with Gaussian noise. Detailed training configurations including mask ratio, the standard deviation of noise, weight scalar α' , and threshold for max degrees are shown in Table 5. Note that we do not include carefully designed implementation mechanisms by BGRL, such as stop gradients, EMA, and batch normalization at the last layer.

Models for Node-Level Datasets. For node-level datasets, we employ a 2-layer GCN (Kipf & Welling, 2017) as the graph encoder \mathcal{E} , and a linear layer or an MLP as the decoder \mathcal{D} . We use a hidden dimension of size 512 at each encoding layer. The neural network is trained using the loss described in Equation (3). We uniformly employ the weight scalar α' of 2 as we observed that the model performance is not sensitive to the selection of α' within the range $[1, 100]$. We obtain the final node representation by concatenating the original feature with the embedding from the last layer of the encoder. The intuition of this is based on the Bayesian rule where the learned encoder provides the prior knowledge (Ulyanov et al., 2018) of data distribution whereas the given graph data serves as the observed samples. And the posteriori should be based on a combination of the priori (encoder output) and the observed data itself (Laine et al., 2019). Node representation is then taken by a logistic regression classifier that is trained using the cross-entropy (CE) loss with a learning rate of 0.01. Detailed training configurations including mask ratio, the standard deviation of noise, number of encoder and decoder layers, learning rate and weight decay of the graph neural network, training epochs, and weight decay of the logistic regression classifier are shown in Table 6. To split train, valid and test sets, we use the public split used in (Shchur et al., 2018) for Coauthor and Amazon, (Zitnik & Leskovec, 2017; Hamilton et al., 2017; Zeng et al., 2020) for PPI, Reddit and Flickr provided by PyTorch Geometric. Note that we do not include implementation mechanisms by BGRL, such as stop gradients, EMA, and batch normalization at the last layer.

Table 6. Model configurations for node-level datasets.

	Am.Computers	Am.Photos	CoauthorCS	CoauthorPhy	PPI	Flickr	Reddit
Mask ratio	0.05	0.05	0.05	0.05	0.05	0.01	0.05
Noise SD	0.5	0.5	0.005	0.5	0.5	0.5	0.5
Decoder layers	1	1	1	2	2	2	2
Learning rate	10^{-4}	10^{-5}	10^{-3}	10^{-3}	10^{-3}	10^{-4}	10^{-3}
Weight decay	0	10^{-4}	0	0	10^{-5}	0	0
LogReg epochs	400	400	400	300	200	200	500
LogReg WD	10^{-3}	10^{-3}	10^{-3}	10^{-3}	0	10^{-3}	10^{-3}

D. Experimental Results under Semi-supervised Setting

For graph-level datasets, we perform semi-supervised experiments with 10% label rate using both GIN and GCN. All experiments are conducted with the same random seed to avoid randomness in data split and initialization. Under the setting of random initialization followed by supervised learning, the GNN is randomly initialized without pre-training. Under the setting of LaGraph followed by supervised learning, the GNN is pre-trained with the proposed LaGraph framework. Weights of GNNs are fine-tuned during supervised learning with 10% labels. For each dataset, the learning rate and epoch number for pre-training are the same as what we use under an unsupervised setting. For fine-tuning, learning rate is selected from $\{10^{-3}, 10^{-4}, 10^{-5}\}$, and epoch number is selected from $\{5, 10, 15, 20\}$. The results shown in Table 7 and Table 8 indicate that our proposed LaGraph framework is also effective for semi-supervised learning with different GNN backbones.

Table 7. GIN results for Semi-supervised learning.

	NCI1	PROTEINS	DD	MUTAG	COLLAB	RDT-B	RDT-M5K	IMDB-B
Rand. Init. + 10% supervised	76.67	75.29	76.66	86.67	77.54	85.45	56.03	72.70
LaGraph + 10% supervised	80.19	76.10	77.93	91.40	78.04	89.65	56.43	74.30

Table 8. GCN results for Semi-supervised learning.

	NCI1	PROTEINS	DD	MUTAG	COLLAB	RDT-B	RDT-M5K	IMDB-B
Rand. Init. + 10% supervised	75.47	74.48	77.33	84.59	79.02	85.30	53.67	72.90
LaGraph + 10% supervised	78.18	76.28	78.86	85.12	80.12	90.35	55.33	75.10

E. Additional Ablation Studies.

Ablation on Optimizing Different Objectives. We empirically compare the effect of different upper bounds on graph-level datasets. In addition to the objectives described in Corollary 2, we further train the graph encoder with the upper bound described in Theorem 1, which applies invariance regularization on the reconstructed node features. In addition, as node attributes in many graph-level datasets are formed as one-hot vectors of the node type, we also provide the results of using two corresponding multinomial versions of the objective. In particular, we replace the reconstruction term by the cross-entropy between $f(\mathbf{A}, \mathbf{X})$ and \mathbf{X} and, if computed on the outputs, the invariance term by the KL-divergence between $f_J(\mathbf{A}, \mathbf{X})$ and $f_J(\mathbf{A}, \mathbf{X}_{J^c})$. Note that the multinomial versions are no longer strictly upper bounds of supervised latent graph prediction. In Table 9, we show the results obtained under the four objectives above, namely, to compute invariance on on-embedding (MSE-Embed), on-output (MSE-Output), and their corresponding multinomial versions (CE-Embed and

Table 5. Model configurations for graph-level datasets.

	NCI1	PROTEINS	DD	MUTAG	COLLAB	RDT-B	RDT-M5K	IMDB-B
Mask ratio	0.05	0.3	0.1	0.05	0.05	0.05	0.05	0.05
Noise SD	0.5	2	0.5	0.5	0.5	0.5	0.5	0.5
Weight scalar α'	10	1	10	10	10	10	10	10
Degree threshold	–	–	–	–	128	–	–	64
Learning rate	10^{-5}	10^{-5}	10^{-5}	10^{-5}	10^{-4}	10^{-3}	10^{-4}	10^{-4}

CE-Output), respectively. Results indicate that there is no significant difference among the four versions on most datasets, while MSE-Embed and CE-Embed generally tend to be more stable and achieve higher performance on MUTAG, RDT-B, and RDT-M5K.

Table 9. Effect of training with different objectives on graph-level datasets.

	NCI1	PROTEINS	DD	MUTAG	COLLAB	RDT-B	RDT-M5K	IMDB-B
MSE-Embed	79.9±0.5	75.2±0.3	78.1±0.3	90.2±1.3	77.6±0.1	90.4±0.9	56.4±0.2	73.7±0.7
MSE-Output	79.9±0.7	75.0±0.4	78.1±0.8	89.2±2.1	77.7±0.1	89.8±0.9	56.0±0.4	73.4±0.6
CE-Embed	79.9±0.5	75.2±0.3	78.1±0.4	90.1±1.0	77.6±0.2	90.5±1.3	56.3±0.4	73.7±0.7
CE-Output	79.9±0.7	75.2±0.4	78.1±0.4	89.3±2.7	77.6±0.1	89.4±1.8	55.7±0.2	73.5±0.5

Ablation on Concatenated Representations for Node-Level Datasets. For the node-level datasets, We obtain the final node representation by concatenating the original feature with the embedding from the last layer of the encoder, due to the intuition discussed in Appendix C. Results in Table 10 compare the performance of representations with or without concatenations. The removal of the concatenation leads to reduced performance on four of the seven datasets and performance gain on the rest datasets including the most challenging PPI. The results indicate that the concatenation generally positively contributes to the final performance. However, the conclusion still holds that, on node-level datasets, LaGraph can provide significant performance gain on challenging datasets where there is a gap between SSL and supervised performance. Meanwhile, the performance of LaGraph is on par with the performance of supervised learning and the SOTA method BGRL on datasets that are less challenging.

Table 10. Effect of performing concatenation with node features.

Dataset	Am. Comp.	Am. Pht.	Co. CS	Co. Phy	PPI	Flickr	Reddit
With concat	88.0±0.3	93.5±0.4	93.3±0.2	95.8±0.1	74.6±0.0	51.3±0.1	95.2±0.0
W/o concat	88.8±0.3	92.7±0.4	92.6±0.2	95.3±0.1	75.2±0.0	51.6±0.1	94.8±0.0

F. Additional Discussion on Potential Limitations and Future directions

In this section, we discuss several limitations of the proposed method and their solutions or related future directions.

Performance comparison with BGRL on transductive tasks. From the experimental perspective, we admit that BGRL is a quite strong baseline method for transductive tasks. As the results are already on the same level as the performance of supervised training, it is very difficult to further obtain significant improvements. However, when it comes to inductive tasks, where there is still a significant gap between the performance of BGRL and supervised learning, our method is able to bring significant improvements in performance. Therefore, we argue that the non-significant performance boost on some transductive datasets does not degrade our main conclusion about the effectiveness of our method and the contribution of our work.

Unattributed graphs. Although, in this work, we mainly focus on graphs with attributed nodes, there exist cases where the nodes are unattributed and all information is contained in the graph topology, especially for some graph-level datasets. In such cases, we follow a common solution to consider the one-hot vectors of node degrees as the node attributes and our objective performs reconstruction on the node degree. To avoid inconsistency between training and testing graphs in their range of degrees, we introduce thresholds to the node degrees, *i.e.*, the degree of a node is considered as k if it exceeds k . The current solution can capture the topological information of a graph to some degree. However, there can be better solutions capturing full topological information of graphs. A potential direction is to perform connectivity reconstruction with the invariance of representations to the changing in the input edge set. Although the described approach does not currently fit into our theoretical framework, it is possible to derive similar objectives (e.g., upper bound to link prediction objective) following a similar idea.

Scaling-up issue. The scaling-up of graph neural networks becomes an emerging topic. Many existing self-supervised methods may suffer from the scaling-up issue when the graph scales up to billions of nodes and edges. Although we do not perform experimental studies on extremely large graphs, we perform ablation studies to demonstrate the robustness of our method to the training schemes of sampling subgraphs (mini-batches of nodes) for each training iteration.

Performance of SSL methods on unsupervised downstream tasks. Performing linear evaluation with supervised downstream tasks on learned representations is the most common way to evaluate the performance of SSL methods. However, evaluation performance on graph-specific unsupervised tasks such as overlapping community detection (Xie et al., 2013) is seldom studied. Further investigations in the unsupervised downstream task are required to fully demonstrate the effectiveness of self-supervised learning methods.