

Journal Pre-proof

GMNI: Achieve good data augmentation in unsupervised graph contrastive learning

Xin Xiong, Xiangyu Wang, Suorong Yang, Furao Shen, Jian Zhao



PII: S0893-6080(24)00728-7
DOI: <https://doi.org/10.1016/j.neunet.2024.106804>
Reference: NN 106804

To appear in: *Neural Networks*

Received date: 23 May 2024
Revised date: 22 September 2024
Accepted date: 10 October 2024

Please cite this article as: X. Xiong, X. Wang, S. Yang et al., GMNI: Achieve good data augmentation in unsupervised graph contrastive learning. *Neural Networks* (2024), doi: <https://doi.org/10.1016/j.neunet.2024.106804>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2024 Published by Elsevier Ltd.

GMNI: Achieve good data augmentation in unsupervised graph contrastive learning

Xin Xiong^{a,b}, Xiangyu Wang^{a,b}, Suorong Yang^{a,c}, Furao Shen^{a,b,*}, Jian Zhao^d

^aState Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, 210023, China

^bSchool of Artificial Intelligence, Nanjing University, Nanjing, 210023, China

^cDepartment of Computer Science and Technology, Nanjing University, Nanjing, 210023, China

^dSchool of Electronic Science and Engineering, Nanjing University, Nanjing, 210023, China

Abstract

Graph contrastive learning (GCL) shows excellent potential in unsupervised graph representation learning. Data augmentation (DA), responsible for generating diverse views, plays a vital role in GCL, and its optimal choice heavily depends on the downstream task. However, it is impossible to measure task-relevant information under an unsupervised setting. Therefore, many GCL methods risk insufficient information by failing to preserve essential information necessary for the downstream task or risk encoding redundant information. In this paper, we propose a novel method called Minimal Noteworthy Information for unsupervised Graph contrastive learning (GMNI), featuring automated DA. It achieves good DA by balancing missing and excessive information, **approximating the optimal views in contrastive learning**. We employ an adversarial training strategy to generate views that share minimal noteworthy information (MNI), reducing nuisance information by minimization optimization and ensuring sufficient information by emphasizing noteworthy information. Besides, we introduce randomness based on MNI to augmentation, thereby enhancing view diversity and stabilizing the model against

*Corresponding author

Email addresses: xiongxin@smail.nju.edu.cn (Xin Xiong), xiangyuwang@smail.nju.edu.cn (Xiangyu Wang), sryang@smail.nju.edu.cn (Suorong Yang), frshen@nju.edu.cn (Furao Shen), jianzhao@nju.edu.cn (Jian Zhao)

perturbations. Extensive experiments on unsupervised and semi-supervised learning over 14 datasets demonstrate the superiority of GMNI over GCL methods with automated and manual DA. GMNI achieves up to a 1.64% improvement over the state-of-the-art in unsupervised node classification, up to a 1.97% improvement in unsupervised graph classification, and up to a 3.57% improvement in semi-supervised graph classification.

Keywords: Graph neural network, Graph contrastive learning, Data augmentation

1. Introduction

Computer technology is successfully applied in various real-world scenarios (Devlin, 2018; Zhou et al., 2024; Fan et al., 2024). Most real-world data, such as traffic networks, social networks, and protein structures, are unstructured and can be abstracted as graphs. Consequently, research on graphs has gained significant attention in recent years. However, labeling graphs is a very challenging and laborious task since it generally requires domain knowledge, and graphs usually have numerous nodes with complex relationships. Thus, unsupervised graph representation learning (Ju et al., 2023; Liu et al., 2023; Mo et al., 2022) has gained significant attention recently, which aims to obtain low-dimensional representations of nodes or graphs without label information. These representations can be used for a wide range of downstream tasks, such as node classification (Xiao et al., 2022), graph classification (Sun et al., 2019), and graph clustering (Pan and Kang, 2021). Graph contrastive learning (GCL) (Lin et al., 2023; Chen and Kou, 2023; Gong et al., 2023; You et al., 2021) shows great potential in unsupervised graph representation learning due to its excellent data efficiency and generalization ability. GCL generally includes two sequential modules: view generation and view comparison. View generation generates views by employing data augmentation (DA) on the original graph. View comparison acquires view representations through an encoder and then optimizes the encoder by pulling view representations from the same distribution closer while pushing away view representations from different distributions.

Nevertheless, GCL heavily depends on DA, and an inappropriate augmentation will lead to severe performance loss. Graph data augmentation (GDA) necessitates the consideration of both the complex graph topology and feature information. Various GDA techniques have been proposed, which

28 can be categorized into three categories based on their augmentation modal-
 29 ity (Ding et al., 2022): structure-oriented (Mishra et al., 2020; Velickovic
 30 et al., 2019; Zhu et al., 2021; Page et al., 1999; Kondor and Lafferty, 2002),
 31 feature-oriented (Feng et al., 2019; Yang et al., 2021; You et al., 2020a), and
 32 label-oriented (You et al., 2020b; Park et al., 2022). Ensuring the selection of
 33 an appropriate GDA technique is critical in GCL (You et al., 2020a), prompt-
 34 ing many GCL methods to rely on trial-and-error or empirical approaches
 35 when choosing GDAs. GraphCL (You et al., 2020a) introduces four GDA
 36 techniques, and demonstrates that good GDA relies on specific characteris-
 37 tics of various graph datasets through pairwise combinations of these GDAs.
 38 MERIT (Jin et al., 2021) achieves GDA by superimposing the methods from
 39 its GDA pool. GCA (Zhu et al., 2021) conducts GDA by employing three
 40 centralities and selecting the optimal one according to its performance on the
 41 downstream task. On the one hand, these methods require a prudent design
 42 of the GDA pool. On the other hand, selecting the best GDA according to
 43 performance on the downstream task of various datasets is entailed, which
 44 is both time-consuming and computationally expensive. Thus, automated
 45 GDA is of significant importance.

46 The mutual information maximization principle (InfoMax) (Linsker, 1988)
 47 is widely used in GCL (Velickovic et al., 2019; Zhu et al., 2020; Thakoor et al.,
 48 2021), which refers to maximizing the agreement (mutual information, MI)
 49 between view representations. Through InfoMax, the contrastive model can
 50 identify the pairs augmented from the same node or graph, thereby learning
 51 the basic topology and feature information on the graph. However, Info-
 52 Max may risk the model learning redundant information irrelevant to the
 53 downstream task, though this information is beneficial for identifying the
 54 pairs. Encoding redundant information results in brittle representations and
 55 may severely degrade the encoder’s performance on the downstream task
 56 (Tschannen et al., 2019). For example, when training an optical character
 57 recognition model, the color information is redundant and harmful to
 58 recognition, and should be removed beforehand. Thus, before implementing
 59 InfoMax, redundant information shared between views should be stripped
 60 out. Now the new question is *what kind of views InfoMax needs*.

61 The InfoMin principle (Tian et al., 2020) in the visual domain indicates
 62 that a good set of views should share the minimal information necessary to
 63 perform well at the downstream task, and we name such information *minimal*
 64 *necessary information*. It points out that information between views should
 65 contain the necessary information needed for the downstream task while ex-

66 cluding nuisance information (Tian et al., 2020). Indeed, GCL methods with
 67 manual DA essentially search in the DA pool for augmentations that yield
 68 views to best satisfy the InfoMin principle. However, many GCL methods
 69 with automated DA fail to follow InfoMin fully. AutoGCL (Yin et al., 2022),
 70 AD-GCL (Suresh et al., 2021), and JOAO (You et al., 2021) adhere to the
 71 min-max optimization framework, seeking views with minimal similarity or
 72 seeking views that are the most challenging. Unfortunately, the views derived
 73 by minimizing the similarity between views or by minimizing the agreement
 74 between view representations are far from optimal. These approaches reduce
 75 nuisance information but do not emphasize the retention of the necessary
 76 information, potentially resulting in excessive information removal and risk-
 77 ing the model suffering from insufficient information. Therefore, applying
 78 InfoMin to acquire optimal views becomes the key: removing redundant
 79 information while maintaining necessary information related to the down-
 80 stream task. However, it is impossible to measure task-relevant information
 81 under an unsupervised setting. Thus, it is hard to define minimal necessary
 82 information in InfoMin. For instance, InfoGCL (Cheng et al., 2022) and LP-
 83 A3 (Yang et al., 2022) rely on task-specific labels Y to guide the acquisition
 84 of views. Limited research has been carried out on achieving optimal views
 85 in InfoMin under an unsupervised setting. Nevertheless, we find minimal
 86 information that is noteworthy for different downstream tasks can approx-
 87 imate minimal necessary information, which avoids requiring task-relevant
 88 information. We name such information *minimal noteworthy information*
 89 (MNI).

90 In this paper, we propose a novel method called Minimal Noteworthy
 91 Information for unsupervised Graph contrastive learning (GMNI). Specifi-
 92 cally, GMNI uses max-min optimization to learn an importance graph. The
 93 set of views derived by the importance graph shares MNI, making these views
 94 approximations of InfoMin’s optimal views. In addition, GMNI introduces
 95 randomness based on MNI to views, thereby enhancing view diversity and
 96 stabilizing the model against perturbations. Finally, an encoder is optimized
 97 by applying InfoMax to view representations. Moreover, GMNI achieves out-
 98 standing performances without using corrupted views, showing that negative
 99 views are unnecessary for GCL.

100 The significant contributions of this paper are summarized as follows:

- 101 • Manual data augmentation GCL methods are both time-consuming and
 102 computationally expensive, while other automated data augmentation

103 **GCL methods risk insufficient information.** In this paper, a novel un-
 104 supervised GCL method with automated data augmentation, GMNI, is
 105 proposed, in which the tedious manual selection of data augmentation
 106 is eliminated **and avoids the risks of both insufficient and redundant**
 107 **information.**

- 108 • **The InfoMin principle defines optimal views in contrastive learning;**
 109 **however, these views depend on downstream tasks, making them unattain-**
 110 **able in unsupervised settings.** In this paper, we achieve good DA in
 111 unsupervised GCL, presenting a more effective approach to approxi-
 112 mating the optimal views of InfoMin under an unsupervised setting. An
 113 adversarial training strategy is employed to generate views that share
 114 minimal noteworthy information, which avoids using task-relevant in-
 115 formation. Applying InfoMax to these views is risk-free, as they reduce
 116 nuisance information and emphasize noteworthy information to ensure
 117 sufficient information.
- 118 • Extensive experiments on node and graph classification **over 14 datasets**
 119 show the effectiveness of our approximately optimal views on different
 120 tasks and the superiority of GMNI over GCL methods with automated
 121 and manual DA.

122 2. Related work

123 Matrix-factorization-based (Ahmed et al., 2013; Cao et al., 2015; Ou et al.,
 124 2016) and random-walk-based methods (Perozzi et al., 2014; Grover and
 125 Leskovec, 2016) are classical approaches for unsupervised graph representa-
 126 tion learning. However, these methods may emphasize topology excessively
 127 while neglecting feature information. Deep unsupervised graph representa-
 128 tion learning (Zhou et al., 2022; Oyallon, 2020) has gradually evolved in
 129 recent years, with GCL emerging as a promising method.

130 InfoMax (Linsker, 1988) is one of the most commonly used principles in
 131 GCL. InfoMax learns the encoder via maximizing the mutual information
 132 between view representations, i.e., $\max_f I(f(V_1); f(V_2))$, where V_1, V_2 denote
 133 the views, and f denotes the encoder. $I(a; b)$ is the mutual information be-
 134 tween a and b . Through InfoMax, the contrastive model can discern pairs
 135 augmented from the same node or graph, facilitating the learning of fun-
 136 damental topology and feature information within the graph. DIM (Hjelm
 137 et al., 2018) learns by maximizing the mutual information between the input

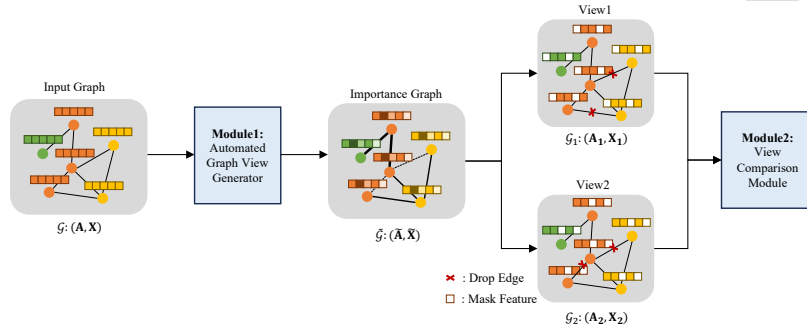


Figure 1: Overview of GMNI, with the automated graph view generator and view comparison modules as the main components.

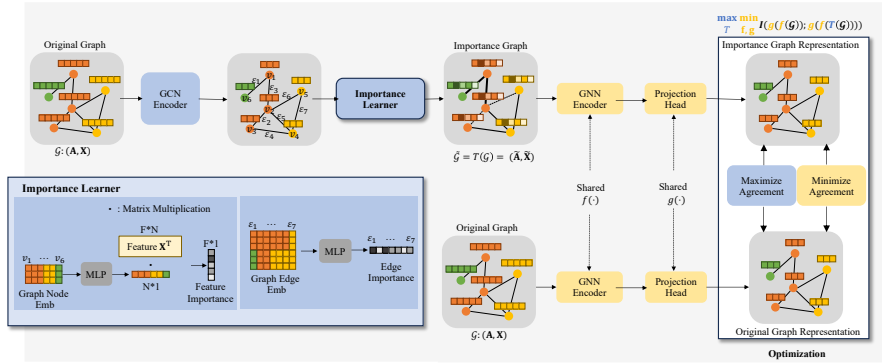


Figure 2: Overview of the automated graph view generator. In the importance graph, darker features and thicker edges are more significant. The original graph is fed into the GCN encoder and importance learner to generate an importance graph. After acquiring representations of the importance and original graph using a shared GNN encoder and projection head, max-min optimization is performed. “Maximize Agreement” optimizes the blue GCN encoder and importance learner, while “Minimize Agreement” optimizes the yellow GNN encoder and projection head.

138 and the output of an encoder. DGI (Velickovic et al., 2019) follows DIM,
 139 applying InfoMax to graph data for the first time. It uses node shuffling to
 140 corrupt the original graph to generate negative pairs, and the positive and
 141 negative pairs are distinguished through a discriminator optimized by Info-
 142 Max. GRACE (Zhu et al., 2020) generates two views through the corruption
 143 of removing edges and masking node features and then learns representations
 144 by applying InfoMax. BGRL (Thakoor et al., 2021) generates views through
 145 stochastic node feature masking and edge masking. It learns by maximizing
 146 the cosine similarity of the view representations at different stages, essen-
 147 tially applying InfoMax. Moreover, BGRL asserts that a corruption method
 148 is unnecessary for contrastive learning (CL). All the above works involve man-
 149 ually selecting DA or corruption method. GCL with manual DA demands
 150 significant time and computational resources to identify an appropriate DA
 151 strategy. Therefore, our work focuses on GCL with automated DA.

152 JOAO (You et al., 2021) searches for the best combination of DAs from
 153 a fixed pool, with the combination coefficients determined by min-max opti-
 154 mization. However, designing a proper DA pool still involves human knowl-
 155 edge, and optimizing only the combination coefficient limits DA’s flexibility.
 156 AD-GCL (Suresh et al., 2021) augments graphs via adversarial training fol-
 157 lowing the information bottleneck (IB) principle (Tishby et al., 2000; Gold-
 158 feld and Polyanskiy, 2020; Alemi et al., 2016), hence reducing redundant
 159 information shared between views. However, AD-GCL only augments on
 160 edges but ignores critical feature information. In InfoGCL (Cheng et al.,
 161 2022), DA is decided mathematically by minimizing mutual information be-
 162 tween views and maximizing mutual information between views and tasks,
 163 which essentially follows IB. InfoGCL relies on task-specific labels Y to guide
 164 the generation of optimal view augmentations, the view encoder, and the
 165 contrastive mode. A manually designed DA pool is entailed in InfoGCL
 166 too, and the best augmentation is acquired by exhaustive search, which is
 167 time-consuming and computationally expensive. In fact, IB and InfoMin are
 168 closely related. IB aims to simultaneously minimize $I(X; Z)$ and maximize
 169 $I(Z; Y)$, where X represents the input, Z is the hidden layer output, and Y
 170 denotes the downstream task information. IB states that the encoder should
 171 minimize the information in the original data while maximizing information
 172 relevant to downstream tasks. AutoGCL (Yin et al., 2022) uses two learn-
 173 able view generators, each of which learns a probability distribution over the
 174 nodes of the input graph. It minimizes the similarity between two views
 175 and maximizes the agreement between the representations of the two views.

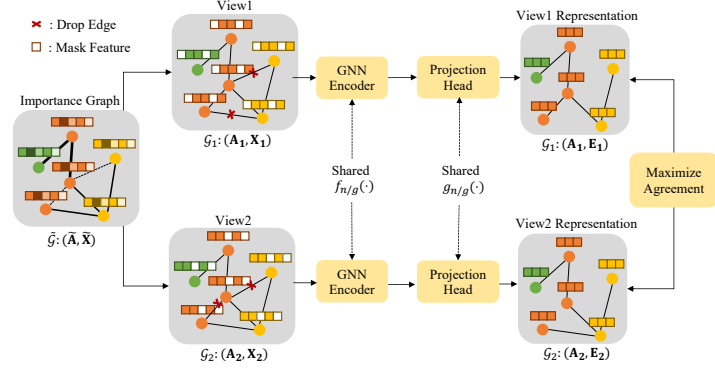


Figure 3: Overview of the view comparison module. Two views are generated by masking features and dropping edges according to the importance graph. Representations of the views are calculated through a shared GNN encoder and projection head. If the downstream task is graph classification, the node representations need to pass through a readout function to get graph representations.

176 Nevertheless, it merely executes node-level view learning without considering
 177 edges.

178 The above GCL methods with automated DA, including JOAO, AD-
 179 GCL, and AutoGCL, adopt a min-max training strategy. They fail to adhere
 180 to InfoMin fully, and their information reduction may lead to a potential
 181 deficiency in essential information. InfoMin, which seeks the optimal views,
 182 states that $(V_1^*, V_2^*) = \operatorname{argmin}_{V_1, V_2} I(V_1; V_2)$, subject to $I(V_1; Y) = I(V_2; Y) =$
 183 $I(X; Y)$. Tian et al. (Tian et al., 2020), the authors of InfoMin, leverage
 184 such a min-max training strategy to complete the unsupervised view learning
 185 of InfoMin too. Given image X , the transformed image $\hat{X} = g(X)$ and
 186 two encoders f_1, f_2 , the objective is: $\min_g \max_{f_1, f_2} I(f_1(g(X)_1); f_2(g(X)_{2:3}))$.
 187 $\{g(X)_1, g(X)_{2:3}\}$ represent the split channels of $g(X)$ and thus serve as the
 188 two views V_1, V_2 . The authors mention that this strategy heavily breaks
 189 constraint $I(V_1; Y) = I(V_2; Y) = I(X; Y)$, and $I(V_1; V_2)$ is overly reduced.
 190 Therefore, a better approximation of InfoMin under an unsupervised setting,
 191 is proposed in this paper.

192 3. Methodology

193 GMNI consists of two sequential modules, the *automated graph view gener-*
 194 *ator* and the *view comparison module*. The automated graph view gener-
 195 ator acquires an importance graph used to generate two views. These two
 196 views are fed into the view comparison module. **Figure 1 illustrates how the**
 197 **two modules cooperate, while Figure 2 and 3 depict their structures.**

198 *Notations.* A graph \mathcal{G} consists of $\mathcal{V} = \{v_1, \dots, v_N\}$, a set of nodes and $\mathcal{E} =$
 199 $\{\varepsilon_1, \dots, \varepsilon_M\}$, a set of edges. The feature matrix, adjacency matrix and degree
 200 matrix of \mathcal{G} are denoted by $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \in \mathbb{R}^{N \times F}$, $\mathbf{A} \in \{0, 1\}^{N \times N}$: $\mathbf{A}_{ij} =$
 201 $\mathbb{I}((v_i, v_j) \in \mathcal{E})$ and \mathbf{D} : $\mathbf{D}_{ii} = \sum_j \mathbf{A}_{ij}$, where F is the feature dimension and
 202 $\mathbf{x}_i \in \mathbb{R}^F$ is the feature of node v_i . We use $\mathcal{G} : (\mathbf{A}, \mathbf{X})$ to represent the graph.
 203 GMNI achieves the approximation of the InfoMin principle without task-
 204 relevant information. To verify the effectiveness of such an approximation
 205 strategy for different downstream tasks, we discuss two downstream tasks,
 206 node classification and graph classification. For node classification, whose
 207 input is a single graph \mathcal{G} , we aim to learn an encoder $f_n : \mathcal{G} \rightarrow \mathbb{R}^{N \times d}$ to obtain
 208 the low-dimensional node embeddings $f_n(\mathcal{G})$, $d \ll F$. For graph classification,
 209 whose input is a set of graphs $\{\mathcal{G}_i\}_{i=1}^Q$, we aim to learn an encoder $f_g : \mathcal{G}_i \rightarrow$
 210 \mathbb{R}^d to obtain the low-dimensional graph embeddings $f_g(\mathcal{G}_i), i = 1, \dots, Q$.

211 3.1. Automated graph view generator

212 **In this section, we present the mathematical motivation for the automated**
 213 **graph view generator and then its implementation details.**

214 **The InfoMin principle indicates that a good set of views in contrastive**
 215 **learning should share the minimal information necessary to perform well at**
 216 **the downstream task. However, applying the InfoMin principle under an un-**
 217 **supervised setting is challenging due to the inability to quantify task-relevant**
 218 **information. Therefore, we propose focusing on views that share the minimal**
 219 **noteworthy information for different downstream tasks as an alternative**
 220 **to the unavailable minimal necessary information. This approach effectively**
 221 **aligns with the InfoMin principle without requiring task-relevant data. Next,**
 222 **we provide the definitions of minimal noteworthy information and the impor-**
 223 **tance graph, with minimal noteworthy information being derived from the**
 224 **importance graph.**

225 **Definition 1** (Importance Graph). **The importance graph $\tilde{\mathcal{G}} = T(\mathcal{G}) =$**
 226 **$(\tilde{\mathbf{A}}, \tilde{\mathbf{X}})$, where $T(\cdot)$ denotes the function that introduces feature and edge im-**

227 *portance to the original graph. The edge and feature importance are repre-*
 228 *sented by $\mathbf{P}_e = \{p_{e,1}, \dots, p_{e,M}\} \in \mathbb{R}^{M \times 1}$ and $\mathbf{P}_f = \{p_{f,1}, \dots, p_{f,F}\} \in \mathbb{R}^{F \times 1}$,*
 229 *respectively. The adjacency matrix $\tilde{\mathbf{A}}$ of $\tilde{\mathcal{G}}$ is formed by replacing the edge*
 230 *values in \mathbf{A} with those from \mathbf{P}_e . Broadcast $\mathbf{P}_f \in \mathbb{R}^{F \times 1}$ to $\mathbf{P}'_f \in \mathbb{R}^{N \times F}$. The*
 231 *feature matrix $\tilde{\mathbf{X}}$ of $\tilde{\mathcal{G}}$ is derived as $\tilde{\mathbf{X}} = \mathbf{X} \odot \mathbf{P}'_f$, where \odot is the Hadamard*
 232 *product.*

233 **Definition 2** (Minimal Noteworthy Information). *Graph information in-*
 234 *cludes both topology and features. Therefore, minimal noteworthy informa-*
 235 *tion (MNI) consists of critical edges and features. For a graph \mathcal{G} and its*
 236 *corresponding importance graph $\tilde{\mathcal{G}}$, MNI refers to the minimal subset of edges*
 237 *and features from \mathcal{G} that exhibit high importance in $\tilde{\mathcal{G}}$, precisely edges with*
 238 *high $p_{e,i}$ and features with high $p_{f,j}$.*

239 3.1.1. Importance graph formulation

240 The importance graph is derived via a max-min optimization,

$$\begin{aligned} & \max_T \min_{f,g} I(g(f(\mathcal{G})); g(f(T(\mathcal{G})))) - |T(\mathcal{G})|, \\ & \text{s. t. } I(g(f(\mathcal{G})); g(f(T(\mathcal{G})))) \geq \zeta, \end{aligned} \quad (1)$$

241 where f is a graph encoder for node embeddings, and g is a projection head
 242 to increase the ability of expression. It should be mentioned that f and g
 243 are not the encoders of the view comparison module. The regularization
 244 term $|T(\mathcal{G})|$ is the normalized sum of edge importance and feature impor-
 245 tance. $I(g(f(\mathcal{G})); g(f(T(\mathcal{G})))) \geq \zeta$ is introduced to prevent the degeneration
 246 of $g(f(\cdot))$ into the corner case, i.e., $I(g(f(\mathcal{G})); g(f(T(\mathcal{G})))) = 0$. Once the
 247 optimization reaches the lower bound ζ , we cease the optimization for both
 248 f and g . However, identifying a corner case of $g(f(\cdot))$ degeneration presents
 249 a challenge for the optimizer. Experiments indicate that the degeneration of
 250 $g(f(\cdot))$ does not occur.

251 **Theorem 1.** *Eqn. (1) is equivalent to the reduction of $I(\mathcal{G}; T(\mathcal{G})|\mathcal{I})$ and the*
 252 *increase of $I(T(\mathcal{G}); \mathcal{I})$, where \mathcal{I} represents the noteworthy information.*

253 **Proof 1.** *Eqn. (1) involves both maximization and minimization; we first*
 254 *introduce the maximization part, i.e.,*

$$\max_T I(g(f(\mathcal{G})); g(f(T(\mathcal{G})))) - |T(\mathcal{G})|. \quad (2)$$

255 $\max_T I(g(f(\mathcal{G})); g(f(T(\mathcal{G}))))$ is achieved by setting $T(\mathcal{G}) = \mathcal{G}$. However,
 256 since $-|T(\mathcal{G})|$ serves as a regularization term, it effectively constrains the
 257 number of emphasized edges and features.

258 Consequently, to maximize the mutual information between the impor-
 259 tance graph $T(\mathcal{G})$ and the original graph \mathcal{G} , it becomes imperative for $T(\mathcal{G})$ to
 260 retain the most vital information within its confined information selectively.

261 *This is a scenario of the greedy algorithm. For clarity, we simplify the*
 262 *problem. Given a set \mathcal{G} with $|\mathcal{G}|$ items, each associated with a value g_i for*
 263 *$i = 1, \dots, |\mathcal{G}|$, the objective is to select a subset $T(\mathcal{G})$ of size $|T(\mathcal{G})|$ that*
 264 *maximizes the total value. The solution is to select the $|T(\mathcal{G})|$ items with the*
 265 *highest values g_i . Suppose we choose a subset $T(\mathcal{G})'$ of size $|T(\mathcal{G})|$ that does*
 266 *not contain the $|T(\mathcal{G})|$ highest values. By replacing items in $T(\mathcal{G})'$ with those*
 267 *from the highest values, the total value can always be increased. Therefore,*
 268 *the subset $T(\mathcal{G})$ that maximizes the total value must consist of the $|T(\mathcal{G})|$*
 269 *highest-value items from \mathcal{G} .*

270 *In other words, optimizing Eqn. (2) boosts the proportion of noteworthy*
 271 *information and reduces that of non-noteworthy information within $T(\mathcal{G})$.*
 272 *This implies an increase in $I(T(\mathcal{G}); \mathcal{I})$ and a decrease in $H(T(\mathcal{G})|\mathcal{I})$, where*
 273 *$H(\cdot)$ represents the information entropy.*

$$\begin{aligned} I(\mathcal{G}; T(\mathcal{G})|\mathcal{I}) &= H(T(\mathcal{G})|\mathcal{I}) - H(T(\mathcal{G})|\mathcal{I}, \mathcal{G}) \\ &= H(T(\mathcal{G})|\mathcal{I}) - 0 \\ &= H(T(\mathcal{G})|\mathcal{I}). \end{aligned} \quad (3)$$

274 $I(\mathcal{G}; T(\mathcal{G})|\mathcal{I}) = H(T(\mathcal{G})|\mathcal{I}) - H(T(\mathcal{G})|\mathcal{I}, \mathcal{G})$ as $I(X; Y) = H(Y) - H(Y|X)$,
 275 and $H(T(\mathcal{G})|\mathcal{I}, \mathcal{G}) = 0$ as $T(\mathcal{G})$ is a function of \mathcal{G} (MacKay, 2003).

276 Therefore, the reduction of $I(\mathcal{G}; T(\mathcal{G})|\mathcal{I})$ and the increase of $I(T(\mathcal{G}); \mathcal{I})$ is
 277 achieved.

278 The aim of $\min_{f,g} I(g(f(\mathcal{G})); g(f(T(\mathcal{G}))))$ is to obtain a challenging in-
 279 formation encoder $g(f(\cdot))$. It is expected that the MNI obtained under a
 280 challenging information encoder $g(f(\cdot))$ will exhibit enhanced robustness. By
 281 treating $g(f(\cdot))$ as a subpar channel with limited data capacity, $T(\mathcal{G})$ must
 282 prioritize transmitting the most crucial information within this limited ca-
 283 pacity. This aligns with the greedy algorithm we previously discussed.

284 Our max-min optimization differs from the min-max optimization used in
 285 methods like AD-GCL. In AD-GCL's min-max optimization, the minimiza-
 286 tion step reduces shared information between views, while the maximization

287 step optimizes the encoder, risking insufficient information in the views. **Con-**
 288 **versely, GMNI preserves noteworthy information in views through maximiza-**
 289 **tion, enhances its robustness through minimization, and reduces redundancy**
 290 **between views via the regularization term.**

291 After deriving the importance graph $\tilde{\mathcal{G}} = T(\mathcal{G})$, MNI can be identified.
 292 MNI represents the minimal subset of edges and features from \mathcal{G} that demon-
 293 strate high importance on $\tilde{\mathcal{G}}$. MNI comprises two components: the minimal
 294 information for $-\max_T |T(\mathcal{G})| = \min_T |T(\mathcal{G})|$ and the noteworthy informa-
 295 tion for $\max_T I(g(f(\mathcal{G})); g(f(T(\mathcal{G}))))$. Those critical edges and features of
 296 $T(\mathcal{G})$ hold universal significance and is highly likely required by various down-
 297 stream tasks. For example, the critical edges connecting different clusters
 298 within a graph play a crucial role in tasks like graph partition (Karypis and
 299 Kumar, 1998; Nazi et al., 2019), graph classification (Kipf and Welling, 2016;
 300 Han et al., 2022), and link prediction (Zhang and Chen, 2018; Yang et al.,
 301 2021). These edges serve as essential connections that provide valuable in-
 302 sights into the underlying relationships and dependencies within the graph.

303 3.1.2. Importance graph construction

304 **This section details the generation of edge and feature importance, leading**
 305 **to the derivation of the importance graph, as shown in the “importance**
 306 **learner” in Figure 2.**

307 First, node representations are computed using a graph convolutional
 308 layer (Welling and Kipf, 2016):

$$309 \mathbf{E} = \sigma(\mathbf{A}\mathbf{X}\mathbf{W}), \quad (4)$$

310 where σ is the activation function, \mathbf{W} is the parameter matrix. \mathbf{A} is the ad-
 311 jacency matrix without normalization, as the node degrees are also essential
 312 for importance computation.

313 We need to get the node importance first to get the feature importance.
 314 The node importance $\mathbf{P}_n \in \mathbb{R}^{N \times 1}$ is given by

$$315 \mathbf{P}_n = \text{Gumbel}(h_\phi(\mathbf{E})),$$

$$\text{Gumbel}(t) = \text{Sigmoid}\left(\frac{\log \eta - \log(1 - \eta) + t}{\tau}\right), \quad \eta \sim \text{Uniform}(0, 1), \quad (5)$$

314 where $h_\phi(\cdot)$ is a simple MLP, and $\text{Gumbel}(t)$ is the Gumbel-Max reparametriza-
 315 tion function (Maddison et al., 2016). τ is the temperature parameter,

316 and the closer τ is to 0, the closer \mathbf{P}_n is to binarization. The Gumbel-
 317 Max reparametrization trick ensures the process can be backpropagated and
 318 makes \mathbf{P}_n probabilistically meaningful. Thus, feature importance $\mathbf{P}_f \in \mathbb{R}^{F \times 1}$
 319 can be calculated through node importance (Zhu et al., 2021), i.e.,

$$\mathbf{P}_f = \mathbf{X}^T \mathbf{P}_n, \quad (6)$$

320 where $\mathbf{X}^T \in \mathbb{R}^{F \times N}$ is the transpose of feature matrix \mathbf{X} . **Combining node**
 321 **features and node importance yields feature importance.** If features are real
 322 numbers between 0 and 1, we perform the Gumbel-Max reparametrization
 323 trick on \mathbf{P}_f rather than on \mathbf{P}_n to avoid \mathbf{P}_f being too small.

324 Inspired by AD-GCL (Suresh et al., 2021), the importance of edge $\varepsilon_k =$
 325 $(v_i, v_j) \in \mathcal{E}$ can be expressed by

$$\mathbf{p}_{e,k} = \text{Gumbel}(h_\psi([\mathbf{E}[v_i]; \mathbf{E}[v_j]])), \quad (7)$$

326 where $[\cdot; \cdot]$ is the concatenation operation, $h_\psi(\cdot)$ is a simple MLP, and $\mathbf{E}[v_i]$
 327 is the embedding of node v_i . The edge importance $\mathbf{P}_e = [\mathbf{p}_{e,1}, \dots, \mathbf{p}_{e,M}]^T \in$
 328 $\mathbb{R}^{M \times 1}$.

329 The importance graph is $\tilde{\mathcal{G}} = (\tilde{\mathbf{A}}, \tilde{\mathbf{X}})$. Broadcast $\mathbf{P}_f \in \mathbb{R}^{F \times 1}$ to $\mathbf{P}'_f \in$
 330 $\mathbb{R}^{N \times F}$. $\tilde{\mathbf{X}}$ is given by $\tilde{\mathbf{X}} = \mathbf{X} \odot \mathbf{P}'_f$. Derive $\tilde{\mathbf{A}}$ by replacing the edge values in
 331 the adjacency matrix \mathbf{A} with the edge values of \mathbf{P}_e .

332 3.2. Importance graph optimization

333 This section instantiates Section 3.1.1 and is illustrated in Figure 2. The
 334 representation of node v_i in \mathcal{G} is defined by $\mathbf{z}_{i,1} = g_\xi(f_\theta(\mathcal{G}(v_i)))$ and the
 335 representation of it in $\tilde{\mathcal{G}}$ by $\mathbf{z}_{i,2} = g_\xi(f_\theta(\tilde{\mathcal{G}}(v_i)))$. For simplicity, f_θ is a
 336 1- or 2- layer graph encoder, and g_ξ is a 1- or 2- layer MLP **servicing as the**
 337 **projection head.**

338 InfoNCE (Oord et al., 2018), a lower bound of mutual information, is ap-
 339 plied to estimate mutual information. Thus, the mutual information between
 340 $\mathbf{z}_{i,1}$ and $\mathbf{z}_{i,2}$ can be estimated by

$$\begin{aligned} \hat{I}(\mathbf{z}_{i,1}; \mathbf{z}_{i,2}) &= \frac{1}{2}(\hat{I}_0(\mathbf{z}_{i,1}; \mathbf{z}_{i,2}) + \hat{I}_0(\mathbf{z}_{i,2}; \mathbf{z}_{i,1})), \\ \hat{I}_0(\mathbf{z}_{i,1}; \mathbf{z}_{i,2}) &= \log \frac{\exp(s(\mathbf{z}_{i,1}, \mathbf{z}_{i,2})/\epsilon)}{\sum_{j=1, j \neq i}^N \exp(s(\mathbf{z}_{i,1}, \mathbf{z}_{j,2})/\epsilon)}, \end{aligned} \quad (8)$$

341 where $s(\cdot, \cdot)$ is the cosine similarity, ϵ is the temperature parameter, and
 342 the symmetrical design considers the same status of \mathcal{G} and $\tilde{\mathcal{G}}$. The mutual
 343 information between the representation of \mathcal{G} and $\tilde{\mathcal{G}} = T(\mathcal{G})$ is estimated by

$$\begin{aligned} I(g_\xi(f_\theta(\mathcal{G})); g_\xi(f_\theta(T(\mathcal{G})))) &\rightarrow \hat{I}(g_\xi(f_\theta(\mathcal{G})); g_\xi(f_\theta(T(\mathcal{G})))) \\ &= \frac{1}{N} \sum_{i=1}^N \hat{I}(\mathbf{z}_{i,1}; \mathbf{z}_{i,2}). \end{aligned} \quad (9)$$

344 The regularization term $|T(\mathcal{G})|$ is defined by

$$|T(\mathcal{G})| = \lambda \left(\sum_{i=1}^F \mathbf{P}_f^{(i)} / F + \sum_{i=1}^M \mathbf{P}_e^{(i)} / M \right), \quad (10)$$

345 where λ is the regularization weight, and $\mathbf{P}_f^{(i)}$ is the i -th dimension of \mathbf{P}_f .
 346 Finally, our optimization strategy to get the importance graph is

$$\begin{aligned} \max_{\mathbf{w}, \phi, \psi} \min_{\theta, \xi} &\hat{I}(g_\xi(f_\theta(\mathcal{G})); g_\xi(f_\theta(T_{\mathbf{w}, \phi, \psi}(\mathcal{G})))) \\ &- \lambda \left(\sum_{i=1}^F \mathbf{P}_f^{(i)} / F + \sum_{i=1}^M \mathbf{P}_e^{(i)} / M \right), \quad (11) \\ \text{s. t. } &\hat{I}(g_\xi(f_\theta(\mathcal{G})); g_\xi(f_\theta(T(\mathcal{G})))) \geq \zeta, \end{aligned}$$

347 which is the instantiation of Eqn. (1).

348 3.2.1. View generation

349 MNI refers to the minimal subset of edges and features from \mathcal{G} that exhibit
 350 high importance in $\tilde{\mathcal{G}}$, specifically edges with high $p_{e,i}$ and features with high
 351 $p_{f,j}$. Views for the comparison module will be generated from MNI, which
 352 contains essential information from the graph. View generation is depicted
 353 in both Figure 1 and Figure 3.

354 CL aims to encourage the model to find consistency across different views.
 355 Employing soft edges and features determined by importance for all views
 356 will constrain view diversity, as these views share graph topology and fea-
 357 tures. Therefore, we introduce randomness based on MNI in views to enhance
 358 diversity and stabilize GMNI against perturbations.

359 Two views are required in CL, $\mathcal{G}_1 : (\mathbf{A}_1, \mathbf{X}_1)$ and $\mathcal{G}_2 : (\mathbf{A}_2, \mathbf{X}_2)$. Take the
 360 generation of \mathcal{G}_1 as an example, and the generation of \mathcal{G}_2 is similar. View
 361 modifies both the topology and features of the input graph. The edge set \mathcal{E}_1

362 of \mathcal{G}_1 is a subset of \mathcal{E} . We use a random variable $p_{\varepsilon_k} \sim \text{Bernoulli}(1 - p_{d,\varepsilon_k})$
 363 to select the edge ε_k , i.e., if $p_{\varepsilon_k} = 1$, then $\varepsilon_k \in \mathcal{E}_1$, else $\varepsilon_k \notin \mathcal{E}_1$. p_{d,ε_k}
 364 represents the drop probability of edge ε_k , which is generated based on the
 365 edge importance $p_{e,k}$. Therefore, \mathbf{A}_1 can be derived via \mathcal{E}_1 . Inspired by
 366 GCA (Zhu et al., 2021), $\mathbf{P}_{d,\varepsilon} = [p_{d,\varepsilon_1}, \dots, p_{d,\varepsilon_M}]$ is given by

$$\mathbf{P}_{d,\varepsilon} = \min \left(\frac{\mathbf{P}_e^{\max} - \mathbf{P}_e}{\mathbf{P}_e^{\max} - \mathbf{P}_e^{\text{avg}}} \cdot p_{s1}, p_t \right), \quad (12)$$

367 where $\min(a, b)$ indicates to select the smaller one of a and b , p_{s1} and p_t are
 368 hyperparameters between 0 and 1, p_{s1} is used to adjust the overall drop rate
 369 of edges, and p_t is used to truncate the drop rate such that every edge has a
 370 chance of being included. \mathbf{P}_e^{\max} , $\mathbf{P}_e^{\text{avg}}$ are the maximum and average value of
 371 \mathbf{P}_e respectively. Thus, edges with higher importance have a greater chance
 372 of being included in \mathcal{E}_1 .

373 The feature matrix of \mathcal{G}_1 can be represented as $\mathbf{X}_1 = \mathbf{X} \odot \mathbf{M}$, where
 374 $\mathbf{M} \in \mathbb{R}^{N \times F}$ is the mask matrix, broadcasted from $\mathbf{m} \in \{0, 1\}^F$. A random
 375 variable $p_{f_k} \sim \text{Bernoulli}(1 - p_{d,f_k})$ is used to select the k -th dimensional
 376 feature, i.e., if $p_{f_k} = 1$, then $\mathbf{m}_k = 1$, else $\mathbf{m}_k = 0$. \mathbf{m}_k represents the k -th
 377 dimensional of \mathbf{m} . p_{d,f_k} represents the probability that the k -th dimensional
 378 feature is masked, which is generated based on the feature importance $p_{f,k}$.
 379 $\mathbf{P}_{d,f} = [p_{d,f_1}, \dots, p_{d,f_F}]$ is given by

$$\mathbf{P}_{d,f} = \min \left(\frac{\mathbf{P}_f^{\max} - \mathbf{P}_f}{\mathbf{P}_f^{\max} - \mathbf{P}_f^{\text{avg}}} \cdot p_{s2}, p_t \right), \quad (13)$$

380 where p_{s2} is a hyperparameter between 0 and 1 to control the overall mask
 381 rate of features. Thus, features with higher importance are less likely to be
 382 masked in \mathbf{X}_1 . The views generated by these random variables effectively
 383 preserve MNI while incorporating randomness.

384 3.3. View comparison module

385 To validate the effectiveness of the views, we employ the simplest view
 386 comparison module, directly pulling view representations closer via InfoMax.
 387 When using views generated through MNI, InfoMax can be applied without
 388 interference from nuisance information or loss of noteworthy information.
 389 The view comparison module is depicted in Figure 3.

390 The node classification task aims to learn an encoder $f_n : \mathcal{G} \rightarrow \mathbb{R}^{N \times d}$. \mathcal{G}_1
 391 and \mathcal{G}_2 are the views of \mathcal{G} , then our objective is

$$\max_{f_n, g_n} I(g_n(f_n(\mathcal{G}_1)); g_n(f_n(\mathcal{G}_2))), \quad (14)$$

392 where f_n is a GCN (Welling and Kipf, 2016), g_n is a simple 1- or 2- layer MLP,
 393 and we use InfoNCE (Eqn. (8)) as the estimator of mutual information.

394 The graph classification task aims to learn an encoder $f_g : \mathcal{G}_i \rightarrow \mathbb{R}^d$.
 395 Given a minibatch of graphs $\{\mathcal{G}_i\}_{i=1}^K$, $\mathcal{G}_{i,1}$ and $\mathcal{G}_{i,2}$ are the views of \mathcal{G}_i , then
 396 the objective is

$$\max_{f_g, g_g} \frac{1}{K} \sum_{i=1}^K I(g_g(f_g(\mathcal{G}_{i,1})); g_g(f_g(\mathcal{G}_{i,2}))), \quad (15)$$

397 where f_g is a GIN (Xu et al., 2018), g_g is a simple 1- or 2- layer MLP
 398 followed by a readout function which is a simple summation. f_n and f_g
 399 can be replaced by any graph encoder. Finally, the node embeddings $f_n(\mathcal{G})$
 400 and graph embeddings $\{f_g(\mathcal{G}_i), i = 1, \dots, Q\}$ are used for downstream testing.
 401 More details of the algorithm and the computational complexity analysis can
 402 be found in Appendix C.

403 4. Experiments

404 4.1. Datasets

405 To demonstrate the superiority of GMNI over other GCL methods with
 406 automated and manual data augmentations, we conduct experiments on 14
 407 datasets, including 5 node classification datasets and 9 graph classification
 408 datasets.

409 The 5 datasets for node classification are Wiki-CS (Mernyei and Cangea,
 410 2020), Amazon-Computers (Shchur et al., 2018), Amazon-Photo (Shchur
 411 et al., 2018), Coauthor-CS (Shchur et al., 2018) and Coauthor-Physics (Shchur
 412 et al., 2018). Their detailed statistics are shown in Table 1. Wiki-CS has
 413 dense real number features, whereas the other datasets have sparse one-hot
 414 features. Following GCA (Zhu et al., 2021), we evaluate models under the
 415 public train, test, and validation sets supplied by Wiki-CS. For the other
 416 four datasets, we randomly split the dataset into three sets: 80% train, 10%
 417 test, and 10% validation.

Table 1: Node classification datasets.

Dataset	#Nodes	#Edges	#Features	#Classes	Category
Wiki-CS	11,701	216,123	300	10	Reference Networks
Amazon-Computers	13,752	245,861	767	10	Co-purchase Networks
Amazon-Photo	7,650	119,081	745	8	Co-purchase Networks
Coauthor-CS	18,333	81,894	6,805	15	Co-authorship Networks
Coauthor-Physics	34,493	247,962	8,415	5	Co-authorship Networks

Table 2: Graph classification datasets.

Dataset	Avg. #Graphs	Avg. #Nodes	Avg. #Edges	#Classes	Category
MUTUG	188	17.93	19.79	2	Biochemical Molecules
PROTEINS	1,113	39.06	72.82	2	Biochemical Molecules
DD	1,178	284.32	715.66	2	Biochemical Molecules
NCI1	4,110	29.87	32.30	2	Biochemical Molecules
COLLAB	5,000	74.49	2457.78	3	Social Networks
GITHUB	12,725	113.79	234.64	2	Social Networks
IMDB-BINARY	1,000	19.77	96.53	2	Social Networks
REDDIT-BINARY	2,000	429.63	497.75	2	Social Networks
REDDIT-MULTI-5K	4,999	508.52	594.87	5	Social Networks

418 We use 9 datasets from TUDataset (Morris et al., 2020) for graph classi-
419 fication, including MUTUG, PROTEINS, DD, NCI1, COLLAB, GITHUB,
420 IMDB-BINARY, REDDIT-BINARY, and REDDIT-MULTI-5K. Their de-
421 tailed statistics are shown in Table 2. Following JOAO (You et al., 2021),
422 we use the entire dataset to learn graph representations and feed them into
423 the downstream classifier using 10-fold cross-validation.

424 4.2. Baselines

425 For node classification, GMNI is compared to (1) GCL methods without
426 data augmentation, including DGI (Velickovic et al., 2019) and GMI (Peng
427 et al., 2020). (2) Manual data augmentation GCL methods, including MV-
428 GRL (Hassani and Khasahmadi, 2020) and GCA (Zhu et al., 2021). (3)
429 Automated data augmentation GCL methods, including AD-GCL (Suresh
430 et al., 2021), JOAOv2 (You et al., 2021), and AutoGCL (Yin et al., 2022).

431 For graph classification, GMNI is compared to (1) GCL method without
432 data augmentation, including InfoGraph (Sun et al., 2019). (2) Manual data
433 augmentation GCL methods, including GraphCL (You et al., 2020a), MV-
434 GRL, GCA, and HGCL (Ju et al., 2023). (3) Automated data augmentation
435 GCL methods, including AD-GCL, JOAOv2, and AutoGCL. (4) A special
436 baseline, SEGA (Wu et al., 2023), is orthogonal to previous works on graph

437 augmentations. It implements the anchor view for GCL.

438 The brief introductions of these baselines are as follows.

- 439 • DGI (Velickovic et al., 2019) learns graph representations by maximiz-
440 ing mutual information between patch representations and correspond-
441 ing high-level graph summaries.
- 442 • GMI (Peng et al., 2020) learns graph representations by maximizing
443 mutual information between the input and output of a graph neural
444 encoder.
- 445 • InfoGraph (Sun et al., 2019) learns graph representations by maxi-
446 mizing mutual information between the graph-level and substructure
447 representations at various scales.
- 448 • GraphCL (You et al., 2020a) designs four types of graph augmentations
449 for view generation and learns graph representations by maximizing
450 mutual information between the view representations.
- 451 • MVGRL (Hassani and Khasahmadi, 2020) learns graph representations
452 through contrasting encodings derived from first-order neighbors and
453 graph diffusion.
- 454 • GCA (Zhu et al., 2021) introduces an adaptive augmentation method
455 for GCL, which integrates diverse priors related to topological and se-
456 mantic aspects of the graph.
- 457 • HGCL (Ju et al., 2023) employs node-level contrastive learning, graph-
458 level contrastive learning, and mutual contrastive learning to capture
459 graph semantics hierarchically. Additionally, the model incorporates a
460 Siamese network and momentum update mechanism.
- 461 • AD-GCL (Suresh et al., 2021) learns by optimizing **an adversarial graph**
462 **augmentation strategy** to avoid capturing redundant information.
- 463 • JOAOv2 (You et al., 2021) proposes a principled bi-level optimization
464 framework to learn the selection of DA methods for GCL.
- 465 • **AutoGCL (Yin et al., 2022) employs an automated data augmentation**
466 **strategy to integrate a collection of trainable graph view generators,**
467 **each learning a probability distribution over the input graph.**

Table 3: Unsupervised learning performance on 5 node classification datasets (average accuracy (%) \pm std. over 5 runs). OOM indicates Out-Of-Memory on a 32GB GPU. The highest and second performances are in **bold** and underlined, respectively.

Type	Model	Wiki-CS	Amaz-Comp	Amaz-Photo	Coauthor-CS	Coauthor-Phy
w/o GDA	DGI	75.35 \pm 0.14	83.95 \pm 0.47	91.61 \pm 0.22	92.15 \pm 0.63	94.51 \pm 0.52
	GMI	74.85 \pm 0.08	82.21 \pm 0.31	90.68 \pm 0.17	OOM	OOM
w/ Manual GDA	MVGRL	77.52 \pm 0.08	87.52 \pm 0.11	91.74 \pm 0.07	92.11 \pm 0.12	95.33 \pm 0.03
	GCA	<u>78.35 \pm 0.05</u>	<u>87.85 \pm 0.31</u>	<u>92.53 \pm 0.16</u>	<u>93.10 \pm 0.01</u>	<u>95.73 \pm 0.03</u>
w/ Automated GDA	AD-GCL	73.46 \pm 0.36	81.32 \pm 0.93	88.75 \pm 0.92	92.16 \pm 0.36	94.57 \pm 0.09
	JOAO-v2	75.36 \pm 0.47	85.96 \pm 0.98	91.15 \pm 0.55	91.33 \pm 0.27	OOM
	AutoGCL	73.66 \pm 0.59	86.44 \pm 1.24	91.98 \pm 0.58	92.26 \pm 0.32	OOM
	GMNI	79.19 \pm 0.13	89.29 \pm 0.05	93.52 \pm 0.33	93.61 \pm 0.15	95.91 \pm 0.10

- **SEGA (Wu et al., 2023) achieves the anchor view for graph contrastive learning by minimizing structural entropy, which preserves essential information from the input graph.**

4.3. Comparison with baselines

This section compares GMNI with baselines in unsupervised and semi-supervised settings for node and graph classification. The details of the experimental setup are provided in Appendix A.

For the baselines, the experimental results are from the original papers. The original papers do not provide results for AD-GCL, JOAOv2, and AutoGCL in unsupervised node classification, nor for GCA in unsupervised graph classification. Therefore, we either remove or add a pooling layer to obtain their node and graph representations, which are subsequently fed into the downstream classifiers. We report the best results for GCA’s 3 variations, GCA-DE, GCA-EVC, and GCA-PR.

Unsupervised learning on node classification.

Settings. We train GMNI and other baselines using unlabeled data to generate representations, then use these representations to train a simple ℓ_2 -regularized logistic regression classifier with a learning rate of 0.01.

Result analysis. Table 3 shows that GMNI achieves outstanding performance on the node classification task. Firstly, it significantly outperforms baselines with automated DA. The reason for the improvement is that the views, defined by retaining minimal noteworthy information, approximate

Table 4: Unsupervised learning performance on graph classification on the benchmark TUDataset (Morris et al., 2020) (average accuracy (%) \pm std. over 5 runs). ‘-’ indicates that results are unavailable in papers. The highest and second performances are in **bold** and underlined, respectively.

Type	Model	MUTAG	PROTEINS	DD	NCI1
w/o GDA	InfoGraph	89.01 \pm 1.13	74.44 \pm 0.31	72.85 \pm 1.78	76.20 \pm 1.06
w/ Manual GDA	GraphCL	86.80 \pm 1.34	74.39 \pm 0.45	78.62 \pm 0.40	77.87 \pm 0.41
	MVGRL	89.70 \pm 1.10	-	-	-
	GCA	<u>90.60 \pm 0.76</u>	75.53 \pm 0.22	79.17 \pm 0.39	75.87 \pm 0.96
	HGCL	90.10 \pm 0.80	75.50 \pm 0.50	<u>79.20 \pm 0.60</u>	-
w/ Automated GDA	AD-GCL	89.30 \pm 1.50	73.59 \pm 0.65	74.49 \pm 0.52	69.67 \pm 0.51
	JOAOv2	87.30 \pm 1.00	71.25 \pm 0.85	66.91 \pm 1.75	72.99 \pm 0.75
	AutoGCL	88.64 \pm 1.08	<u>75.80 \pm 0.36</u>	77.57 \pm 0.60	<u>82.00 \pm 0.29</u>
	GMNI	91.57 \pm 0.57	76.58 \pm 0.29	79.32 \pm 0.26	83.12 \pm 0.20
Type	Model	COLLAB	IMDB-B	REDDIT-B	REDDIT-M-5K
w/o GDA	InfoGraph	70.65 \pm 1.13	73.03 \pm 0.87	82.50 \pm 1.42	53.46 \pm 1.03
w/ Manual GDA	GraphCL	71.36 \pm 1.15	71.14 \pm 0.44	<u>89.53 \pm 0.84</u>	55.99 \pm 0.28
	MVGRL	-	74.20 \pm 0.70	84.50 \pm 0.60	-
	GCA	76.67 \pm 0.43	<u>74.97 \pm 0.40</u>	87.03 \pm 0.90	56.04 \pm 0.28
	HGCL	75.80 \pm 0.40	73.90 \pm 0.70	-	-
w/ Automated GDA	AD-GCL	73.32 \pm 0.61	71.57 \pm 1.01	85.52 \pm 0.79	53.00 \pm 0.82
	JOAOv2	70.40 \pm 2.21	71.60 \pm 0.86	78.35 \pm 1.38	45.57 \pm 2.86
	AutoGCL	70.12 \pm 0.68	73.30 \pm 0.40	88.58 \pm 1.49	<u>56.75 \pm 0.18</u>
	GMNI	<u>76.19 \pm 0.71</u>	75.54 \pm 0.29	91.29 \pm 0.41	57.02 \pm 0.14

491 the optimal views in InfoMin while other baselines fail to. Utilizing such
 492 views can prevent the contrastive module from learning nuisance informa-
 493 tion and make it focuses on noteworthy information, thus achieving better
 494 performance. Secondly, GMNI outperforms manual augmentation baselines
 495 on all datasets. GCL methods with manual augmentation heavily rely on
 496 the predefined DA pool, resulting in a limited approximation of the InfoMin
 497 principle. Considering factors such as time cost, computational cost, and the
 498 ability to approximate optimal views in InfoMin, GMNI emerges as a prefer-
 499 able alternative to GCL with manual GDA. Thus, GMNI liberates GCL from
 500 the tedious manual selection of DA.

501 *Unsupervised learning on graph classification.*

502 **Settings.** The downstream classifier for graph classification is an SVM with
 503 parameter C grid searching in [0.001, 0.01, 0.1, 1, 10, 100, 1000].

504 **Result analysis.** Table 4 shows that GMNI achieves the best performance
 505 on 7 out of 8 datasets and exhibits sub-optimal performance on the remain-

Table 5: Unsupervised learning performance on graph classification (average accuracy (%) \pm std) of SEGA, SEGA in collaboration with SOTAs, and GMNI. The highest and second performances are in **bold** and underlined, respectively.

Model	MUTAG	PROTEINS	DD	NC11
SEGA	90.21 \pm 0.66	76.01 \pm 0.42	78.76 \pm 0.57	79.00 \pm 0.72
SEGA+AD-GCL	89.89 \pm 0.69	74.61 \pm 0.81	75.84 \pm 0.64	70.38 \pm 0.76
SEGA+JOAOv2	88.53 \pm 2.45	75.94 \pm 0.88	78.37 \pm 1.26	78.04 \pm 0.19
SEGA+AutoGCL	89.03 \pm 1.01	<u>76.43 \pm 0.67</u>	78.31 \pm 1.37	<u>81.84 \pm 0.53</u>
GMNI	91.57 \pm 0.57	76.58 \pm 0.29	79.32 \pm 0.26	83.12 \pm 0.20
	COLLAB	IMDB-B	REDDIT-B	REDDIT-M-5K
SEGA	74.12 \pm 0.47	73.58 \pm 0.44	<u>90.21 \pm 0.65</u>	56.13 \pm 0.30
SEGA+AD-GCL	<u>75.03 \pm 0.36</u>	72.32 \pm 0.49	87.74 \pm 0.39	54.29 \pm 0.54
SEGA+JOAOv2	72.76 \pm 0.27	72.12 \pm 0.79	87.98 \pm 0.29	56.15 \pm 0.29
SEGA+AutoGCL	72.68 \pm 0.23	<u>73.95 \pm 0.87</u>	89.88 \pm 1.21	57.43 \pm 0.37
GMNI	76.19 \pm 0.71	75.54 \pm 0.29	91.29 \pm 0.41	<u>57.02 \pm 0.14</u>

ing 1. Furthermore, GMNI outperforms automated DA methods across all datasets. The experimental results presented in Table 5 demonstrate that GMNI outperforms both SEGA and SEGA with collaborative models on 7 out of 8 datasets. The results of the graph and node classification tasks prove that our strategy of using minimal noteworthy information to approximate minimal necessary information is effective on different downstream tasks. GMNI approximates the optimal views of InfoMin, without requiring task-relevant information. Lastly, Table 3 and Table 4 illustrate that methods with DA generally outperform methods without DA, highlighting the significance of DA in GCL.

Semi-supervised learning on graph classification.

Settings. Following AutoGCL (Yin et al., 2022), we perform semi-supervised learning on TUDataset using 10-fold cross-validation.

We apply three methods for GMNI on semi-supervised training: GMNI-Label, GMNI-Alter and GMNI-Unlabel-Fint. GMNI-Label is trained and tested using 10% labeled data. GMNI-Alter involves several alternate training steps, where each step comprises contrastive learning (CL) with 80% unlabeled data and fine-tuning (FT) with 10% labeled data. Finally, it is tested on 10% labeled data. GMNI-Unlabel-Fint is trained on 80% unlabeled data, fine-tuned on 10% labeled data, and tested on 10% labeled data. GMNI-Alter alternates between CL and FT in a sequential manner, i.e., (CL1-FT1-CL2-FT2-....), whereas GMNI-Unlabel-Fint performs FT after completing CL.

Result analysis. Table 6 shows that GMNI outperforms baselines on 7 out

Table 6: Semi-supervised learning performance on graph classification (average accuracy (%) \pm std. over 5 runs). The highest and second performances are in **bold** and underlined, respectively.

Dataset	GCA	GraphCL	JOAOv2	AD-GCL	AutoGCL
PROTEINS	73.85 \pm 5.56	74.21 \pm 4.50	73.31 \pm 0.48	73.96 \pm 0.47	75.65 \pm 2.40
DD	76.74 \pm 4.09	76.65 \pm 5.12	75.81 \pm 0.73	77.91 \pm 0.73	77.50 \pm 4.41
NC11	68.73 \pm 2.36	73.16 \pm 2.90	74.86 \pm 0.39	75.18 \pm 0.31	73.75 \pm 2.25
COLLAB	74.32 \pm 2.30	75.50 \pm 2.15	75.53 \pm 0.18	75.82 \pm 0.26	77.16 \pm 1.48
GITHUB	59.24 \pm 3.21	63.51 \pm 1.02	<u>66.66 \pm 0.60</u>	-	62.46 \pm 1.51
IMDB-B	<u>73.70 \pm 4.88</u>	68.10 \pm 5.15	-	-	71.90 \pm 4.79
REDDIT-B	77.15 \pm 6.96	78.05 \pm 2.65	88.79 \pm 0.65	90.10 \pm 0.15	79.80 \pm 3.47
REDDIT-M-5K	32.95 \pm 10.89	48.09 \pm 1.74	52.71 \pm 0.28	<u>53.49 \pm 0.28</u>	49.91 \pm 2.70

Dataset	SEGA	GMNI-Label	GMNI-Alter	GMNI-Unlabel-Fint	Ranks
PROTEINS	74.65 \pm 0.54	78.31 \pm 1.31	<u>76.46 \pm 0.67</u>	74.71 \pm 0.59	1
DD	76.33 \pm 0.43	80.69 \pm 0.47	<u>79.30 \pm 1.04</u>	76.38 \pm 0.91	1
NC11	75.09 \pm 0.22	76.43 \pm 0.50	<u>76.13 \pm 0.59</u>	74.32 \pm 0.19	1
COLLAB	75.18 \pm 0.22	75.56 \pm 0.30	78.50 \pm 0.48	<u>77.21 \pm 0.47</u>	1
GITHUB	66.01 \pm 0.66	65.62 \pm 0.15	68.46 \pm 0.23	66.32 \pm 0.41	1
IMDB-B	-	74.16 \pm 0.57	71.96 \pm 0.84	69.66 \pm 1.00	1
REDDIT-B	<u>89.40 \pm 0.23</u>	82.68 \pm 0.57	87.96 \pm 0.35	85.91 \pm 1.94	4
REDDIT-M-5K	<u>53.73 \pm 0.28</u>	50.16 \pm 0.77	53.96 \pm 0.51	52.77 \pm 0.57	1

529 of 8 datasets and gets an average rank of 1.38. Intriguingly, unlabeled data is
530 not always effective, as the performance of GMNI-Unlabel-Fint is sometimes
531 worse than that of GMNI-Label, which is trained with only 10% labeled data.
532 Using labels directly to fine-tune the model, which was initially trained with
533 unlabeled data, can interfere with the already learned distribution, resulting
534 in a performance loss. GMNI-Alter outperforms GMNI-Unlabel-Fint by us-
535 ing unlabeled and labeled data alternately. The alternating training of CL
536 and FT helps alleviate the divergence of learned knowledge between unla-
537 belled and labeled data, promoting mutual learning between the two sources.
538 GMNI’s performance on REDDIT-BINARY is inferior to the two baseline
539 methods due to the limited generalization ability of the shared view gener-
540 ator across all graphs. Thus, training a dedicated view generator for each
541 batch can significantly boost performance, with scores of GMNI-Label at
542 91.10 ± 0.36 (Rank 1), GMNI-Alter at 89.28 ± 1.83 , and GMNI-Unlabel-Fint
543 at 88.35 ± 0.23 . However, due to computational constraints, we do not imple-
544 ment this training strategy in practice.

545 4.4. Ablation studies

546 Ablation studies on 10 datasets, conducted with identical hyperparame-
547 ters, are presented in Table 7 and Figure 4 to validate the rationale behind

Table 7: Ablation studies on GMNI (average accuracy (%) \pm std. over 3 runs). GMNI-Uni, GMNI-UniE-F, GMNI-E-UniF, and GMNI are one group of ablation experiments, while GMNI-Simp, GMNI-ViewM, GMNI-Simult, and GMNI are another group of ablation experiments. The highest performance on the first and the second ablation group are underlined and in **bold**, respectively.

Dataset	GMNI-Uni	GMNI-UniE-F	GMNI-E-UniF	GMNI
Wiki-CS	<u>79.30 \pm 0.00</u>	76.43 \pm 0.15	79.16 \pm 0.07	79.29 \pm 0.05
Amazon-Computers	88.12 \pm 0.25	86.74 \pm 0.28	88.87 \pm 0.17	<u>89.29 \pm 0.08</u>
Amazon-Photo	93.27 \pm 0.47	91.21 \pm 0.49	93.57 \pm 0.26	<u>93.68 \pm 0.35</u>
Coauthor-CS	93.32 \pm 0.05	93.41 \pm 0.10	93.33 \pm 0.02	<u>93.62 \pm 0.10</u>
MUTUG	90.79 \pm 0.28	91.13 \pm 0.77	90.24 \pm 0.42	91.83 \pm 0.62
PROTEINS	76.28 \pm 0.09	75.17 \pm 0.37	76.25 \pm 0.51	<u>76.61 \pm 0.37</u>
DD	79.37 \pm 0.59	76.62 \pm 0.48	<u>79.62 \pm 0.09</u>	79.43 \pm 0.34
COLLAB	72.30 \pm 1.06	74.00 \pm 0.28	<u>77.02 \pm 1.46</u>	76.08 \pm 0.86
IMDB-B	75.20 \pm 0.70	72.67 \pm 0.06	75.07 \pm 0.76	75.70 \pm 0.26
REDDIT-B	90.05 \pm 0.61	<u>92.00 \pm 0.48</u>	90.05 \pm 0.30	91.25 \pm 0.75
Dataset	GMNI-Simp	GMNI-ViewM	GMNI-Simult	GMNI
Wiki-CS	79.10 \pm 0.15	79.34 \pm 0.05	79.13 \pm 0.12	79.29 \pm 0.05
Amazon-Computers	89.12 \pm 0.16	89.04 \pm 0.08	89.00 \pm 0.32	89.29 \pm 0.08
Amazon-Photo	93.65 \pm 0.32	93.44 \pm 0.39	93.80 \pm 0.27	93.68 \pm 0.35
Coauthor-CS	93.44 \pm 0.08	93.48 \pm 0.02	93.39 \pm 0.12	93.62 \pm 0.10
MUTUG	91.53 \pm 1.02	90.45 \pm 0.54	91.35 \pm 0.61	91.83 \pm 0.62
PROTEINS	76.49 \pm 0.19	76.53 \pm 0.28	77.30 \pm 0.21	76.61 \pm 0.37
DD	79.20 \pm 0.22	79.00 \pm 0.38	78.83 \pm 0.05	79.43 \pm 0.34
COLLAB	70.39 \pm 0.95	75.91 \pm 0.10	74.12 \pm 0.63	76.08 \pm 0.86
IMDB-B	70.23 \pm 1.50	75.33 \pm 1.07	75.60 \pm 0.66	75.70 \pm 0.26
REDDIT-B	91.63 \pm 0.25	90.90 \pm 0.38	89.92 \pm 0.40	91.25 \pm 0.75

548 the components of GMNI. Table 8 provides a comparison to clearly illustrate
 549 the differences between the variants of GMNI.

550 *Effectiveness of automated view generator.*

551 **Settings.** We use 6 variants to verify the effectiveness of the generator.
 552 GMNI-Uni employs uniformly distributed edge and feature importance. GMNI-
 553 UniE-F employs uniformly distributed edge importance but retains the auto-
 554 mated generation of feature importance. GMNI-E-UniF employs uniformly
 555 distributed feature importance but retains the automated generation of edge
 556 importance. The previously mentioned variants perform DA on both topol-
 557 ogy and features simultaneously. GMNI-F masks only features based on au-
 558 tomatically generated feature importance, while GMNI-E drops only edges
 559 based on automatically generated edge importance. GMNI-EF applies masks
 560 separately to edges and features in two views according to automatically gen-
 561 erated importance.

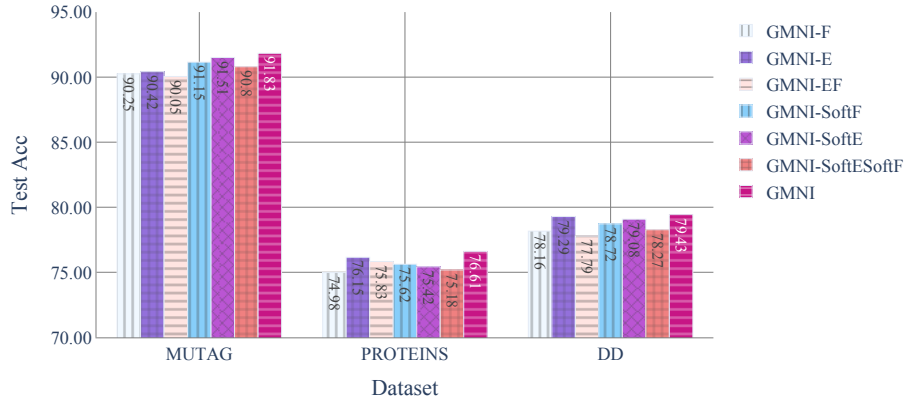


Figure 4: Ablation studies on GMNI (average accuracy (%) over 3 runs). GMNI-F, GMNI-E, GMNI-EF, GMNI-SoftF, GMNI-SoftE, and GMNI-SoftESoftF are different variations of GMNI.

Table 8: Different variations of GMNI for ablation studies.

	GMNI-UniE-F	GMNI-F	GMNI-E-UniF	GMNI-E	GMNI-EF
View1 Edge	Uniform Distribution	None	Edge Importance	Edge Importance	None
View1 Feat	Feature Importance	Feature Importance	Uniform Distribution	None	Feature Importance
View2 Edge	Uniform Distribution	None	Edge Importance	Edge Importance	Edge Importance
View2 Feat	Feature Importance	Feature Importance	Uniform Distribution	None	None
	GMNI-SoftF	GMNI-SoftE	GMNI-SoftESoftF	GMNI	
View1 Edge	Edge Importance	Soft Edges	Soft Edges	Edge Importance	
View1 Feat	Soft Features	Feature Importance	Soft Features	Feature Importance	
View2 Edge	Edge Importance	Soft Edges	Soft Edges	Edge Importance	
View2 Feat	Soft Features	Feature Importance	Soft Features	Feature Importance	

562 **Result analysis.** According to Table 7, GMNI outperforms GMNI-Uni on
563 all datasets except Wiki-CS, and performs similarly on Wiki-CS, demonstrat-
564 ing the effectiveness of the automated view generator. GMNI outperforms
565 GMNI-UniE-F on 9 datasets and outperforms GMNI-E-UniF on 8 datasets.
566 The results indicate that combining edge and feature importance is more
567 effective than using only one. According to Figure 4, the performance of
568 GMNI-F, GMNI-E, and GMNI-EF is inferior to that of GMNI. Therefore,
569 GDA should simultaneously incorporate topological and feature augmenta-
570 tion on both views to achieve optimal performance.

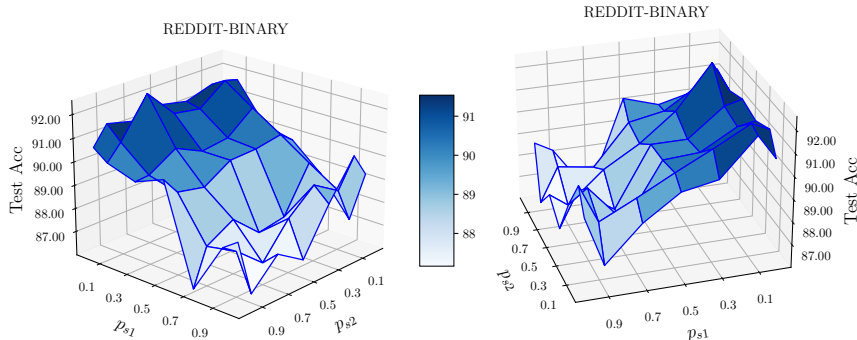


Figure 5: Graph classification performance of GMNI with different p_{s1} and p_{s2} on REDDIT-BINARY.

571 *Effectiveness of randomness.*

572 **Settings.** We employ 3 variations to validate the rationality of introduc-
 573 ing randomness in view generation. GMNI-SoftE denotes a GMNI variant
 574 that uses edge importance as soft edges, while GMNI-SoftF signifies a vari-
 575 ant that uses feature importance as soft features. **GMNI-SoftE and GMNI-**
 576 **SoftF eliminate randomness from topology and feature data augmentation,**
 577 **respectively. GMNI-SoftESoftF combines both soft edges and soft features,**
 578 **removing randomness from both topology and feature data augmentation.**

579

580 **Result analysis.** According to Figure 4, the experimental results indi-
 581 cate that GMNI outperforms both GMNI-SoftE and GMNI-SoftF, while
 582 GMNI-SoftE and GMNI-SoftF outperform GMNI-SoftESoftF. It indicates
 583 that GMNI, which incorporates randomness in both feature and topology
 584 DA, performs the best, while GMNI-SoftESoftF, which removes all random-
 585 ness, performs the worst.

586 *Rationality of max-min optimization.*

587 **Settings.** We employ 3 variations to validate the rationality of adversar-
 588 ial optimization in the automated view generator. GMNI-Simp replaces the
 589 learnable GCN encoder in the view generator with a parameterless GCN, i.e.,
 590 $\mathbf{W} = \mathbf{I}$. GMNI-ViewM optimizes the view generator through max-max op-
 591 timization. GMNI-Simult trains the view generator and comparison module
 592 simultaneously, using the GNN encoder from the view comparison module
 593 as the GNN encoder for the view generator.

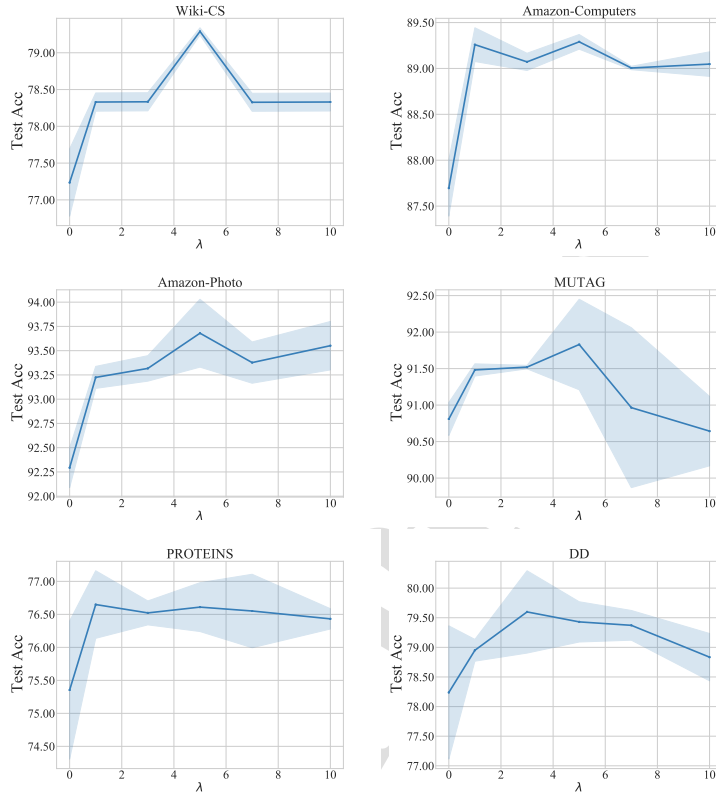


Figure 6: Node and graph classification performance of GMNI with different λ .

594 **Result analysis.** According to Table 7, the following conclusions can be
 595 drawn. Firstly, GMNI performs better than GMNI-Simp on 9 datasets, proving
 596 that the learnable GCN encoder in the generator is effective. Secondly,
 597 GMNI outperforms GMNI-ViewM on 9 datasets, showing that the GNN
 598 encoder and the projection head trained by minimizing are better than max-
 599 imizing. **Minimizing the objective results in an aggressive graph encoder and**
 600 **projection head, which contribute to producing a more generalized impor-**
 601 **tance graph and a more precise approximation of InfoMin’s optimal views.**
 602 Finally, GMNI beats GMNI-Simult on 8 datasets, as training the view gen-
 603 erator and the comparison module simultaneously prevents the generation of

Table 9: Ablation studies on ℓ_2 distance minimization, cosine similarity maximization, and InfoMax across 13 datasets (average accuracy (%) \pm std. over 3 runs). The highest performance is indicated in **bold**.

Dataset	ℓ_2 Distance	Cosine Similarity	InfoMax
Wiki-CS	75.69 \pm 0.21	77.28 \pm 0.11	79.19 \pm 0.13
Amazon-Computers	66.85 \pm 0.89	85.87 \pm 0.57	89.29 \pm 0.05
Amazon-Photo	84.07 \pm 0.68	91.66 \pm 0.24	93.52 \pm 0.33
Coauthor-CS	89.22 \pm 0.58	93.42 \pm 0.05	93.61 \pm 0.15
Coauthor-Physics	85.75 \pm 0.93	95.79 \pm 0.03	95.91 \pm 0.10
MUTAG	90.62 \pm 0.26	91.51 \pm 0.39	91.57 \pm 0.57
PROTEINS	76.01 \pm 0.26	76.28 \pm 0.25	76.58 \pm 0.29
DD	77.70 \pm 0.21	77.54 \pm 0.20	79.32 \pm 0.26
NC11	79.91 \pm 0.36	81.33 \pm 0.10	83.12 \pm 0.20
COLLAB	69.75 \pm 0.27	73.99 \pm 0.44	76.19 \pm 0.71
IMDB-B	72.77 \pm 0.17	73.17 \pm 0.33	75.54 \pm 0.29
REDDIT-B	87.65 \pm 0.50	88.37 \pm 0.49	91.29 \pm 0.41
REDDIT-M-5K	54.05 \pm 0.34	53.71 \pm 0.53	57.02 \pm 0.14

604 superior views and representations. The view generator emphasizes remain-
 605 ing noteworthy information and discarding nuisance information, whereas
 606 the comparison module highlights the quality of the graph encoder. They
 607 impede each other when training together.

608

609 *Effectiveness of InfoMax.*

610 **Settings.** First, we use two variants to validate the rationale of InfoMax by
 611 modifying the optimization objectives in the GMNI view comparison module
 612 to ℓ_2 distance minimization and cosine similarity maximization. Second, Info-
 613 Max consists of two parts: pulling the representations of the same node closer
 614 and pushing those of different nodes further apart. Thus, we decompose this
 615 objective into two distinct optimization goals: one for solely maximizing the
 616 similarity between representations of the same node, and another for solely
 617 minimizing the similarity between representations of different nodes.

618 **Result analysis.** First, the experimental results in Table 9 demonstrate
 619 that InfoMax consistently outperforms both ℓ_2 distance and cosine similar-
 620 ity. Second, the experimental results in Table 10 indicate that InfoMax is
 621 optimal. Focusing solely on bringing representations of the same node closer
 622 may cause all representations to converge into a uniform state, while solely
 623 pushing representations of different nodes apart may result in overly dis-
 624 persed representations, failing to achieve the classification goal.

Table 10: Ablation studies on the decomposition of the InfoMax over 13 datasets (average accuracy (%) \pm std. over 3 runs). The highest performance is indicated in **bold**.

Dataset	OnlyClose	OnlyFar	InfoMax
Wiki-CS	75.47 \pm 0.03	77.28 \pm 0.11	79.19 \pm 0.13
Amazon-Computers	66.84 \pm 0.88	85.87 \pm 0.57	89.29 \pm 0.05
Amazon-Photo	84.07 \pm 0.69	91.66 \pm 0.24	93.52 \pm 0.33
Coauthor-CS	89.30 \pm 0.60	93.42 \pm 0.05	93.61 \pm 0.15
Coauthor-Physics	85.89 \pm 0.91	95.79 \pm 0.03	95.91 \pm 0.10
MUTAG	90.61 \pm 0.27	91.51 \pm 0.39	91.57 \pm 0.57
PROTEINS	75.35 \pm 0.79	76.28 \pm 0.25	76.58 \pm 0.29
DD	77.14 \pm 0.63	77.54 \pm 0.20	79.32 \pm 0.26
NCI1	80.84 \pm 0.24	81.33 \pm 0.10	83.12 \pm 0.20
COLLAB	69.45 \pm 0.23	73.99 \pm 0.44	76.19 \pm 0.71
IMDB-B	73.10 \pm 0.24	73.17 \pm 0.33	75.54 \pm 0.29
REDDIT-B	86.45 \pm 0.57	88.37 \pm 0.49	91.29 \pm 0.41
REDDIT-M-5K	52.32 \pm 0.75	53.71 \pm 0.53	57.02 \pm 0.14

625 4.5. Additional experiments

626 *Hyperparameter analysis.* We analyze critical hyperparameters p_{s1}, p_{s2} of GMNI.
 627 To simplify the analysis, p_{s1} and p_{s2} of the two views are set to be identical,
 628 and they are utilized to adjust the overall drop rate of edges and the overall
 629 mask rate of features. p_{s1} and p_{s2} are selected from [0.0, 0.1, 0.3, 0.5, 0.7,
 630 0.9, 1.0]. The results are shown in Figure 5.

631 Firstly, it can be observed that GMNI is more sensitive to p_{s1} than p_{s2} ;
 632 altering p_{s1} results in a more noticeable performance change. This observa-
 633 tion suggests that topology information plays a more crucial role in GMNI
 634 than feature information. Secondly, it is worth noting that excessively large
 635 values of p_{s1} or p_{s2} have a detrimental effect on performance. When p_{s1} be-
 636 comes excessively large, the topology information of the graph is almost com-
 637 pletely destroyed, resulting in isolated nodes that lack connections. When
 638 p_{s2} becomes excessively large, the graph almost degenerates into a featureless
 639 graph. The performance tends to be extremely poor when both are large.
 640 Thirdly, points in Figure 5 surrounding (0,0) perform better than (0,0). This
 641 observation suggests that utilizing importance to generate views, even with
 642 small drop or mask rates, is more advantageous than solely employing the
 643 original graph for comparison.

644 Subsequently, we analyze the hyperparameter λ used for regularization.
 645 The results are shown in Figure 6. Firstly, the model’s performance deteri-
 646 orates when the regularization weight λ is set to zero. This occurs because

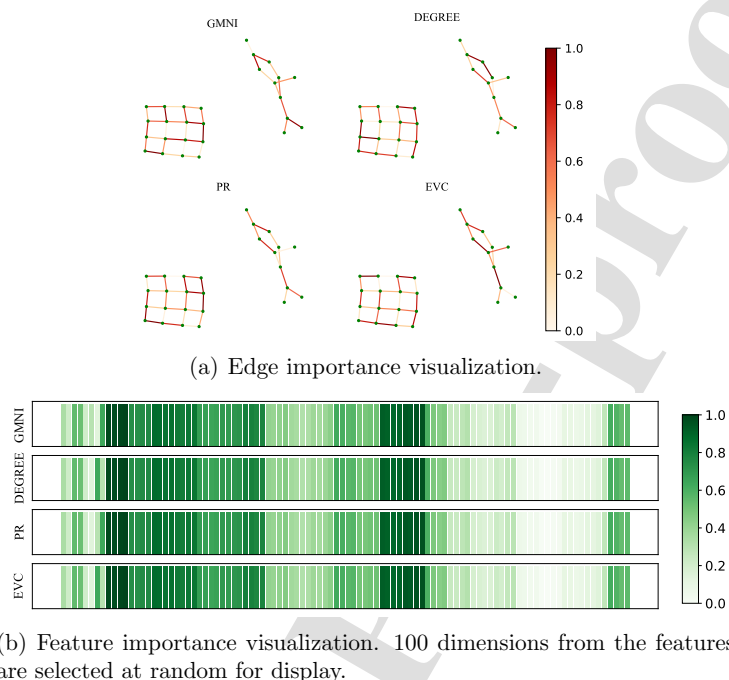


Figure 7: DEGREE, PR, and EVC are degree, PageRank and eigenvector centrality, respectively.

647 of the absence of the ‘minimal condition’ required for ‘minimal noteworthy
 648 information,’ resulting in excessive redundant information in the views. Sec-
 649 ondly, on the MUTAG and DD datasets, we observe a decrease in model
 650 performance with a large regularization weight λ . This is attributed to the
 651 overly restrictive constraints on \mathbf{P}_f and \mathbf{P}_e , which hinder the inclusion of
 652 noteworthy information in the importance graph.

653 *Importance visualization.* Centrality reflects the significance of a node in a
 654 graph. In GCA, edge and feature importance are defined by degree, PageR-
 655 ank, and eigenvector centrality. To demonstrate the reasonableness of the
 656 importance gained by GMNI, we visualize the importance defined by those
 657 centralities and obtained by GMNI on a graph of MUTAG. MUTAG is a
 658 collection of nitroaromatic compounds, where the graph’s vertices stand for
 659 atoms, and the edges between vertices represent bonds between the corre-

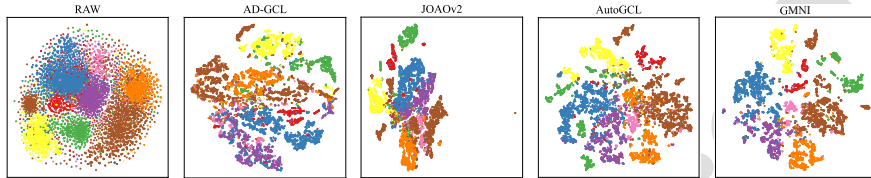


Figure 8: t-SNE visualization of raw features and node representations on Amazon-Photo. Different colors represent different node classes.

660 sponding atoms. In Figure 7(b), the feature importance gained by GMNI
 661 aligns with those acquired by centrality, demonstrating the rationality of the
 662 automatically generated importance.

663 *Embedding visualization.* To demonstrate the quality of the representations,
 664 we employ two components PCA (Abdi and Williams, 2010) and t-SNE (Van der
 665 Maaten and Hinton, 2008) to visualize the raw features and representations
 666 of AD-GCL, JOAOv2, AutoGCL, and GMNI on Amazon-Photo. Figure 8
 667 demonstrates that GMNI achieves a more distinct classification boundary,
 668 showcasing the superior representations learned by GMNI. These represen-
 669 tations prove to be effective for the downstream classification task.

670 5. Conclusion

671 In this paper, we propose GMNI, a novel GCL method featuring auto-
 672 mated DA. GMNI approximates InfoMin’s optimal views by replacing min-
 673 imal necessary information with minimal noteworthy information, without
 674 requiring task-relevant information. Applying InfoMax to these views can
 675 avoid the risk of redundant information and insufficient information. In ad-
 676 dition, GMNI introduces randomness to augmentation, thus enhancing view
 677 diversity and stabilizing the model against perturbations. Extensive experi-
 678 ments on node and graph classification tasks demonstrate GMNI’s superior-
 679 ity, which outperforms both automated and manual DA GCL methods.

680 **In the future, we will further validate the effectiveness of our approx-**
 681 **imate optimal views on more downstream tasks and the strong correlation**
 682 **between noteworthy and necessary information. Furthermore, approximating**
 683 **InfoMin’s optimal views in the spectral domain is another promising research**
 684 **approach.**

685 **Declaration of competing interest**

686 The authors declare the following financial interests/personal relationships
 687 which may be considered as potential competing interests: Furao Shen
 688 reports financial support was provided by the STI 2030-Major Projects of
 689 China under Grant 2021ZD0201300, the National Science Foundation of
 690 China under Grant 62276127.

691 **Data availability**

692 Data will be made available on request.

693 **Acknowledgments**

694 This work was supported in part by the STI 2030-Major Projects of China
 695 under Grant 2021ZD0201300, and by the National Science Foundation of
 696 China under Grant 62276127.

697 **Appendix A. Experimental settings**

698 Table A.11 details the hyperparameters of GMNI for unsupervised node
 699 classification. The GNN encoder and projection head of the view generator
 700 have the same hidden dimensions as the GCN, and the projection head is
 701 1-layer. The GNN encoder of the view comparison module is 2-layer, and
 702 the projection head is 1-layer.

703 Table A.12 details the hyperparameters of GMNI for unsupervised graph
 704 classification. The projection head of the view generator is 2-layer, whereas
 705 the view comparison uses a single layer.

706 For semi-supervised learning, we follow AutoGCL (Yin et al., 2022) and
 707 use ResGCN (Chen et al., 2019) with a hidden size of 128 as the GNN encoder
 708 for view comparison. All other hyperparameters remain the same as those
 709 used for unsupervised graph classification. For GITHUB, the parameters are
 710 set as follows: (view 1 p_{s1} , view 2 p_{s1} , view 1 p_{s2} , view 2 p_{s2}) = (0.3, 0.2,
 711 0.3, 0.2). The GNN encoder in the view comparison module has a hidden
 712 dimension of 64 and is 3-layer. The learning rate for both the view generator
 713 and the view comparison module is set to 0.01.

714 In GMNI, the sum of node representations of the two views and the original
 715 graph: $f_n(\mathcal{G}_1) + f_n(\mathcal{G}_2) + 2 * f_n(\mathcal{G})$ or the sum of node representations
 716 of the two views: $f_n(\mathcal{G}_1) + f_n(\mathcal{G}_2)$ is used for node classification, considering

Table A.11: Hyperparameters of GMNI for unsupervised node classification. We substitute ‘dim’ for ‘dimension’ and ‘proj’ for ‘projection head.’ The symbol ” is the ditto mark.

GMNI View Generator	GCN hidden dim	h_ψ structure	h_ϕ structure	GNN #layers	τ	λ	activation	learning rate
Wiki-CS	128	1	1	1	1	5	PReLU	0.001
Amazon-Computers	128	64-1	32-1	2	”	5	RReLU	0.001
Amazon-Photo	256	128-1	64-1	2	”	5	ReLU	0.001
Coauthor-CS	512	256-1	64-1	1	”	5	RReLU	0.001
Coauthor-Physics	128	64-1	32-1	2	”	0.5	RReLU	0.01

GMNI View Comparison	view 1 p_{s1}	view 2 p_{s1}	view 1 p_{s2}	view 2 p_{s2}	p_t	GNN hidden dim	GNN #layers	proj hidden dim	ϵ	learning rate
Wiki-CS	0.2	0.4	0.1	0.1	0.6	128	2	128	0.4	0.01
Amazon-Computers	0.6	0.3	0.2	0.1	0.7	128	”	128	0.2	0.001
Amazon-Photo	0.3	0.5	0.1	0.1	0.7	256	”	64	0.3	0.001
Coauthor-CS	0.1	0.1	0.3	0.2	0.7	256	”	256	0.4	0.0001
Coauthor-Physics	0.4	0.4	0.1	0.4	0.7	128	”	64	0.5	0.001

Table A.12: Hyperparameters of GMNI for unsupervised graph classification. We substitute ‘VG’ for ‘View Generator’, and ‘VC’ for ‘View Comparison.’

GMNI	VG’s GCN hidden dim	VG’s h_ψ structure	VG’s h_ϕ structure	VG’s GNN #layers	τ	λ	activation	VG’s learning rate
All Datasets	64	32-1	16-1	1	1	5	ReLU	0.001

	view 1 p_{s1}	view 2 p_{s1}	view 1 p_{s2}	view 2 p_{s2}	p_t	VC’s GNN hidden dim	VC’s GNN #layers	ϵ	VC’s learning rate
MUTAG	0.3	0.2	0.3	0.2	0.7	64	5	0.2	0.001
PROTEINS	0.3	0.2	0.3	0.2	”	128	5	”	”
DD	0.3	0.2	0.3	0.2	”	128	5	”	”
NCI	0.1	0.2	0.4	0.2	”	128	5	”	”
COLLAB	0.3	0.2	0.3	0.2	”	64	2	”	”
IMDB-B	0.3	0.2	0.3	0.2	”	128	5	”	”
REDDIT-B	0.3	0.2	0.3	0.2	”	64	2	”	”
REDDIT-M-5K	0.4	0.2	0.3	0.2	”	128	5	”	”

Table A.13: Average training time per epoch of GCL methods with automated DA.

Dataset	AD-GCL	JOAOv2	AutoGCL	GMNI VG	GMNI VC
Wiki-CS	0.37 ± 0.09	1.53 ± 0.15	0.29 ± 0.09	0.26 ± 0.03	0.10 ± 0.00
Amaz-Comp	0.42 ± 0.09	2.28 ± 0.12	0.40 ± 0.08	0.42 ± 0.05	0.13 ± 0.00
Amaz-Photo	0.29 ± 0.10	0.94 ± 0.10	0.27 ± 0.08	0.33 ± 0.06	0.06 ± 0.00
Coauthor-CS	0.76 ± 0.12	7.25 ± 0.28	2.01 ± 0.84	1.53 ± 0.05	0.48 ± 0.01
Coauthor-Phy	0.40 ± 0.16	-	-	2.53 ± 0.20	0.93 ± 0.04
MUTAG	0.27 ± 0.00	7.83 ± 0.17	0.06 ± 0.00	0.03 ± 0.00	0.03 ± 0.00
PROTEINS	1.76 ± 0.15	15.99 ± 0.96	0.54 ± 0.17	0.24 ± 0.02	0.15 ± 0.02
DD	2.52 ± 0.22	171.16 ± 2.73	0.71 ± 0.15	4.33 ± 0.13	3.96 ± 0.05
NCII	6.82 ± 0.54	106.52 ± 1.98	1.10 ± 0.04	0.66 ± 0.04	0.49 ± 0.02
COLLAB	13.78 ± 1.47	419.96 ± 10.61	2.58 ± 0.41	8.01 ± 0.07	3.66 ± 0.05
IMDB-B	1.48 ± 0.01	11.95 ± 0.53	0.26 ± 0.03	0.19 ± 0.02	0.10 ± 0.00
RDT-B	4.79 ± 0.41	153.11 ± 6.83	1.09 ± 0.09	1.55 ± 0.05	0.44 ± 0.01
RDT-M-5K	12.85 ± 1.73	520.64 ± 27.55	3.37 ± 0.46	2.47 ± 0.07	2.14 ± 0.02

717 that incorporating view representations can lead to more generalized repre-
718 sentations. Graph representations $\{f_g(\mathcal{G}_i), i = 1, \dots, Q\}$ are used for graph

719 classification. On all datasets, the Xavier initialization (Glorot and Bengio,
720 2010) and Adam optimizer (Kingma and Ba, 2014) are utilized.

721 All our experiments are conducted on four Tesla V100-SXM2 (32GB)
722 with Intel(R) Xeon(R) Gold 6248 CPU @ 2.50GHz, four NVIDIA GeForce
723 RTX 2080 Ti (11GB) with Intel(R) Xeon(R) CPU E5-1650 v4 @ 3.60GHz,
724 and two TITAN Xp (12GB) GPUs with Intel(R) Xeon(R) CPU E5-2643 v4
725 @ 3.40GHz. Python version is 3.8.8, and PyTorch version is 1.12.0+cu102.

726 Appendix B. Notations and abbreviations

Notation	Meaning
\mathcal{G}	The graph.
$\mathcal{V} = \{v_1, \dots, v_N\}$	The set of nodes.
$\mathcal{E} = \{\varepsilon_1, \dots, \varepsilon_M\}$	The set of edges.
$\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \in \mathbb{R}^{N \times F}$	The feature matrix. F is the feature dimension.
$\mathbf{A} \in \{0, 1\}^{N \times N} : \mathbf{A}_{ij} = \mathbb{I}((v_i, v_j) \in \mathcal{E})$	The adjacency matrix.
$\mathbf{D} : \mathbf{D}_{ii} = \sum_j \mathbf{A}_{ij}$	The degree matrix.
$f_n : \mathcal{G} \rightarrow \mathbb{R}^{N \times d}$	The encoder for node classification. d is the output dimension.
$f_g : \mathcal{G}_i \rightarrow \mathbb{R}^d$	The encoder for graph classification.
$\mathbf{P}_e \in \mathbb{R}^{M \times 1}$	The edge importance.
$\mathbf{P}_f \in \mathbb{R}^{F \times 1}$	The feature importance.
$\mathbf{P}_n \in \mathbb{R}^{N \times 1}$	The node importance.
$\tilde{\mathbf{A}}$	The adjacency matrix with edge importance.
$\tilde{\mathbf{X}}$	The feature matrix with feature importance.
$\tilde{\mathcal{G}} = T(\mathcal{G}) = (\tilde{\mathbf{A}}, \tilde{\mathbf{X}})$	The importance graph. T is the transformation function that maps \mathcal{G} to $\tilde{\mathcal{G}}$.
$ T(\mathcal{G}) $	The normalized sum of edge importance and feature importance.
ζ	The mutual information lower bound of our objective.
\mathcal{I}	The noteworthy information.
\mathbf{W}	The parameter matrix of the GCN.
\mathbf{E}	The node representation.
$\sigma(\cdot)$	The activation function.
h_ϕ, h_ψ	Simple MLPs.
τ, ϵ	The temperature parameter.
Gumbel(\cdot)	The Gumbel-Max reparametrization function.
η	A random variable following the uniform distribution $(0, 1)$.
$g_\xi(\cdot)$	A 1- or 2- layer MLP as the projection head.

$f_{\theta}(\cdot)$	A 1- or 2- layer graph encoder.
λ	The regularization weight.
$\hat{I}(\cdot; \cdot)$	The mutual information estimator.
$\mathbf{P}_{d,\varepsilon}$	The edge drop rate.
$\mathbf{P}_{d,f}$	The feature mask rate.
$\mathbf{M} \in \mathbb{R}^{N \times F}$	The mask matrix.
$\mathcal{G}_1 : (\mathbf{A}_1, \mathbf{X}_1), \mathcal{G}_2 : (\mathbf{A}_2, \mathbf{X}_2)$	Two views.
p_{s1}	The hyperparameter used to adjust the overall drop rate of edges.
p_{s2}	The hyperparameter used to adjust the overall mask rate of features.
p_t	The hyperparameter used to truncate the drop rate of edges and features.
g_n, g_g	Projection heads.
Abbreviation	Expanded Version
CL	Contrastive Learning
GCL	Graph Contrastive Learning
DA	Data Augmentation
GDA	Graph Data Augmentation
MI	Mutual Information
MNI	Minimal Noteworthy Information

727 Appendix C. Algorithm

728 Algorithm 1 describes the details of the automated graph view generator
 729 for GMNI. Algorithm 2 and 3 describe the details of the view comparison
 730 module for GMNI on node and graph classification, respectively.

731 For a graph with N nodes, F features, and M edges, the view generator
 732 of GMNI requires $O(N^2F + (M + N)F^2)$ floating-point operations (FLOPs),
 733 while the view comparison requires $O(N^2F + NF^2)$ FLOPs. The FLOPs of
 734 AD-GCL, JOAOv2, and AutoGCL are $O(N^2F + (M + N)F^2)$, $O(N^2F +$
 735 $NF^2)$, and $O(N^2F + NF^2)$, respectively. Thus, GMNI’s view generator has
 736 an asymptotic complexity comparable to that of AD-GCL, and GMNI’s view
 737 comparison has an asymptotic complexity comparable to that of JOAO and
 738 AutoGCL.

739 Table A.13 illustrates the average training time per epoch for various
 740 GCL methods with automated DA on a Tesla V100-SXM2 (32GB). GMNI is
 741 efficient in terms of computing resources and time. Firstly, the view gener-
 742 ator training can be completed in a relatively small number of epochs. Sec-
 743 ondly, the view comparison module in GMNI focuses solely on maximization,

⁷⁴⁴ whereas the other three models require both maximization and minimization
⁷⁴⁵ processes.

Journal Pre-proof

Algorithm 1 Automated graph view generator

Require: graph $\mathcal{G} = (\mathbf{A}, \mathbf{X})$; edge set \mathcal{E} , $|\mathcal{E}| = M$; feature dimension F ; GCN encoder $f_{\mathbf{W}}(\cdot)$; feature importance MLP $h_{\phi}(\cdot)$; edge importance MLP $h_{\psi}(\cdot)$; Gumbel-Max function $\text{Gumbel}(\cdot)$; GNN encoder $f_{\theta}(\cdot)$; projection head $g_{\xi}(\cdot)$; mutual information estimator $\hat{I}(\cdot; \cdot)$.

HyperParas: regularization weight λ ; scale parameters for view 1 and view 2 $p_{1,s1}, p_{2,s1}, p_{1,s2}, p_{2,s2}$; truncate parameter p_t ; learning rate α .

Ensure: view 1 $\mathcal{G}_1 = (\mathbf{A}_1, \mathbf{X}_1)$; view 2 $\mathcal{G}_2 = (\mathbf{A}_2, \mathbf{X}_2)$.

- 1: **for** ep in #epochs **do**
- 2: $\mathbf{E} = f_{\mathbf{W}}(\mathbf{A}, \mathbf{X})$;
- 3: $\mathbf{P}_n = \text{Gumbel}(h_{\phi}(\mathbf{E}))$;
- 4: $\mathbf{P}_f = \mathbf{X}^T \mathbf{P}_n$;
- 5: **for** $\forall \varepsilon_k \in \mathcal{E}$ **do**
- 6: $\mathbf{p}_{e,k} = \text{Gumbel}(h_{\psi}([\mathbf{E}[v_i]; \mathbf{E}[v_j]]))$, $\varepsilon_k = (v_i, v_j)$;
- 7: **end for**
- 8: $\mathbf{P}_e = [\mathbf{p}_{e,1}, \dots, \mathbf{p}_{e,M}]^T$;
- 9: Broadcast $\mathbf{P}_f \in \mathbb{R}^{F \times 1}$ to $\mathbf{P}'_f \in \mathbb{R}^{N \times F}$;
- 10: $\tilde{\mathbf{X}} = \mathbf{X} \odot \mathbf{P}'_f$;
- 11: Derive $\tilde{\mathbf{A}}$ via \mathbf{P}_e ;
- 12: $\tilde{\mathcal{G}} = (\tilde{\mathbf{A}}, \tilde{\mathbf{X}})$;
- 13: $\mathcal{L} = -\hat{I}(g_{\xi}(f_{\theta}(\mathcal{G})); g_{\xi}(f_{\theta}(\tilde{\mathcal{G}})))$;
- 14: $\mathcal{R} = \frac{\sum_{i=1}^F \mathbf{P}_f^{(i)}}{F} + \frac{\sum_{i=1}^M \mathbf{P}_e^{(i)}}{M}$;
- 15: /* Update the importance learner */
- 16: $\mathbf{W} \leftarrow \mathbf{W} - \alpha \nabla_{\mathbf{W}}(\mathcal{L} + \lambda * \mathcal{R})$;
- 17: $\phi \leftarrow \phi - \alpha \nabla_{\phi}(\mathcal{L} + \lambda * \mathcal{R})$;
- 18: $\psi \leftarrow \psi - \alpha \nabla_{\psi}(\mathcal{L} + \lambda * \mathcal{R})$;
- 19: /* Update the GNN encoder and projection head */
- 20: $\theta \leftarrow \theta + \alpha \nabla_{\theta}(\mathcal{L})$;
- 21: $\xi \leftarrow \xi + \alpha \nabla_{\xi}(\mathcal{L})$;
- 22: **end for**
- 23: $\mathbf{P}_{d,\varepsilon} = [p_{d,\varepsilon_1}, \dots, p_{d,\varepsilon_M}] = \min \left(\frac{\mathbf{P}_e^{\max} - \mathbf{P}_e}{\mathbf{P}_e^{\max} - \mathbf{P}_e^{\text{avg}}} \cdot p_{1,s1}, p_t \right)$;
- 24: Sample p_{ε_k} from Bernoulli($1 - p_{d,\varepsilon_k}$), $k = 1, \dots, M$;
- 25: Use p_{ε_k} , $k = 1, \dots, M$ to generate \mathbf{A}_1 ;
- 26: $\mathbf{P}_{d,f} = [p_{d,f_1}, \dots, p_{d,f_F}] = \min \left(\frac{\mathbf{P}_f^{\max} - \mathbf{P}_f}{\mathbf{P}_f^{\max} - \mathbf{P}_f^{\text{avg}}} \cdot p_{1,s2}, p_t \right)$;
- 27: Sample p_{f_k} from Bernoulli($1 - p_{d,f_k}$), $k = 1, \dots, F$;
- 28: Use p_{f_k} , $k = 1, \dots, F$ to generate \mathbf{X}_1 ;
- 29: Generate $\mathbf{A}_2, \mathbf{X}_2$ similarly;
- 30: **return** $\mathcal{G}_1 = (\mathbf{A}_1, \mathbf{X}_1)$; $\mathcal{G}_2 = (\mathbf{A}_2, \mathbf{X}_2)$;

Algorithm 2 View comparison module for node classification

Input: graph \mathcal{G} ; GNN encoder $f_\omega(\cdot)$; projection head $g_\rho(\cdot)$; mutual information estimator $\hat{I}(\cdot; \cdot)$.

HyperParams: learning rate β .

Ensure: trained $f_\omega(\cdot)$.

- 1: $(\mathcal{G}_1, \mathcal{G}_2) = \text{Algorithm 1}(\mathcal{G})$;
 - 2: **for** ep in #epochs **do**
 - 3: $\mathcal{L} = -\hat{I}(g_\rho(f_\omega(\mathcal{G}_1)); g_\rho(f_\omega(\mathcal{G}_2)))$;
 - 4: $\omega \leftarrow \omega - \beta \nabla_\omega \mathcal{L}$; $\rho \leftarrow \rho - \beta \nabla_\rho \mathcal{L}$;
 - 5: **end for**
 - 6: **return** trained $f_\omega(\cdot)$;
-

Algorithm 3 View comparison module for graph classification

Input: graph set $\{\mathcal{G}_i\}, i = 1, \dots, Q$; GNN encoder $f_\omega(\cdot)$; projection head $g_\rho(\cdot)$; mutual information estimator $\hat{I}(\cdot; \cdot)$.

HyperParams: learning rate β ; batch size B ; number of batches BN .

Ensure: trained $f_\omega(\cdot)$.

- 1: **for** $i = 1$ to BN **do**
 - 2: **for** \mathcal{G}_j in sampled minibatch $\{\mathcal{G}_j\}_{j=1}^B$ **do**
 - 3: $(\mathcal{G}_{i,(j,1)}, \mathcal{G}_{i,(j,2)}) = \text{Algorithm 1}(\mathcal{G}_{i,j})$;
 - 4: **end for**
 - 5: **end for**
 - 6: **for** ep in #epochs **do**
 - 7: **for** $i = 1$ to BN **do**
 - 8: $\mathcal{L} = -\frac{1}{B} \sum_{j=1}^B \hat{I}(g_\rho(f_\omega(\mathcal{G}_{i,(j,1)})); g_\rho(f_\omega(\mathcal{G}_{i,(j,2)})))$;
 - 9: $\omega \leftarrow \omega - \beta \nabla_\omega \mathcal{L}$; $\rho \leftarrow \rho - \beta \nabla_\rho \mathcal{L}$;
 - 10: **end for**
 - 11: **end for**
 - 12: **return** trained $f_\omega(\cdot)$;
-

746 **References**

- 747 J. Devlin, Bert: Pre-training of deep bidirectional transformers for language
748 understanding, arXiv preprint arXiv:1810.04805 (2018).
- 749 G. Zhou, S. Bu, T. Kirubarajan, Simultaneous spatiotemporal bias com-
750 pensation and data fusion for asynchronous multisensor systems, Chinese
751 Journal of Information Fusion 1 (2024) 16–32.
- 752 H. Fan, D. Lu, Y. Jiang, A. J. Lilienthal, Extraction of motion information
753 from occupancy grid map using keystone transform, Chinese Journal of
754 Information Fusion 1 (2024) 63–78.
- 755 W. Ju, Y. Gu, X. Luo, Y. Wang, H. Yuan, H. Zhong, M. Zhang, Unsu-
756 pervised graph-level representation learning with hierarchical contrasts,
757 Neural Networks 158 (2023) 359–368.
- 758 Y. Liu, Y. Zheng, D. Zhang, V. C. Lee, S. Pan, Beyond smoothing: Un-
759 supervised graph representation learning with edge heterophily discrimi-
760 nating, in: Proceedings of the AAAI conference on artificial intelligence,
761 volume 37, 2023, pp. 4516–4524.
- 762 Y. Mo, L. Peng, J. Xu, X. Shi, X. Zhu, Simple unsupervised graph repre-
763 sentation learning, in: Proceedings of the AAAI Conference on Artificial
764 Intelligence, volume 36, 2022, pp. 7797–7805.
- 765 S. Xiao, S. Wang, Y. Dai, W. Guo, Graph neural networks in node classifi-
766 cation: survey and evaluation, Machine Vision and Applications 33 (2022)
767 1–19.
- 768 F.-Y. Sun, J. Hoffmann, V. Verma, J. Tang, Infograph: Unsupervised and
769 semi-supervised graph-level representation learning via mutual information
770 maximization, arXiv preprint arXiv:1908.01000 (2019).
- 771 E. Pan, Z. Kang, Multi-view contrastive graph clustering, Advances in neural
772 information processing systems 34 (2021) 2148–2159.
- 773 M. Lin, T. Xiao, E. Dai, X. Zhang, S. Wang, Certifiably robust graph
774 contrastive learning, in: A. Oh, T. Naumann, A. Globerson, K. Saenko,
775 M. Hardt, S. Levine (Eds.), Advances in Neural Information Processing
776 Systems, volume 36, Curran Associates, Inc., 2023, pp. 17008–17037.

- 777 J. Chen, G. Kou, Attribute and structure preserving graph contrastive learn-
778 ing, in: Proceedings of the AAAI Conference on Artificial Intelligence,
779 volume 37, 2023, pp. 7024–7032.
- 780 X. Gong, C. Yang, C. Shi, Ma-gcl: Model augmentation tricks for graph
781 contrastive learning, in: Proceedings of the AAAI Conference on Artificial
782 Intelligence, volume 37, 2023, pp. 4284–4292.
- 783 Y. You, T. Chen, Y. Shen, Z. Wang, Graph contrastive learning auto-
784 mated, in: International Conference on Machine Learning, PMLR, 2021,
785 pp. 12121–12132.
- 786 K. Ding, Z. Xu, H. Tong, H. Liu, Data augmentation for deep graph learning:
787 A survey, ACM SIGKDD Explorations Newsletter 24 (2022) 61–77.
- 788 P. Mishra, A. Piktus, G. Goossen, F. Silvestri, Node masking: Making graph
789 neural networks generalize and scale better, CoRR abs/2001.07524 (2020).
790 arXiv:2001.07524.
- 791 P. Velickovic, W. Fedus, W. L. Hamilton, P. Liò, Y. Bengio, R. D. Hjelm,
792 Deep graph infomax., ICLR (Poster) 2 (2019) 4.
- 793 Y. Zhu, Y. Xu, F. Yu, Q. Liu, S. Wu, L. Wang, Graph contrastive learning
794 with adaptive augmentation, in: Proceedings of the Web Conference 2021,
795 2021, pp. 2069–2080.
- 796 L. Page, S. Brin, R. Motwani, T. Winograd, The PageRank citation ranking:
797 Bringing order to the web., Technical Report, Stanford InfoLab, 1999.
- 798 R. I. Kondor, J. Lafferty, Diffusion kernels on graphs and other discrete
799 structures, in: Proceedings of the 19th international conference on machine
800 learning, volume 2002, 2002, pp. 315–322.
- 801 F. Feng, X. He, J. Tang, T.-S. Chua, Graph adversarial training: Dynamically
802 regularizing based on graph structure, IEEE Transactions on Knowl-
803 edge and Data Engineering 33 (2019) 2493–2504.
- 804 L. Yang, L. Zhang, W. Yang, Graph adversarial self-supervised learning, Ad-
805 vances in Neural Information Processing Systems 34 (2021) 14887–14899.

- 806 Y. You, T. Chen, Y. Sui, T. Chen, Z. Wang, Y. Shen, Graph contrastive
807 learning with augmentations, *Advances in Neural Information Processing*
808 *Systems* 33 (2020a) 5812–5823.
- 809 Y. You, T. Chen, Z. Wang, Y. Shen, When does self-supervision help graph
810 convolutional networks?, in: *international conference on machine learning*,
811 PMLR, 2020b, pp. 10871–10880.
- 812 J. Park, H. Shim, E. Yang, Graph transplant: Node saliency-guided graph
813 mixup with local structure preservation, in: *Proceedings of the AAAI*
814 *Conference on Artificial Intelligence*, volume 36, 2022, pp. 7966–7974.
- 815 M. Jin, Y. Zheng, Y.-F. Li, C. Gong, C. Zhou, S. Pan, Multi-scale contrastive
816 siamese networks for self-supervised graph representation learning, *arXiv*
817 *preprint arXiv:2105.05682* (2021).
- 818 R. Linsker, Self-organization in a perceptual network, *Computer* 21 (1988)
819 105–117.
- 820 Y. Zhu, Y. Xu, F. Yu, Q. Liu, S. Wu, L. Wang, Deep graph contrastive
821 representation learning, *arXiv preprint arXiv:2006.04131* (2020).
- 822 S. Thakoor, C. Tallec, M. G. Azar, R. Munos, P. Veličković, M. Valko, Boot-
823 strapped representation learning on graphs, in: *ICLR 2021 Workshop on*
824 *Geometrical and Topological Representation Learning*, 2021.
- 825 M. Tschannen, J. Djolonga, P. K. Rubenstein, S. Gelly, M. Lucic, On mu-
826 tual information maximization for representation learning, *arXiv preprint*
827 *arXiv:1907.13625* (2019).
- 828 Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, P. Isola, What makes
829 for good views for contrastive learning?, *Advances in Neural Information*
830 *Processing Systems* 33 (2020) 6827–6839.
- 831 Y. Yin, Q. Wang, S. Huang, H. Xiong, X. Zhang, Autogcl: Automated
832 graph contrastive learning via learnable view generators, in: *Proceedings*
833 *of the AAAI Conference on Artificial Intelligence*, volume 36, 2022, pp.
834 8892–8900.
- 835 S. Suresh, P. Li, C. Hao, J. Neville, Adversarial graph augmentation to
836 improve graph contrastive learning, *Advances in Neural Information Pro-*
837 *cessing Systems* 34 (2021) 15920–15933.

- 838 W. Cheng, D. Xu, H. Chen, Information-aware graph contrastive learning,
839 2022. US Patent App. 17/728,071.
- 840 K. Yang, Y. Sun, J. Su, F. He, X. Tian, F. Huang, T. Zhou, D. Tao, Ad-
841 versarial auto-augment with label preservation: A representation learning
842 principle guided approach, *Advances in Neural Information Processing*
843 *Systems* 35 (2022) 22035–22048.
- 844 A. Ahmed, N. Shervashidze, S. Narayanamurthy, V. Josifovski, A. J. Smola,
845 Distributed large-scale natural graph factorization, in: *Proceedings of the*
846 *22nd international conference on World Wide Web*, 2013, pp. 37–48.
- 847 S. Cao, W. Lu, Q. Xu, Grarep: Learning graph representations with global
848 structural information, in: *Proceedings of the 24th ACM international on*
849 *conference on information and knowledge management*, 2015, pp. 891–900.
- 850 M. Ou, P. Cui, J. Pei, Z. Zhang, W. Zhu, Asymmetric transitivity pre-
851 serving graph embedding, in: *Proceedings of the 22nd ACM SIGKDD*
852 *international conference on Knowledge discovery and data mining*, 2016,
853 pp. 1105–1114.
- 854 B. Perozzi, R. Al-Rfou, S. Skiena, Deepwalk: Online learning of social rep-
855 resentations, in: *Proceedings of the 20th ACM SIGKDD international*
856 *conference on Knowledge discovery and data mining*, 2014, pp. 701–710.
- 857 A. Grover, J. Leskovec, node2vec: Scalable feature learning for networks,
858 in: *Proceedings of the 22nd ACM SIGKDD international conference on*
859 *Knowledge discovery and data mining*, 2016, pp. 855–864.
- 860 Z. Zhou, Y. Hu, Y. Zhang, J. Chen, H. Cai, Multiview deep graph infomax to
861 achieve unsupervised graph embedding, *IEEE Transactions on Cybernetics*
862 (2022).
- 863 E. Oyallon, Interferometric graph transform: a deep unsupervised graph
864 representation, in: *International Conference on Machine Learning*, PMLR,
865 2020, pp. 7434–7444.
- 866 R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman,
867 A. Trischler, Y. Bengio, Learning deep representations by mutual infor-
868 mation estimation and maximization, *arXiv preprint arXiv:1808.06670*
869 (2018).

- 870 N. Tishby, F. C. Pereira, W. Bialek, The information bottleneck method,
871 arXiv preprint physics/0004057 (2000).
- 872 Z. Goldfeld, Y. Polyanskiy, The information bottleneck problem and its
873 applications in machine learning, *IEEE Journal on Selected Areas in In-*
874 *formation Theory* 1 (2020) 19–38.
- 875 A. A. Alemi, I. Fischer, J. V. Dillon, K. Murphy, Deep variational information
876 bottleneck, arXiv preprint arXiv:1612.00410 (2016).
- 877 D. J. MacKay, *Information theory, inference and learning algorithms*, Cam-
878 *bridge university press*, 2003.
- 879 G. Karypis, V. Kumar, A fast and high quality multilevel scheme for parti-
880 tioning irregular graphs, *SIAM Journal on scientific Computing* 20 (1998)
881 359–392.
- 882 A. Nazi, W. Hang, A. Goldie, S. Ravi, A. Mirhoseini, Gap: Generalizable ap-
883 proximate graph partitioning framework, arXiv preprint arXiv:1903.00614
884 (2019).
- 885 T. N. Kipf, M. Welling, Semi-supervised classification with graph convolu-
886 tional networks, arXiv preprint arXiv:1609.02907 (2016).
- 887 X. Han, Z. Jiang, N. Liu, X. Hu, G-mixup: Graph data augmentation for
888 graph classification, in: *International Conference on Machine Learning*,
889 *PMLR*, 2022, pp. 8230–8248.
- 890 M. Zhang, Y. Chen, Link prediction based on graph neural networks, *Ad-*
891 *vances in neural information processing systems* 31 (2018).
- 892 S. Yang, B. Hu, Z. Zhang, W. Sun, Y. Wang, J. Zhou, H. Shan, Y. Cao, B. Ye,
893 Y. Fang, et al., Inductive link prediction with interactive structure learning
894 on attributed graph, in: *Machine Learning and Knowledge Discovery in*
895 *Databases. Research Track: European Conference, ECML PKDD 2021,*
896 *Bilbao, Spain, September 13–17, 2021, Proceedings, Part II* 21, Springer,
897 2021, pp. 383–398.
- 898 M. Welling, T. N. Kipf, Semi-supervised classification with graph convolu-
899 tional networks, in: *J. International Conference on Learning Representa-*
900 *tions (ICLR 2017)*, 2016.

- 901 C. J. Maddison, A. Mnih, Y. W. Teh, The concrete distribution: A
902 continuous relaxation of discrete random variables, arXiv preprint
903 arXiv:1611.00712 (2016).
- 904 A. v. d. Oord, Y. Li, O. Vinyals, Representation learning with contrastive
905 predictive coding, arXiv preprint arXiv:1807.03748 (2018).
- 906 K. Xu, W. Hu, J. Leskovec, S. Jegelka, How powerful are graph neural
907 networks?, arXiv preprint arXiv:1810.00826 (2018).
- 908 P. Mernyei, C. Cangea, Wiki-cs: A wikipedia-based benchmark for graph
909 neural networks, arXiv preprint arXiv:2007.02901 (2020).
- 910 O. Shchur, M. Mumme, A. Bojchevski, S. Günnemann, Pitfalls of graph
911 neural network evaluation, arXiv preprint arXiv:1811.05868 (2018).
- 912 C. Morris, N. M. Kriege, F. Bause, K. Kersting, P. Mutzel, M. Neumann,
913 Tudataset: A collection of benchmark datasets for learning with graphs,
914 arXiv preprint arXiv:2007.08663 (2020).
- 915 Z. Peng, W. Huang, M. Luo, Q. Zheng, Y. Rong, T. Xu, J. Huang, Graph
916 representation learning via graphical mutual information maximization,
917 in: Proceedings of The Web Conference 2020, 2020, pp. 259–270.
- 918 K. Hassani, A. H. Khasahmadi, Contrastive multi-view representation learn-
919 ing on graphs, in: International Conference on Machine Learning, PMLR,
920 2020, pp. 4116–4126.
- 921 W. Ju, Y. Gu, X. Luo, Y. Wang, H. Yuan, H. Zhong, M. Zhang, Unsu-
922 pervised graph-level representation learning with hierarchical contrasts,
923 Neural Networks 158 (2023) 359–368.
- 924 J. Wu, X. Chen, B. Shi, S. Li, K. Xu, Sega: Structural entropy guided
925 anchor view for graph contrastive learning, in: International Conference
926 on Machine Learning, PMLR, 2023, pp. 37293–37312.
- 927 H. Abdi, L. J. Williams, Principal component analysis, Wiley interdisci-
928 plinary reviews: computational statistics 2 (2010) 433–459.
- 929 L. Van der Maaten, G. Hinton, Visualizing data using t-sne., Journal of
930 machine learning research 9 (2008).

- 931 T. Chen, S. Bian, Y. Sun, Are powerful graph neural nets necessary? a
932 dissection on graph classification, arXiv preprint arXiv:1905.04579 (2019).
- 933 X. Glorot, Y. Bengio, Understanding the difficulty of training deep feed-
934 forward neural networks, in: Proceedings of the thirteenth international
935 conference on artificial intelligence and statistics, JMLR Workshop and
936 Conference Proceedings, 2010, pp. 249–256.
- 937 D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv
938 preprint arXiv:1412.6980 (2014).

Declaration of interests

- The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.
- The author is an Editorial Board Member/Editor-in-Chief/Associate Editor/Guest Editor for [*Journal name*] and was not involved in the editorial review or the decision to publish this article.
- The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Furao Shen reports financial support was provided by the STI 2030-Major Projects of China under Grant 2021ZD0201300, the National Science Foundation of China under Grant 62276127.