

# GDM: Dual Mixup for Graph Classification with Limited Supervision

Abdullah Alchihabi<sup>1</sup> and Yuhong Guo<sup>1,2</sup>

<sup>1</sup> School of Computer Science, Carleton University, Ottawa, Canada

<sup>2</sup> Canada CIFAR AI Chair, Amii, Canada

abdullahalchihabi@cmail.carleton.ca, yuhong.guo@carleton.ca

**Abstract.** Graph Neural Networks (GNNs) require a large number of labeled graph samples to obtain good performance on the graph classification task. The performance of GNNs degrades significantly as the number of labeled graph samples decreases. To reduce the annotation cost, it is therefore important to develop graph augmentation methods that can generate new graph instances to increase the size and diversity of the limited set of available labeled graph samples. In this work, we propose a novel mixup-based graph augmentation method, Graph Dual Mixup (GDM), that leverages both functional and structural information of the graph instances to generate new labeled graph samples. GDM employs a graph structural auto-encoder to learn structural embeddings of the graph samples, and then applies mixup to the structural information of the graphs in the learned structural embedding space and generates new graph structures from the mixup structural embeddings. As for the functional information, GDM applies mixup directly to the input node features of the graph samples to generate functional node feature information for new mixup graph instances. Jointly, the generated input node features and graph structures yield new graph samples which can supplement the set of original labeled graphs. Furthermore, we propose two novel Balanced Graph Sampling methods to enhance the balanced difficulty and diversity for the generated graph samples. Experimental results on the benchmark datasets demonstrate that our proposed method substantially outperforms the state-of-the-art graph augmentation methods when the labeled graphs are scarce.

**Keywords:** Graph Augmentation · Graph Classification · Limited Supervision.

## 1 Introduction

Graph Neural Networks (GNNs) have successfully tackled a wide range of graph related tasks such as node classification, knowledge graph completion, and graph classification. In particular, the graph classification task has been addressed using various GNN models such as Graph Convolution Networks (GCNs) [11], Graph Attention Networks (GATs) [17], and Graph Isomorphism Networks (GINs) [22]. The effectiveness of such GNN models can be attributed to their natural ability

to leverage both the functional information (nodes input features) and structural information (graph adjacency matrix) of graph data using message passing and message aggregation operations.

The success of GNNs in addressing the graph classification task nevertheless has been contingent on the availability of a large set of labeled graph samples, which induces a significant annotation burden in many domains where the labels are scarce and require substantial domain-expertise to generate. To tackle this problem, various graph augmentation methods have been proposed to increase the size and diversity of the training set by generating additional new graph samples. Most common graph augmentation methods, such as DropEdge [14], DropNode [25], and SoftEdge [5], involve perturbing the nodes, edges, or sub-graphs of a given graph sample to generate a new graph. However such methods assume that the employed graph-augmentation operations are label invariant, which is difficult to guarantee in many cases. Additionally, these methods use a single graph sample to generate new graph instances, which limits the diversity of the generated graphs. Although mixup-based augmentation methods have demonstrated tremendous success in improving the generalization capacity of deep neural networks on image-based [26] and text-based tasks [16], it remains an open challenge to apply mixup to graph-based tasks given the irregular, discrete and not well-aligned nature of graph data. Few works have proposed methods to adapt mixup to graph data, including G-Mixup [7] and M-Mixup [21]. However, these methods either are computationally expensive and need a relatively large number of graph samples to obtain good performance or generate new graph samples in the manifold space and offer limited improvement in performance.

In this work, we propose a novel mixup-based graph augmentation method named Graph Dual Mixup (GDM) for graph classification, which applies parallel mixup to the functional and structural information of the graph samples to generate new graph instances in the input space. Given the discrete nature of the graph structures, GDM employs a Graph Structural Auto-Encoder (GSAE) to learn a structural embedding of the graph nodes. It then applies mixup to the learned structural node embeddings of existing graphs to generate structural node embeddings for new mixup graph samples, which are subsequently used to produce the graph structures (i.e., adjacency matrices) of the new graph samples using the Graph Structural Decoder. Regarding the functional information, GDM applies mixup directly to the input node features of existing graphs to obtain the input node features of the corresponding mixup graph samples. The new graph instances generated through the parallel mixup over both the input features and graph structures are thereafter used to supplement the original set of labeled graph samples, reduce overfitting, and help GNNs generalize better with scarce graph labels. Furthermore, we propose two Balanced Graph Sampling methods to guide the mixup procedure to achieve balanced difficulty and diversity for the generated graph instances. We conduct comprehensive experiments on six graph classification benchmark datasets. The experimental results demonstrate that our proposed method substantially outperforms state-of-the-

art graph augmentation methods in the literature when the number of labeled graphs is limited.

## 2 Related Works

### 2.1 Graph Classification

Earlier works have addressed the graph classification task using graph-kernel based methods where the graph samples are decomposed into small subgraphs [8,15,23]. More recently, Graph Neural Networks (GNNs) have been successfully adopted in tackling the graph classification task. Many GNN models such as Graph Convolution Networks (GCNs) [11], Graph Attention Networks (GATs) [17], GraphSAGE [6] and Graph Isomorphism Networks (GINs) [22] have been shown to possess strong capacity to represent the graph data using message passing and message aggregation operations, and facilitate graph classification. Moreover, some works have developed novel graph readout methods to obtain discriminative graph-level representations from the node-level representation learned by various GNN models [24,1].

### 2.2 Graph Augmentation

Data augmentation methods play a crucial role in regularizing the training of deep models. Common graph augmentation methods are perturbation-based methods that augment graph samples by applying perturbations to graph nodes [25,9], edges [5,14], or subgraphs [20,13]. DropEdge randomly drops a number of edges from the graph structure during training [14]. SoftEdge selects a random subset of edges and assigns random weights to them to generate augmented graphs while preserving the connectivity patterns of the input graphs [5]. DropNode randomly deletes a subset of the nodes in the graph together with their connections to generate augmented graph samples [25]. GraphCrop augments the graphs with sub-structure deletion, which motivates GNNs to learn a robust global-view of the graph samples [20]. Graph Transplant uses subgraph transplantation to augment graphs where node saliency is used to select the transplanted subgraphs [13]. These methods however operate under the strong assumption that the applied graph perturbations are label-invariant insofar the augmented graph shares the same ground-truth label as the original graph. Such an assumption is hard to guarantee in many cases. Meanwhile, although there has been tremendous success of Mixup-based methods in regularizing deep models in domains where the data is regular, well-aligned and continuous such as images [26] and text [16], few works have attempted to adapt mixup to graph data. M-Mixup applies mixup to the graph-level representation in the manifold space learned by GNNs in a similar way to manifold mixup [21]. G-Mixup performs mixup to the graphons of different classes which are learned from the graph samples, and generates augmented graphs by sampling from the mixed graphons [7]. GraphMix is a node-level augmentation method where manifold mixup is applied to a fully-connected network that is trained jointly with a GNN [18]. Further details on graph augmentation methods can be found in [27].

### 3 Method

#### 3.1 Problem Setup

We consider the following graph classification setting. The input is a set of  $N$  labeled graphs:  $\mathcal{G} = \{(G_1, \mathbf{y}_1), \dots, (G_N, \mathbf{y}_N)\}$ . Each graph  $G$  is made up of a pair  $(V, E)$ , where  $V$  is the set of graph nodes with size  $|V| = n$  and  $E$  is the set of edges.  $E$  is represented by an adjacency matrix  $A$  of size  $n \times n$ . The adjacency matrix can have either binary or weighted values, be symmetric (in the case of undirected graphs) or asymmetric (in the case of directed graphs). Each node in the graph  $G$  is associated with a corresponding feature vector of size  $d$ . The feature vectors of all the nodes in the graph are represented by an input feature matrix  $X \in \mathbb{R}^{n \times d}$ . The graphs in the training set  $\mathcal{G}$  may potentially have different sizes (different number of nodes), while the feature vectors of the nodes of all graphs have the same size  $d$ . The graph label vector  $\mathbf{y}$  is a one-hot label indicator vector of size  $C$ , where  $C$  is the number of classes.

#### 3.2 Graph Classification

GNNs address the graph classification task by utilizing both the graph adjacency matrix and the input node features, which correspond to the structural and functional information of the graphs, respectively. GNN models in the literature are commonly made up of three components: a node representation learning function  $f_\theta$ , a graph readout function, and a graph classification function  $g_\phi$ . The node representation function  $f_\theta$  typically consists of multiple (e.g.,  $L$ ) GNN layers, each of which performs message propagation and message aggregation at the node level to learn new node embedding as follows:

$$\mathbf{h}_u^l = \text{AGGREGATE}(\mathbf{h}_u^{l-1}, \mathbf{h}_v^{l-1} | v \in \mathcal{N}(u), \theta^l) \quad (1)$$

where  $\mathbf{h}_u^l \in \mathbb{R}^{d_l \times 1}$  is the learned embedding of node  $u$  with size  $d_l$  at layer  $l$ ,  $\mathcal{N}(u)$  is the set of neighboring nodes of node  $u$ ,  $\theta^l$  is the learnable parameters of the  $l$ -th GNN layer, and AGGREGATE is the message aggregation function which can be any permutation invariant function (sum, average, max, etc.). The initial node embedding  $\mathbf{h}_u^0$  is the input node feature vector  $\mathbf{x}_u$ . The graph readout function is a permutation-invariant function used to obtain the graph-level embedding from the learned node-level embedding as follows:

$$\mathbf{h}_G = \text{READOUT}(\mathbf{h}_u^L | u \in V) \quad (2)$$

where  $\mathbf{h}_u^L \in \mathbb{R}^{d_L \times 1}$  is the embedding of node  $u$  obtained from the top layer  $L$  of  $f_\theta$  and  $\mathbf{h}_G \in \mathbb{R}^{d_G \times 1}$  is the graph-level embedding. The graph classification function  $g_\phi$  takes the graph-level embedding  $\mathbf{h}_G$  as input to produce the predicted class probability vector for the given graph  $G$  as follows:

$$\mathbf{p}_G = g(\mathbf{h}_G | \phi). \quad (3)$$

All the components are trained end-to-end by minimizing the following cross-entropy loss over the labeled graphs in the training set:

$$\mathcal{L} = \sum_{G \in \mathcal{G}} \ell(\mathbf{p}_G, \mathbf{y}_G) \quad (4)$$

where  $\ell(\cdot, \cdot)$  is the cross-entropy loss function,  $\mathbf{p}_G$  and  $\mathbf{y}_G$  are the predicted class probability vector and the ground-truth label indicator vector for graph  $G$ , respectively.

### 3.3 Mixup

Mixup is an interpolation-based augmentation method that has demonstrated significant success in reducing overfitting and improving the generalization of deep neural networks [26,16]. Mixup generates augmented training samples  $(\tilde{x}, \tilde{y})$  by applying linear interpolation between a randomly sampled pair of input instances and their corresponding labels as follows:

$$\tilde{x} = \lambda x_i + (1 - \lambda)x_j, \quad \tilde{y} = \lambda y_i + (1 - \lambda)y_j \quad (5)$$

where  $\lambda$  is a scalar mixing coefficient sampled from a Beta distribution  $\text{Beta}(\alpha, \beta)$  with hyper-parameters  $\alpha$  and  $\beta$ .  $(\tilde{x}, \tilde{y})$  is the new sample generated by mixing the input labeled samples  $(x_i, y_i)$  and  $(x_j, y_j)$ . Mixup can be readily applied to any classification task where the input data is regular, continuous and well-aligned such as images, text and time-series data. However, mixup cannot be applied directly to graph data given that: (1) graph data is irregular where different graphs may potentially have different sizes (different number of nodes). (2) graphs do not have a natural-ordering of their nodes, therefore aligning a pair of graphs is a non-trivial task. (3) graph structures may be discrete where the edges are binary whereas mixup generates continuous samples. Therefore, it is important to develop new methods that adapt mixup to the discrete, irregular and not well-aligned graph data.

### 3.4 Graph Dual Mixup

In this section, we introduce our proposed Graph Dual Mixup (GDM) method which generates new graph samples by applying parallel structural (i.e., structure-based) mixup and functional (i.e., feature-based) mixup over each selected pair of existing graph samples. In particular, GDM employs a Graph Structural Auto-Encoder (GSAE) to learn a structural embedding of the graph nodes based on the adjacency matrix. The structural mixup is then applied on the structural node embeddings of the input pair of graphs to produce a new set of node embeddings, which is used to generate the adjacency matrix (i.e., graph structure) of the mixup graph sample using the Graph Structural Decoder of the GSAE. As for the functional information encoded with node features, GDM applies mixup directly to the input node features to obtain the node features of the generated mixup graph sample. In the remainder of this section, we elaborate on the dual mixup procedure of this GDM methodology.

**Structural Graph Node Representation Learning** Given the discrete nature of graph structures, mixup cannot be directly applied to the structures of a pair of graphs (represented by their corresponding adjacency matrices) to generate a new graph structure. Therefore, we propose to employ a Graph Structural Auto-Encoder (GSAE) to learn a structural embedding of the graph nodes and support mixup in the learned structural embedding space. This allows us to evade the difficulties associated with applying mixup to the original graph structures. GSAE is made up of a structural encoder  $\mathcal{E}_s$  and a structural decoder  $\mathcal{D}_s$ . The structural encoder  $\mathcal{E}_s$  consists of multiple GNN layers that learn the structural node embeddings by propagating and aggregating messages across the graph structure, where the messages reflect solely the structural information of the nodes. The goal is to learn a structural embedding of all the nodes in the graph that would enable us to reconstruct the graph adjacency matrix. Specifically, for a given graph sample  $G = (X, A)$ ,  $\mathcal{E}_s$  takes the adjacency matrix  $A \in \mathbb{R}^{n \times n}$  and the node degree matrix  $D \in \mathbb{R}^{n \times n}$  (represent the initial node structural features) computed from  $A$  as input to learn the structural node embeddings as follows:

$$H_s = \mathcal{E}_s(D, A), \quad \text{where } D[i, i] = \sum_j A[i, j], \quad (6)$$

where the node degree matrix  $D$  is an identity matrix whose main diagonal values correspond to the degrees of the associated nodes;  $H_s \in \mathbb{R}^{n \times d_s}$  is the learned structural embedding of the nodes in the graph with size  $d_s$ .  $H_s$  holds solely the structural information of all the nodes in the graph, from which one can reconstruct the connections/edges between the nodes and therefore the original adjacency matrix  $A$  using the structural decoder  $\mathcal{D}_s$  of the GSAE. In particular, we adopt a simple inner product similarity based decoder as the structural decoder  $\mathcal{D}_s$ , which takes the learned structural node embeddings as input to reconstruct the graph adjacency matrix  $A$  as follows:

$$\hat{A} = \mathcal{D}_s(H_s) = \sigma(H_s H_s^T) \quad (7)$$

where  $\sigma$  is the sigmoid activation function and  $\hat{A} \in \mathbb{R}^{n \times n}$  is the decoded/reconstructed adjacency matrix. The GSAE is trained end-to-end to minimize the following graph structure reconstruction loss:

$$\mathcal{L}_{\text{re}}^s = - \sum_{G \in \mathcal{G}} \left[ \sum_{(i,j) \in E_G} \log(\hat{A}_G[i, j]) + \sum_{(i,j) \in S_G^{\text{neg}}} \log(1 - \hat{A}_G[i, j]) \right] \quad (8)$$

where  $E_G$  is the set of edges for graph  $G$  and  $S_G^{\text{neg}}$  is the set of randomly sampled negative edges of graph  $G$  (i.e. edges that do not exist in the original graph). It is important to note that GSAE does not access/use the input node features (functional graph information) as it replaces the input node features with the corresponding node degrees calculated from the adjacency matrix. GSAE also does not make use of the graph class labels as it is learned in a completely self-supervised/unsupervised fashion.

**Graph Generation via Dual Mixup** After training the GSAE, our proposed Graph Dual Mixup is ready to apply Structural Mixup and Functional Mixup to the structural and functional information of the graphs respectively to generate new graph samples. To achieve that, for a given pair of graphs and their corresponding label vectors  $(G_i, \mathbf{y}_i)$  and  $(G_j, \mathbf{y}_j)$ , where the two graphs are made up of input node feature matrices and graph adjacency matrices such as  $G_i = (X_i, A_i)$  and  $G_j = (X_j, A_j)$ , GDM randomly aligns the nodes of the graph pair. When  $G_i$  and  $G_j$  have different sizes ( $n_i \neq n_j$ ), we pad the input node feature matrix and adjacency matrix of the smaller graph with zeros to match the size of the larger graph. Then we apply functional mixup directly to the input node features and the label vectors of the graph pair to generate the node features of the new graph sample  $\tilde{G}$  and its corresponding label vector  $\tilde{\mathbf{y}}$  as follows:

$$\tilde{X} = \lambda X_i + (1 - \lambda)X_j, \quad \tilde{\mathbf{y}} = \lambda \mathbf{y}_i + (1 - \lambda)\mathbf{y}_j \quad (9)$$

To obtain the structural information of the generated new graph sample  $\tilde{G}$ , GDM applies structural mixup in the structural embedding space learned by the GSAE as follows:

$$\tilde{H}_s = \lambda \mathcal{E}_s(D_i, A_i) + (1 - \lambda) \mathcal{E}_s(D_j, A_j) \quad (10)$$

where  $D_i$  and  $D_j$  are the degree matrices of  $G_i$  and  $G_j$ , respectively;  $\tilde{H}_s \in \mathbb{R}^{\max(n_i, n_j) \times d_s}$  is the structural node embedding matrix of the generated graph  $\tilde{G}$ . The graph structural decoder is then used to reconstruct the adjacency matrix of graph  $\tilde{G}$  from the mixed structural node embeddings:

$$\tilde{A} = \mathcal{D}_s(\tilde{H}_s) = \sigma(\tilde{H}_s \tilde{H}_s^T) \quad (11)$$

The obtained matrix  $\tilde{A} \in \mathbb{R}^{\max(n_i, n_j) \times \max(n_i, n_j)}$  is a weighted adjacency matrix with edge weights between 0 and 1. In order to filter out the noise in the edge weights and sparsify the structure of generated graph sample, we prune the adjacency matrix by dropping off the weak edges with weights smaller than a pre-defined threshold  $\epsilon$  as follows:

$$\tilde{A}[i, j] = \begin{cases} \tilde{A}[i, j], & \text{if } \tilde{A}[i, j] \geq \epsilon \\ 0, & \text{otherwise.} \end{cases} \quad (12)$$

Moreover, in order for the structure of the generated graph sample  $\tilde{G}$  to match the structural properties of the original graph samples, we post-process  $\tilde{A}$  accordingly. In the case that the original graph samples have weighted edges, no post-processing is required. As for the case of the original graph samples being unweighted/binary graphs, we binarize  $\tilde{A}$  by replacing all its non-zero values with value 1 as follows:

$$\tilde{A}[i, j] = \begin{cases} 1, & \text{if } \tilde{A}[i, j] > 0 \\ 0, & \text{otherwise} \end{cases} \quad (13)$$

In this manner, we obtain a new generated graph  $\tilde{G}$  with its mixup node features  $\tilde{X}$ , adjacency matrix  $\tilde{A}$  and label vector  $\tilde{\mathbf{y}}$ .

### 3.5 Balanced Graph Sampling

Given the limited number of available labeled graph instances, randomly sampling pairs of graphs to generate new graph instances might be inadequate for improving model generalization and reducing overfitting as random sampling does not take the difficulty or diversity of the generated graph instances into consideration. Therefore, we propose two novel Balanced Graph Sampling methods to enhance the diversity and balanced difficulty of the generated graph samples. The proposed methods can separately: (1) generate low difficulty graphs by applying GDM to randomly sampled pairs of low difficulty graphs; (2) generate medium difficulty graphs by applying GDM to mix randomly sampled low difficulty graphs with high difficulty graphs; and (3) generate high difficulty graphs by applying GDM to randomly sampled pairs of high difficulty graphs. The advantage of balanced graph sampling over random sampling is that it guarantees that the generated graph samples have 3 subsets with *equal sizes*: a low difficulty subset, a medium difficulty subset, and a high difficulty subset.

To achieve that, we need to assess/estimate the difficulty level of the original graph instances. This is accomplished by pre-training a GNN model on the original set of labeled graph instances to minimize the classification loss shown in Eq.(1)—Eq.(4). Then the pre-trained GNN model is used to evaluate the difficulty level of each graph  $G$  based on its predicted class probability vector  $\mathbf{p}_G$ . The first balanced graph sampling method is an Accuracy-based method (Acc), which determines the level of difficulty for graph  $G$  based on the accuracy/correctness of its predicted class label:

$$\text{Diff}_{\text{Acc}}(G) = \begin{cases} \text{low}, & \text{if } \operatorname{argmax} \mathbf{p}_G = \operatorname{argmax} \mathbf{y}_G \\ \text{high}, & \text{otherwise} \end{cases} \quad (14)$$

The second balanced graph sampling method is an Uncertainty-based method (Unc), which uses the uncertainty/entropy of the model prediction on a sample graph  $G$  to determine its level of difficulty. In particular, we sort the graphs from the training set  $\mathcal{G}$  based on the entropy of their corresponding predicted class probability vectors, then consider the graphs with the lowest half of entropy scores to be low difficulty ones while the other half of the graphs are taken as high difficulty ones:

$$\text{Diff}_{\text{Unc}}(G) = \begin{cases} \text{low}, & \text{if } \text{Ent}(\mathbf{p}_G) \leq \text{Med}(\{\text{Ent}(\mathbf{p}_1), \dots, \text{Ent}(\mathbf{p}_N)\}) \\ \text{high}, & \text{otherwise} \end{cases} \quad (15)$$

where  $\text{Ent}(\cdot)$  is the entropy function and  $\text{Med}(\cdot)$  is the median function. Therefore, GDM can be applied with Accuracy-based Balanced Graph Sampling (GDM Acc) or Uncertainty-based Balanced Graph Sampling (GDM Unc) to generate a new set of diverse graph samples  $\mathcal{G}_{\text{GDM}}$  with balanced difficulty.

**Algorithm 1** Augmentation and Training Procedure of Graph Dual Mixup

---

**Input:** Graph set  $\mathcal{G}$ ; hyper-parameters  $\alpha, \beta, \epsilon, \lambda_{\text{GDM}}$   
**Output:** Learned model parameters  $\theta, \phi$   
Pre-train a GNN Model on  $\mathcal{G}$  to determine the graph difficulty levels  
Train GSAE on  $\mathcal{G}$  using Eq.(6), (7), (8).  
 $\mathcal{G}_{\text{low}}$  = Generate low difficulty samples with GDM  
 $\mathcal{G}_{\text{med}}$  = Generate medium difficulty samples with GDM  
 $\mathcal{G}_{\text{high}}$  = Generate high difficulty samples with GDM  
 $\mathcal{G}_{\text{GDM}} = \mathcal{G}_{\text{low}} \cup \mathcal{G}_{\text{med}} \cup \mathcal{G}_{\text{high}}$   
Train the final GNN Model on  $\mathcal{G}$  and  $\mathcal{G}_{\text{GDM}}$  using Eq.(1), (2), (3), and (16).

---

**3.6 Augmented Training Procedure**

The combination of Balanced Graph Sampling and Graph Dual Mixup generates a diverse set of new graph instances, which can supplement the limited number of original labeled graph samples. Finally, we train the GNN model using the original graph set  $\mathcal{G}$  and the generated graph set  $\mathcal{G}_{\text{GDM}}$  by minimizing the following loss function:

$$\mathcal{L}_{\text{total}} = \sum_{G \in \mathcal{G}} \ell_{CE}(\mathbf{p}_G, \mathbf{y}_G) + \lambda_{\text{GDM}} \sum_{G \in \mathcal{G}_{\text{GDM}}} \ell_{CE}(\mathbf{p}_G, \tilde{\mathbf{y}}_G) \quad (16)$$

where  $\lambda_{\text{GDM}}$  is a trade-off hyper-parameter controlling the contribution of the generated graph set  $\mathcal{G}_{\text{GDM}}$ . An overview of the graph augmentation process and the GNN augmented training procedure is presented in Algorithm 1.

**4 Experiments****4.1 Experimental Setup**

**Datasets & Baselines** We evaluate our proposed method on 6 graph classification benchmark datasets from the TUDatasets [12], including 3 chemical datasets and 3 social datasets. The chemical datasets are D&D [3], Proteins [2] and NCI1 [19], while the social datasets are IMDB-Binary, IMDB-Multi and Reddit-5K [23]. We employ the same 10-fold train/validation/test split provided by [4]. We apply our proposed Graph Dual Mixup on the Graph Convolution Network (GCN) baseline [11] and compare our proposed method against 5 other graph augmentation methods from the literature: DropNode [25], DropEdge [14], M-Mixup [21], SoftEdge [5] and G-Mixup [7].

**Implementation Details** The node representation function  $f_{\theta}$  of the GNN model is made up of 4 message passing layers, followed by Global Mean Pooling as the Readout function. The graph classification function  $g_{\phi}$  is made up of 2 fully connected layers followed by a softmax function. Each message passing layer and fully connected layer is followed by a Rectified Linear Unit (ReLU)

Table 1: Mean classification accuracy (standard deviation is within brackets) on 6 graph classification benchmark datasets with 10 labeled graphs per class.

Dataset	Proteins	NCI1	D&D	IMDB-B	IMDB-M	Reddit
GCN	59.3 <sub>(6.8)</sub>	51.0 <sub>(1.6)</sub>	59.5 <sub>(2.7)</sub>	54.5 <sub>(3.9)</sub>	36.9 <sub>(3.7)</sub>	25.1 <sub>(5.1)</sub>
DropNode	61.0 <sub>(8.5)</sub>	52.9 <sub>(3.4)</sub>	62.1 <sub>(2.9)</sub>	59.0 <sub>(5.7)</sub>	36.9 <sub>(4.6)</sub>	30.8 <sub>(8.4)</sub>
DropEdge	59.4 <sub>(5.8)</sub>	53.1 <sub>(3.7)</sub>	62.6 <sub>(4.5)</sub>	57.6 <sub>(5.5)</sub>	37.2 <sub>(4.1)</sub>	26.7 <sub>(8.4)</sub>
SoftEdge	58.9 <sub>(7.2)</sub>	52.0 <sub>(3.2)</sub>	59.5 <sub>(2.4)</sub>	55.3 <sub>(6.6)</sub>	36.2 <sub>(3.0)</sub>	25.0 <sub>(4.9)</sub>
M-Mixup	59.0 <sub>(7.2)</sub>	51.9 <sub>(3.3)</sub>	59.1 <sub>(5.3)</sub>	57.1 <sub>(6.4)</sub>	37.4 <sub>(5.3)</sub>	23.0 <sub>(2.8)</sub>
G-Mixup	60.8 <sub>(2.1)</sub>	51.8 <sub>(3.2)</sub>	58.7 <sub>(4.2)</sub>	55.1 <sub>(8.5)</sub>	36.9 <sub>(4.3)</sub>	24.1 <sub>(7.3)</sub>
GDM Acc	<b>66.0</b> <sub>(5.3)</sub>	<b>57.5</b> <sub>(2.6)</sub>	62.1 <sub>(3.7)</sub>	<b>61.3</b> <sub>(6.7)</sub>	<b>40.9</b> <sub>(5.4)</sub>	<b>36.3</b> <sub>(8.0)</sub>
GDM Unc	65.1 <sub>(6.1)</sub>	56.8 <sub>(3.9)</sub>	<b>64.0</b> <sub>(4.2)</sub>	61.0 <sub>(7.0)</sub>	39.8 <sub>(5.5)</sub>	34.9 <sub>(9.1)</sub>

activation function. The structural encoder  $\mathcal{E}_s$  is made up of 2 GCN message passing layers. The GNN model is pre-trained on the original graph set for 100 epochs and subsequently trained on the original graph set and augmented graph set for 800 epochs, both using the Adam optimizer with learning rate of 1e-2. The Graph Structural Auto-Encoder is trained for 200 epochs using the Adam optimizer with learning rate of 1e-2. The loss trade-off hyperparameter  $\lambda_{\text{GDM}}$  and the weak edge pruning threshold  $\epsilon$  take values 1 and 0.1, respectively. The mixing scalar coefficient  $\lambda$  is sampled from distribution  $\text{Beta}(\alpha, \beta)$  with hyperparameters  $\alpha = \beta = 1.0$ . We use a dropout rate of 0.25 for SoftEdge, DropNode and DropEdge. For G-Mixup, we use the same hyper-parameters reported in [7].

## 4.2 Comparison Results

We investigate the performance of our proposed GDM with limited numbers of labeled graphs. We aim to use a small number of labeled graphs per class, e.g.,  $\{2, 3, 5, 10, 25, 50\}$ , as the training set. To achieve that, we randomly sampled graphs from the training set of each fold in the 10-fold split provided [4] to match the desired label rates. For each label rate, we repeat our experiments 3 times on all the 10-folds and average the test accuracy over all folds and all runs. We evaluate GDM in combination with the proposed two Balanced Graph Sampling methods to obtain: (1) ‘‘GDM Acc’’, where GDM is applied with the Accuracy-based Balanced Graph Sampling; and (2) ‘‘GDM Unc’’, where GDM is applied with Uncertainty-based Balanced Graph Sampling. We report the obtained test accuracy results with 10 labeled graphs per class in Table 1, while the test accuracy results for all label rates are presented in Figure 1.

The results in Table 1 clearly demonstrate that both variants of our proposed GDM greatly outperform the underlying GCN baseline and the other 5 graph augmentation methods across all 6 datasets. GDM improves the performance of the underlying GCN baseline by 6.7%, 7.5%, 6.8% and 11.2% on the Proteins, NCI1, IMDB-Binary and Reddit-5K datasets, respectively. The performance gain over the other graph augmentation methods is also notable, exceeding 5%, 4.4% and 6.3% on Proteins, NCI1 and Reddit-5K, respectively.

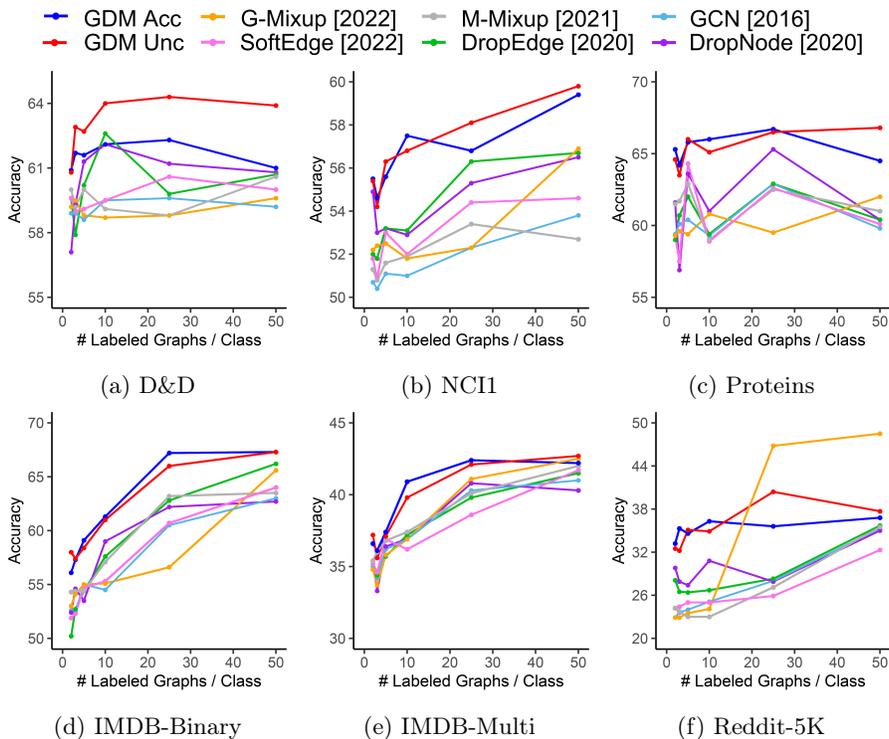


Fig. 1: Mean classification accuracy on 6 graph classification benchmark datasets with few labeled graphs per class (2, 3, 5, 10, 25, 50).

Moreover, Figure 1 clearly shows that our proposed GDM consistently outperforms the GCN baseline and the 5 comparison graph augmentation methods on 5 datasets across almost all label rates. Only in the case of the Reddit-5K dataset with label rates of larger than 25 labeled graphs per class, G-Mixup outperforms our proposed method. Nevertheless, GDM consistently improves the performance of the underlying GCN baseline across all the label rates on all the datasets, achieving performance gains over 6%, 5%, 5% and 11% on Proteins, NCI1, IMDB-Binary and Reddit-5K, respectively, in the case of 2 labeled graphs per class. Furthermore, GDM yields remarkable performance gains over the other graph augmentation methods, exceeding 4% on Proteins, Reddit-5K and IMDB-Binary in the case of 2 labeled graphs per class. This highlights the superior performance of the proposed GDM over the existing state-of-the-art graph augmentation methods for graph classification with limited supervision.

### 4.3 Ablation Study

**Impact of Balanced Graph Sampling** We conduct an ablation study to investigate the impact of our balanced graph sampling methods on the pro-

Table 2: Ablation study results on the impact of Balanced Graph Sampling in terms of mean classification accuracy (standard deviation is within brackets) with a few labeled graphs per class (2, 3, 5, 10).

	D&D				IMDB-Multi			
	2	3	5	10	2	3	5	10
GDM Rand	59.6 <sub>(5.3)</sub>	61.2 <sub>(3.9)</sub>	61.5 <sub>(3.5)</sub>	63.0 <sub>(4.5)</sub>	35.4 <sub>(3.8)</sub>	35.5 <sub>(4.0)</sub>	36.4 <sub>(5.5)</sub>	39.7 <sub>(5.3)</sub>
GDM Acc	<b>61.0</b> <sub>(2.6)</sub>	<b>61.7</b> <sub>(3.3)</sub>	<b>61.6</b> <sub>(3.5)</sub>	62.1 <sub>(3.7)</sub>	<b>36.6</b> <sub>(4.9)</sub>	<b>36.1</b> <sub>(3.8)</sub>	37.4 <sub>(4.5)</sub>	<b>40.9</b> <sub>(5.4)</sub>
w/o Low Diff	59.2 <sub>(1.4)</sub>	57.7 <sub>(5.7)</sub>	59.2 <sub>(2.1)</sub>	58.2 <sub>(4.0)</sub>	35.0 <sub>(3.6)</sub>	34.4 <sub>(2.8)</sub>	34.5 <sub>(1.8)</sub>	36.5 <sub>(3.1)</sub>
w/o Med Diff	60.2 <sub>(3.9)</sub>	61.6 <sub>(4.1)</sub>	59.9 <sub>(2.1)</sub>	59.6 <sub>(3.1)</sub>	34.9 <sub>(3.4)</sub>	35.4 <sub>(4.7)</sub>	38.3 <sub>(4.0)</sub>	39.9 <sub>(5.8)</sub>
w/o High Diff	60.1 <sub>(3.2)</sub>	60.6 <sub>(4.4)</sub>	60.6 <sub>(2.7)</sub>	61.3 <sub>(4.1)</sub>	33.8 <sub>(3.2)</sub>	34.4 <sub>(3.9)</sub>	<b>38.4</b> <sub>(4.3)</sub>	39.0 <sub>(5.5)</sub>
GDM Unc	<b>60.8</b> <sub>(2.8)</sub>	<b>62.9</b> <sub>(3.7)</sub>	<b>62.7</b> <sub>(3.2)</sub>	<b>64.0</b> <sub>(4.2)</sub>	<b>37.2</b> <sub>(4.9)</sub>	<b>35.6</b> <sub>(3.6)</sub>	<b>37.1</b> <sub>(5.4)</sub>	<b>39.8</b> <sub>(5.5)</sub>
w/o Low Diff	59.3 <sub>(1.0)</sub>	58.8 <sub>(1.0)</sub>	58.9 <sub>(1.0)</sub>	59.7 <sub>(1.8)</sub>	35.8 <sub>(2.9)</sub>	34.7 <sub>(3.7)</sub>	34.4 <sub>(2.7)</sub>	37.0 <sub>(3.4)</sub>
w/o Med Diff	<b>60.8</b> <sub>(2.9)</sub>	60.4 <sub>(4.3)</sub>	<b>62.7</b> <sub>(3.3)</sub>	62.8 <sub>(4.5)</sub>	35.3 <sub>(3.6)</sub>	34.8 <sub>(3.4)</sub>	36.6 <sub>(4.9)</sub>	38.8 <sub>(4.6)</sub>
w/o High Diff	60.4 <sub>(2.6)</sub>	61.5 <sub>(3.7)</sub>	62.0 <sub>(4.2)</sub>	63.2 <sub>(3.4)</sub>	35.7 <sub>(3.2)</sub>	35.0 <sub>(3.1)</sub>	<b>37.1</b> <sub>(4.0)</sub>	39.6 <sub>(3.8)</sub>

posed GDM method. Specifically, we consider four variants of the balanced graph sampling: (1) w/o Low Diff: we do not generate low difficulty samples. (2) w/o Med Diff: we do not generate medium difficulty samples. (3) w/o High Diff: we do not generate high difficulty samples. (4) GDM Rand: we drop the proposed balanced graph sampling method and use random sampling for mixup. We evaluate the first three variants using both the GDM Acc and GDM Unc methods of balanced graph sampling. The comparison results with different label rates—{2, 3, 5, 10} labeled graphs per class—on the D&D and IMDB-Multi datasets are reported in Table 2.

From Table 2, we can see that all variants have a performance drop from the full balanced graph sampling on both datasets with almost all label rates for both GDM Acc and GDM Unc. The w/o Low Diff variant produces the most notable performance degradation, which can be attributed to the GNN models’ needs for low difficulty and confident samples to improve generalization and prevent underfitting when learning with very low label rates. The w/o Med Diff and w/o High Diff variants also suffer performance degradations, indicating the importance of medium difficulty and high difficulty samples for inducing better generalization and reducing overfitting. Additionally, the GDM Rand variant also demonstrates notable performance drops compared to both GDM Acc and GDM Unc with almost all label rates, which highlights the importance of ensuring the diversity and balanced difficulty of the generated graph samples. These results validate the contribution of each component in balanced graph sampling.

**Impact of Graph Structural Auto-Encoder** We further conduct an ablation study to investigate the impact of the Graph Structural Auto-Encoder on the proposed GDM. Specifically, we compare our proposed GSAE with a Variational Graph Structural Auto-Encoder (VGSAE). The VGSAE learns the parameters of a Gaussian distribution (mean and variance) to represent the underlying structure of the graph [10]. The comparison results with different label

Table 3: Ablation study results on the impact of Graph Structural Auto-encoder in terms of mean classification accuracy (standard deviation is within brackets).

	Proteins				IMDB-Binary			
	2	3	5	10	2	3	5	10
GCN	59.4 <sub>(6.7)</sub>	60.1 <sub>(1.0)</sub>	60.4 <sub>(4.6)</sub>	59.3 <sub>(6.8)</sub>	52.6 <sub>(2.7)</sub>	52.5 <sub>(2.5)</sub>	55.0 <sub>(6.4)</sub>	54.5 <sub>(3.9)</sub>
GSAE Acc	<b>65.3</b> <sub>(5.5)</sub>	<b>64.2</b> <sub>(7.9)</sub>	<b>65.8</b> <sub>(6.2)</sub>	<b>66.0</b> <sub>(5.3)</sub>	<b>56.1</b> <sub>(4.6)</sub>	57.3 <sub>(6.8)</sub>	<b>59.1</b> <sub>(6.2)</sub>	<b>61.3</b> <sub>(6.7)</sub>
VGSAE Acc	64.6 <sub>(4.2)</sub>	63.4 <sub>(5.1)</sub>	65.6 <sub>(6.0)</sub>	65.3 <sub>(7.0)</sub>	55.9 <sub>(4.0)</sub>	<b>58.7</b> <sub>(5.2)</sub>	57.8 <sub>(6.5)</sub>	59.7 <sub>(7.4)</sub>
GSAE Unc	<b>64.6</b> <sub>(4.0)</sub>	63.5 <sub>(7.2)</sub>	<b>66.0</b> <sub>(5.3)</sub>	<b>65.1</b> <sub>(6.1)</sub>	<b>58.0</b> <sub>(6.0)</sub>	57.4 <sub>(7.1)</sub>	<b>58.4</b> <sub>(6.1)</sub>	<b>61.0</b> <sub>(7.0)</sub>
VGSAE Unc	61.8 <sub>(9.4)</sub>	<b>63.7</b> <sub>(8.4)</sub>	65.2 <sub>(5.2)</sub>	63.7 <sub>(5.7)</sub>	55.2 <sub>(3.9)</sub>	<b>57.7</b> <sub>(6.6)</sub>	56.9 <sub>(6.5)</sub>	59.4 <sub>(5.8)</sub>

Table 4: Ablation study results on the impact of GNN baselines in terms of mean classification accuracy (standard deviation is within brackets).

	IMDB-Binary				IMDB-Multi			
	2	3	5	10	2	3	5	10
GIN	57.9 <sub>(8.0)</sub>	54.5 <sub>(6.1)</sub>	54.3 <sub>(6.1)</sub>	56.7 <sub>(9.7)</sub>	32.6 <sub>(4.8)</sub>	31.6 <sub>(6.0)</sub>	32.8 <sub>(4.6)</sub>	36.0 <sub>(5.2)</sub>
GDM Acc	<b>58.8</b> <sub>(6.9)</sub>	57.0 <sub>(5.5)</sub>	57.8 <sub>(6.3)</sub>	<b>59.4</b> <sub>(6.3)</sub>	35.6 <sub>(4.2)</sub>	35.1 <sub>(4.6)</sub>	36.9 <sub>(3.9)</sub>	38.3 <sub>(3.5)</sub>
GDM Unc	58.2 <sub>(5.0)</sub>	<b>58.0</b> <sub>(6.5)</sub>	<b>60.5</b> <sub>(6.0)</sub>	57.7 <sub>(6.2)</sub>	<b>36.6</b> <sub>(4.2)</sub>	<b>37.2</b> <sub>(4.8)</sub>	<b>37.1</b> <sub>(3.9)</sub>	<b>39.5</b> <sub>(4.5)</sub>
GAT	51.2 <sub>(2.1)</sub>	50.6 <sub>(1.2)</sub>	55.1 <sub>(5.7)</sub>	54.4 <sub>(5.5)</sub>	31.8 <sub>(2.0)</sub>	34.0 <sub>(1.5)</sub>	32.6 <sub>(1.9)</sub>	33.6 <sub>(2.4)</sub>
GDM Acc	<b>55.7</b> <sub>(4.3)</sub>	56.2 <sub>(6.2)</sub>	55.0 <sub>(6.5)</sub>	<b>60.0</b> <sub>(7.1)</sub>	35.6 <sub>(2.9)</sub>	<b>36.4</b> <sub>(3.8)</sub>	<b>37.2</b> <sub>(4.3)</sub>	38.5 <sub>(5.0)</sub>
GDM Unc	54.5 <sub>(3.1)</sub>	<b>58.2</b> <sub>(6.7)</sub>	<b>56.0</b> <sub>(6.9)</sub>	59.2 <sub>(5.0)</sub>	<b>35.8</b> <sub>(2.8)</sub>	34.4 <sub>(3.6)</sub>	36.8 <sub>(5.4)</sub>	<b>38.8</b> <sub>(4.3)</sub>

rates on the Proteins and IMDB-Binary datasets are reported in Table 3. From the table, it is clear that GSAE outperforms VGSAE on both datasets across almost all label rates. The performance gain of GSAE decreases as the label rate increases, which highlights that VGSAE requires more training samples to obtain good performance. Therefore GSAE is more suitable for the case of learning with limited supervision as it is able to obtain good generalization performance with few samples due to its simple architecture and smaller number of learnable parameters. Nevertheless, the proposed GDM greatly and consistently outperforms the underlying GCN baseline across all different label rates with both GSAE and VGSAE on both datasets.

**Impact of GNN Baseline** We also conduct an ablation study to investigate the performance of our proposed GDM on additional GNN baselines. In particular, we applied GDM on the Graph Attention Network (GAT) [17] and Graph Isomorphism Network (GIN) [22] baselines. The comparison results with multiple label rates,  $\{2, 3, 5, 10\}$ , on the IMDB-Binary and IMDB-Multi datasets are reported in Table 4. The table clearly shows that GDM significantly improves the performance of both the GAT and GIN baselines across all label rates for both datasets. The performance gains are notable, exceeding 6%, 5% for GAT with label rates 5 and 3 for IMDB-Binary and IMDB-Multi, respectively. Similarly, GDM yields notable performance boost over GIN, exceeding 7%, 5% with label rates 3 and 5, respectively, for the IMDB-Binary dataset.

Table 5: Ablation study results on the impact of graph readout function in terms of mean classification accuracy (standard deviation is within brackets).

	Proteins				IMDB-Binary			
	2	3	5	10	2	3	5	10
Acc Mean	<b>65.3</b> <sub>(5.5)</sub>	<b>64.5</b> <sub>(4.8)</sub>	<b>65.8</b> <sub>(5.2)</sub>	<b>66.0</b> <sub>(5.3)</sub>	<b>56.1</b> <sub>(4.6)</sub>	<b>57.3</b> <sub>(6.8)</sub>	<b>59.1</b> <sub>(6.2)</sub>	<b>61.3</b> <sub>(6.7)</sub>
Acc Add	64.9 <sub>(5.1)</sub>	63.2 <sub>(4.9)</sub>	65.3 <sub>(4.8)</sub>	63.3 <sub>(4.3)</sub>	53.0 <sub>(2.8)</sub>	55.2 <sub>(5.3)</sub>	56.9 <sub>(6.1)</sub>	59.9 <sub>(6.7)</sub>
Acc Max	63.8 <sub>(4.7)</sub>	64.4 <sub>(4.8)</sub>	63.2 <sub>(7.0)</sub>	63.5 <sub>(7.4)</sub>	53.9 <sub>(4.7)</sub>	56.0 <sub>(7.8)</sub>	58.6 <sub>(7.0)</sub>	60.6 <sub>(6.3)</sub>
Unc Mean	<b>64.6</b> <sub>(4.0)</sub>	<b>63.6</b> <sub>(6.0)</sub>	<b>65.1</b> <sub>(5.3)</sub>	66.0 <sub>(5.0)</sub>	<b>58.0</b> <sub>(6.0)</sub>	<b>57.4</b> <sub>(7.1)</sub>	<b>58.4</b> <sub>(6.1)</sub>	<b>61.0</b> <sub>(7.0)</sub>
Unc Add	61.7 <sub>(6.5)</sub>	<b>63.6</b> <sub>(6.1)</sub>	63.0 <sub>(7.0)</sub>	63.4 <sub>(5.5)</sub>	55.6 <sub>(4.9)</sub>	54.8 <sub>(3.6)</sub>	57.4 <sub>(8.2)</sub>	60.5 <sub>(6.7)</sub>
Unc Max	63.0 <sub>(1.1)</sub>	63.4 <sub>(6.6)</sub>	61.9 <sub>(7.0)</sub>	<b>66.8</b> <sub>(5.9)</sub>	54.0 <sub>(4.2)</sub>	54.4 <sub>(7.6)</sub>	57.9 <sub>(7.1)</sub>	60.4 <sub>(7.1)</sub>

**Impact of Graph Readout Method** We conduct an ablation study to investigate the impact of the graph Readout function employed in our GNN model. Specifically, in addition to Global Mean Pooling, we consider the following two variants: (1) Add, where Global Add Pooling is used to obtain the graph-level embedding. (2) Max, where Global Max Pooling is used to obtain the graph-level embedding. The comparison results with different label rates on the Proteins and IMDB-Binary datasets are reported in Table 5. From the table, we can see that the Global Max Pooling and Global Add Pooling variants have performance drops compared to the Global Mean Pooling with almost all label rates for both the Proteins and IMDB-Binary datasets. Global Add Pooling suffers from obtaining un-normalized graph-level embeddings which causes generalization issues given that the graphs in each dataset have different sizes. Global Max Pooling only considers one feature per node corresponding to the feature with max value, causing the obtained graph-level embeddings to omit discriminative information present in the other features of the node-level embeddings.

## 5 Conclusion

In this paper, we proposed a novel Graph Dual Mixup (GDM) augmentation method for graph classification with limited labeled data. The proposed method employs a Graph Structural Auto-encoder to learn the structural embedding of the nodes, and then applies dual mixup on the structural node embeddings and the original node features of a pair of existing graphs in parallel to generate the structural and functional information of a new graph instance. The generated graph samples can augment the set of original graphs to alleviate overfitting and improve the generalizability of the GNN models. Additionally, we further propose two novel Balanced Graph Sampling methods to support GDM and enhance the balanced difficulty and diversity of the generated graph samples. We conducted experiments on six graph benchmark datasets, the experimental results demonstrate that the proposed method improves the generalization performance of the underlying GNNs when the labeled graphs are scarce and outperforms the state-of-the-art graph augmentation methods.

## References

1. Bianchi, F.M., Grattarola, D., Alippi, C.: Spectral clustering with graph neural networks for graph pooling. In: International Conference on Machine Learning (ICML) (2020)
2. Borgwardt, K.M., Ong, C.S., Schönauer, S., Vishwanathan, S., Smola, A.J., Kriegel, H.P.: Protein function prediction via graph kernels. *Bioinformatics* **21**(suppl\_1), i47–i56 (2005)
3. Dobson, P.D., Doig, A.J.: Distinguishing enzyme structures from non-enzymes without alignments. *Journal of molecular biology* **330**(4), 771–783 (2003)
4. Errica, F., Podda, M., Bacciu, D., Micheli, A.: A fair comparison of graph neural networks for graph classification. In: International Conference on Learning Representations (ICLR) (2020)
5. Guo, H., Sun, S.: Softedge: Regularizing graph classification with random soft edges. arXiv preprint arXiv:2204.10390 (2022)
6. Hamilton, W., Ying, Z., Leskovec, J.: Inductive representation learning on large graphs. In: Advances in Neural Information Processing Systems (NIPS) (2017)
7. Han, X., Jiang, Z., Liu, N., Hu, X.: G-mixup: Graph data augmentation for graph classification. In: International Conference on Machine Learning (ICML) (2022)
8. Haussler, D., et al.: Convolution kernels on discrete structures. Tech. rep., Citeseer (1999)
9. Huang, W., Zhang, T., Rong, Y., Huang, J.: Adaptive sampling towards fast graph representation learning. In: Advances in Neural Information Processing Systems (NeurIPS) (2018)
10. Kipf, T.N., Welling, M.: Variational graph auto-encoders. arXiv preprint arXiv:1611.07308 (2016)
11. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. In: International Conference on Learning Representations (ICLR) (2017)
12. Morris, C., Kriege, N.M., Bause, F., Kersting, K., Mutzel, P., Neumann, M.: TUDataset: A collection of benchmark datasets for learning with graphs. In: ICML Workshop on Graph Representation Learning and Beyond (2020)
13. Park, J., Shim, H., Yang, E.: Graph transplant: Node saliency-guided graph mixup with local structure preservation. In: AAAI Conference on Artificial Intelligence (2022)
14. Rong, Y., Huang, W., Xu, T., Huang, J.: Dropedge: Towards deep graph convolutional networks on node classification. In: International Conference on Learning Representations (ICLR) (2020)
15. Shervashidze, N., Schweitzer, P., Van Leeuwen, E.J., Mehlhorn, K., Borgwardt, K.M.: Weisfeiler-lehman graph kernels. *Journal of Machine Learning Research* **12**(9) (2011)
16. Sun, L., Xia, C., Yin, W., Liang, T., Yu, P.S., He, L.: Mixup-transformer: dynamic data augmentation for nlp tasks. arXiv preprint arXiv:2010.02394 (2020)
17. Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., Bengio, Y.: Graph attention networks. In: International Conference on Learning Representations (ICLR) (2018)
18. Verma, V., Qu, M., Kawaguchi, K., Lamb, A., Bengio, Y., Kannala, J., Tang, J.: Graphmix: Improved training of gnn for semi-supervised learning. In: AAAI Conference on Artificial Intelligence (2021)
19. Wale, N., Watson, I.A., Karypis, G.: Comparison of descriptor spaces for chemical compound retrieval and classification. *Knowledge and Information Systems* **14**, 347–375 (2008)

20. Wang, Y., Wang, W., Liang, Y., Cai, Y., Hooi, B.: Graphcrop: Subgraph cropping for graph classification. arXiv preprint arXiv:2009.10564 (2020)
21. Wang, Y., Wang, W., Liang, Y., Cai, Y., Hooi, B.: Mixup for node and graph classification. In: International World Wide Web Conference (WWW) (2021)
22. Xu, K., Hu, W., Leskovec, J., Jegelka, S.: How powerful are graph neural networks? In: International Conference on Learning Representations (ICLR) (2019)
23. Yanardag, P., Vishwanathan, S.: Deep graph kernels. In: ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD) (2015)
24. Ying, Z., You, J., Morris, C., Ren, X., Hamilton, W., Leskovec, J.: Hierarchical graph representation learning with differentiable pooling. In: Advances in Neural Information Processing Systems (NeurIPS) (2018)
25. You, Y., Chen, T., Sui, Y., Chen, T., Wang, Z., Shen, Y.: Graph contrastive learning with augmentations. In: Advances in Neural Information Processing Systems (NeurIPS) (2020)
26. Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. In: International Conference on Learning Representations (ICLR) (2018)
27. Zhao, T., Liu, G., Günnemann, S., Jiang, M.: Graph data augmentation for graph machine learning: A survey. arXiv preprint arXiv:2202.08871 (2022)