

Graph Domain Adaptation: A Generative View

Ruichu Cai*, *Member, IEEE*, Fengzhu Wu, Zijian Li, Pengfei Wei, Lingling Yi, Kun Zhang

Abstract—Recent years have witnessed tremendous interest in deep learning on graph-structured data. Due to the high cost of collecting labeled graph-structured data, domain adaptation is important to supervised graph learning tasks with limited samples. However, current graph domain adaptation methods are generally adopted from traditional domain adaptation tasks, and the properties of graph-structured data are not well utilized. For example, the observed social networks on different platforms are controlled not only by the different crowd or communities but also by the domain-specific policies and the background noise. Based on these properties in graph-structured data, we first assume that the graph-structured data generation process is controlled by three independent types of latent variables, i.e., the semantic latent variables, the domain latent variables, and the random latent variables. Based on this assumption, we propose a disentanglement-based unsupervised domain adaptation method for the graph-structured data, which applies variational graph auto-encoders to recover these latent variables and disentangles them via three supervised learning modules. Extensive experimental results on two real-world datasets in the graph classification task reveal that our method not only significantly outperforms the traditional domain adaptation methods and the disentangled-based domain adaptation methods but also outperforms the state-of-the-art graph domain adaptation algorithms.

Index Terms—Graph Neural Network, Graph Generative Models, Domain Adaptation

I. INTRODUCTION

Though deep learning on graph-structured data has achieved great success, similar to other deep learning methods, it heavily depends on labeled data. However, the high cost of collecting and labeling graph-structured data in real-world applications is unacceptable. Fortunately, unsupervised domain adaptation can figure out the aforementioned problem. Therefore, unsupervised domain adaptation on graph-structured data is important to various graph learning tasks. One example

Ruichu Cai is with the School of Computing, Guangdong University of Technology and Guangdong Provincial Key Laboratory of Public Finance and Taxation with Big Data Application, Guangzhou China, 510006. E-mail: cairuichu@gmail.com

Fengzhu Wu is with the School of Computing, Guangdong University of Technology, Guangzhou China, 510006. E-mail: fzwu97@gmail.com

Zijian Li is with the School of Computing, Guangdong University of Technology, Guangzhou China, 510006. E-mail: leizigin@gmail.com

Pengfei Wei is with National University of Singapore. E-mail: wpf89928@gmail.com

Lingling Yi is with Tencent Technology (SZ) Co., Ltd. E-mail: chrisyi@tencent.com

Kun Zhang is with the Department of Philosophy, Carnegie Mellon University, Pittsburgh, PA 15213 USA. E-mail: kunz1@cmu.edu

Manuscript received XX; revised XX; accepted XX. Date of publication XX XX, 2019; date of current version XX XX, 2019. Ruichu Cai and Zijian Li was supported in part by Natural Science Foundation of China (61876043, 61976052), Science and Technology Planning Project of Guangzhou (201902010058) and Guangdong Provincial Science and Technology Innovation Strategy Fund (2019B121203012). (Corresponding author: Ruichu Cai.)

is the influence prediction on social networks with different communities [1], [2]¹, which aims to predict whether a given user will retweet.

Numerous unsupervised domain adaptation works have been designed for images data [3], texts data [4], [5] and time series data [6], while only a few methods have been proposed for graph data. For instance, [7] achieves domain-adaptive network embedding with the help of the shared graph convolutional network and the Generative Adversarial Networks (GAN) based adversarial regularization [8]. [9] applies the attention mechanism to integrate global and local consistency and extract cross-domain node embedding by deceiving the domain discriminator with the help of Gradient Reversal Layer (GRL) [10]. As summary, these adversarial methods are generally adopted from traditional domain adaptation tasks, and the properties of graph-structured data are not explicitly utilized.

Different from traditional data, the graph data is featured with the high uncertainty of the generation process and the high complexity of the domain-specific structures. We take the social network in Figure 1 (d) as an example. The observed social networks are not only controlled by the latent semantic factors (e.g., the shared interests among the nodes), but also by the domain-specific communities (e.g., raised by the domain-specified network policies) and the uncertainty (e.g., two nodes with the same interests may not know each other or the zombie fans in the social networks). The ignorance of these important properties for graph-structured data hinders the performance of the graph domain adaptation model.

In order to learn a graph representation that is robust across domains, two kinds of important information should be included. The first is the domain-invariant structure/topological information and the second is the domain-invariant node information. However, it is not a trivial task to extract this domain-invariant information at the same time. The main difficulty of extracting the domain-invariant structure information from the observed graph-structured data is the uncertainty. Due to the high uncertainty of observed graph data, even for the users surrounded by the same communities, the observed graph structures between different platforms/domains are totally different, e.g., the cases are given in Figure 1 (a) and (c). And the main obstacle of extracting the domain-invariant node information the different communities from different platforms/domains, which is usually led by the domain-specified network policies. Give a detailed example shown in Figure 1 (a) and (b), users with a similar observed graph structure, can be surrounded by different social communities. In this case, since the distribution and the types of communities vary across domains, it is hard to extract the domain-invariant node representation solely using

¹A community is a group of users with the same interest.

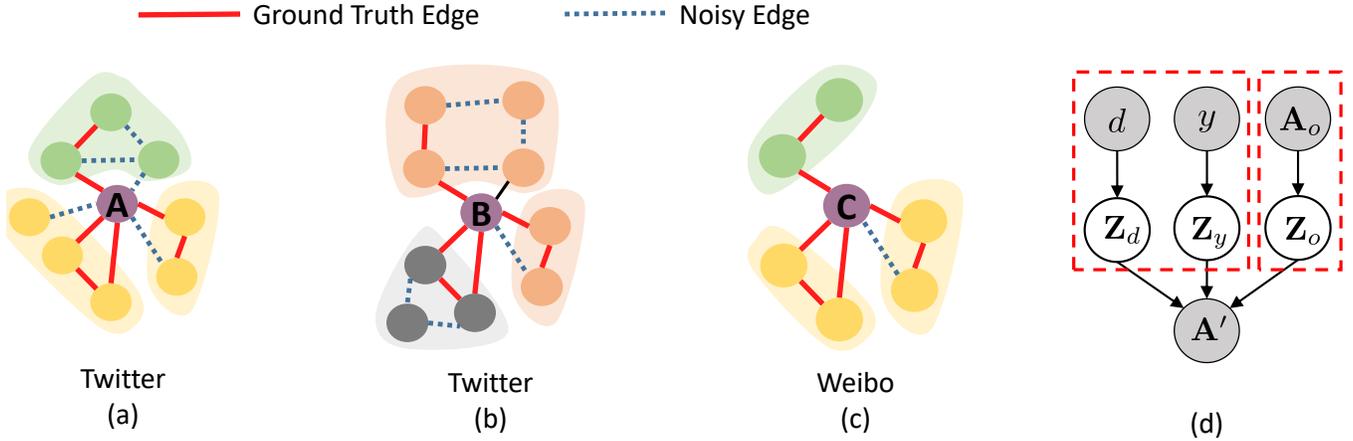


Fig. 1. (a)-(c) Three observed graph structures from Twitter and Weibo. Each node denotes a user. The purple nodes denote the users who will retweet since they share the same structures. The irregular blocks in different colors denote different communities. Communities from different domains are not identical. The blue dashed lines and the red lines denote the noisy edges and the true edges, respectively. (d) The causal model of the graph-structured data generation process, which is controlled by three independent types of latent variables Z_d , Z_y and Z_o . Here, d , y , and A_o denote the domain information, semantic information, and the uncertainty of graph data. (*best view in color*)

the conventional domain adaptation restrictions like gradient reversal layer and MMD.

Therefore, we propose a Disentanglement-based Graph Domain Adaptation Model (**DGDA** in short), which is motivated by the graph data generation process as shown in Figure 1 (d). In detail, we assume that the graph data generation process is controlled by three independent types of latent variables, i.e., the semantic latent variables Z_y , the domain latent variables Z_d , and the random latent variables Z_o . These three groups of latent variables are encoded by the semantic information y , the domain information d , and the uncertainty A_o , respectively. By reconstructing and disentangling these three latent variables, we can easily classify the label of different graphs based on Z_y . Technically, we employ a variational graph auto-encoders to reconstruct these latent variables. We disentangle the semantic and domain latent variables by employing a label classifier and a domain classifier. We further disentangle the random latent variables by reconstructing the uncertainty of the graph data. Furthermore, we apply the latent variables regularization on these three types of latent variables for better disentanglement. The extensive experimental studies demonstrate that our method outperforms state-of-the-art unsupervised domain adaptation methods on graph data.

The rest of the article is organized as follows. Section II reviews existing studies on graph representation learning, domain adaptation, and graph domain adaptation. Section III provides the problem definition on graph domain adaptation and the disentanglement-based graph domain adaptation model. Section IV presents the experiment results and analysis on two real-world datasets. Section V concludes the article with future work discussion.

II. RELATED WORK

Our work is closely related to graph representation learning and domain adaptation, so we first review the existing techniques on graph representation learning and domain adapta-

tion, and then we give an introduction about domain adaptation on graph-structured data.

A. Graph Representation Learning

Representation learning [11] has been one of the hotspots of deep learning research fields. Graph representation learning aims to map edges and nodes of a graph into a low-dimensional vector space with original graph structures and properties being well preserved. Current methods can be categorized as inductive methods and transductive methods. Inductive methods [12]–[15] learn functions that take the graph topology structure and edge or node features as input and output their representation vectors. Transductive methods [16]–[20] optimize the representation vectors directly. Though these methods do not provide a representation of the entire graph, these methods are enough for our user-oriented prediction tasks.

Recently, there have been several attempts to learn latent representations of sub-structures for graphs via kernel technique [21] or pooling operation. Global pooling methods aggregate the node representations either via simple readout functions such as averaging the node embeddings [22] or more complex set operations [23], [24]. On another line, hierarchical pooling methods [25]–[31] coarsen the node representations over the network’s layers [32], and finally achieve the entire graph representation. However, these methods mainly focus on learning graph representation from a single domain and ignores the transferable sub-structures in the graph.

B. Domain Adaptation

Domain adaptation studies how to learn a model that is transferable on different but related domains. Previous studies mainly focus on learning domain invariant representation via neural networks, so that deep learning models trained on a labeled source domain can be transferred to a target domain with few or no labeled samples. They can be classified into

statistic-based methods and adversarial methods. Statistic-based methods [6], [33], [34] utilize maximum mean discrepancy (MMD) to achieve the domain alignment. Inspired by GAN [8], adversarial domain adaptation methods [5], [10], [35]–[38] employ a domain adversarial layer to minimize the domain discrepancy, where a feature extractor and a domain classifier compete against each other. These methods mainly focus on grid data like images. In this paper, we focus on graph-structured data, which is more complicated and challenging than grid data.

C. Graph Domain Adaptation

Recently, few methods have been proposed to address the graph domain adaption tasks. [7] achieves domain adaptive network embedding via shared weight graph convolutional network and adversarial learning regularization. [9] applies an attention mechanism to integrate global and local consistency and extract cross-domain node embedding with the help of Gradient Reversal Layer (GRL) [10]. [39] utilizes two feature extractors to preserve attributed affinity and topological proximity between nodes.

However, these works mainly focus on the node classification, by simply replacing the feature extractor with the graph convolutional network, without taking the property of the graph-structured data into account. Recently, Elif et al [40], [41] handle graph domain adaptation via learning aligned graph bases. In this paper, we not only focus on the challenging graph classification task but also well utilize the property of graph-structured data via the generation process of graph-structured data.

Please note that our work is related to but different from DIVA [42]. The similarity is that both our method and DIVA aim to reconstruct three types of latent variables. The difference is that DIVA only reconstructs the residual variations in an unsupervised way, while our method obtains the random latent variables by explicitly reconstructing three latent variables explicitly with the help of the supervised information as well as the synthetic random noise.

III. DISENTANGLEMENT BASED GRAPH DOMAIN ADAPTATION MODEL

In this paper, we focus on the unsupervised domain adaptation problem on graph-structured data, which aims to use the labeled samples $G_S = \{g_l^S, y_l^S\}_{l=1}^{n_S}$ on source domain to classify the unlabeled samples $G_T = \{g_l^T\}_{l=1}^{n_T}$ on target domain, where g and y denote the graph sample and its label. Herein, n_S and n_T denote the number of samples in source and target domain, respectively. Each graph sample g with N nodes contains an adjacent matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ and a node feature matrix $\mathbf{X} \in \mathbb{R}^{N \times K_x}$, i.e., $g = \{\mathbf{A}, \mathbf{X}\}$. The goal of this paper is to model the generation process of graph-structured data and design an effective graph domain adaptation framework. The mathematical notations used in this paper are summarized in Table I.

The main motivation of our model is to address the high uncertainty of graph data and extract the semantic-related

TABLE I
NOTATION TABLE.

Symbols	Definitions and Descriptions
G_S, n_S	The source domain dataset and its size.
G_T, n_T	The target domain dataset and its size.
$g = \{\mathbf{A}, \mathbf{X}\}$	A graph sample in G_S or G_T , its adjacency matrix, and its node features.
d	The domain label of g .
y	The graph label of g .
N	The number of nodes in g .
K_x	The dimension of node features.
\mathbf{Z}_d, E_d	The semantic latent variables and its encoder.
\mathbf{Z}_y, E_y	The domain latent variables and its encoder.
\mathbf{Z}_o, E_o	The random latent variables and its encoder.
F	The node feature extractor.
\mathbf{H}	The node representation outputs by F .
μ	The mean of latent variables.
σ	The standard deviation of latent variables.
C_d	The domain disentanglement module.
C_y	The label disentanglement module.
D_o	The noise reconstruction module.
D_g	The graph reconstruction module.
\mathbf{A}'	The augmented adjacency matrix of g .
\mathbf{A}_o	The randomly generated noise matrix.
\mathbf{D}	The diagonal degree matrix of \mathbf{A}' .
$\hat{\mathbf{A}}$	The symmetric normalized \mathbf{A}' .
\mathbf{W}, \mathbf{b}	Weights and biases in neural networks.
$\phi_f, \phi_d, \phi_y, \phi_o$	The parameters in encoding block.
$\theta_d, \theta_y, \theta_o$	The parameters in disentanglement block.
θ_g	The parameters of graph decoder D_g .
\oplus	The concatenation operator of any two vectors.

information behind the graphs. We start from the graph-structured data generation process as shown in Figure 1 (d). Given a graph sample, it assumed to be generated from three independent types of latent variables, i.e., \mathbf{Z}_d encoding the domain information, \mathbf{Z}_y encoding the semantic information, and \mathbf{Z}_o encoding the uncertainty information. And d , y , and \mathbf{A}_o respectively denote the domain information, the semantic information, and the uncertainty of graph data. We use $\mathbf{Z}_d \in \mathbb{R}^{N \times K_d}$, $\mathbf{Z}_y \in \mathbb{R}^{N \times K_y}$ and $\mathbf{Z}_o \in \mathbb{R}^{N \times K_o}$ to denote domain latent variables, semantic latent variables and random latent variables. Considering that the domain information varies with domains, we induce that the semantic latent variables play an important role in extracting the domain-invariant representation. We further assume that these latent variables are independent.

Based on the aforementioned graph-structured data generative mechanism, we design a graph domain adaptation framework. The model architecture is shown in Figure 2. We first extract features via a feature extractor $F(\cdot; \phi_f)$, and then use a VAE-like architecture to reconstruct three independent latent variables \mathbf{Z}_d , \mathbf{Z}_y and \mathbf{Z}_o . However, unlike the vanilla VAE, we further design a disentanglement block as shown in the red dashed box in Figure 2. In this architecture, three learning modules are placed under \mathbf{Z}_d , \mathbf{Z}_y , and \mathbf{Z}_o respectively. For the domain disentanglement module $C_d(\mathbf{Z}_d; \theta_d)$, it aims to extract all the domain information into \mathbf{Z}_d and exclude other information. For the label disentanglement module $C_y(\mathbf{Z}_y; \theta_y)$, it aims to extract all the semantic information into \mathbf{Z}_y and exclude other information. For the noise variables reconstruction module $D_o(\mathbf{Z}_d; \theta_o)$, it aims to model the uncertainty of graph-structured data from \mathbf{Z}_o and exclude other information. As a

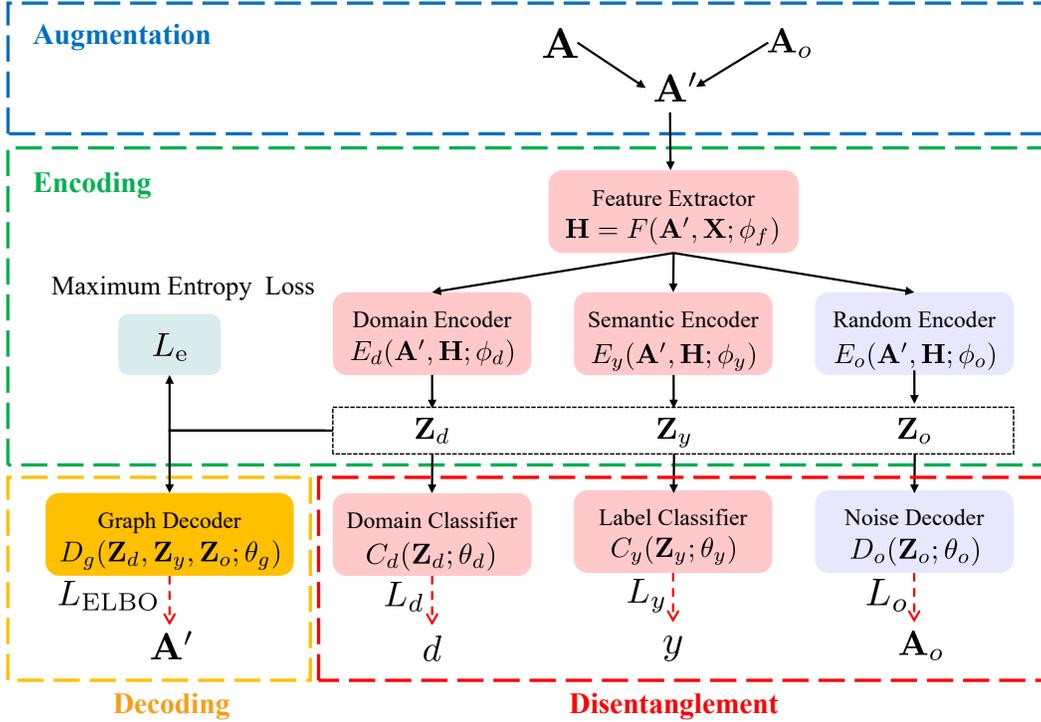


Fig. 2. The framework of DGDA. First, in the augmentation block, augmented samples are generated. Second, in the encoding block, a VGAE is applied to recover \mathbf{Z}_d , \mathbf{Z}_y , \mathbf{Z}_o . Third, in the decoding block, a graph decoder D_g integrates information from all three latent variables to reconstruct \mathbf{A}' . Final, in the disentanglement block, three types of reconstruction are used to disentangle the latent variables. C_y and C_d are the classifiers for the label y and the domain d , respectively. D_o is used to reconstruct \mathbf{A}_o . L_e is the latent variables regularization. (best view in color)

result, we obtain domain-invariant label information, which is also robust to uncertainty. We will give a detailed description of our model in the following section.

A. Uncertainty Modeling

Real-world data is usually noisy, which leads to the degeneration of model performance. Data augmentation can mitigate this problem to a certain extent, but the unseen manipulations still can not be handle because it is impossible to see all different manipulations at the training stage. Motivated by [43], we apply a random manipulation-based data augmentation on graph-structured data. We explicitly model the distribution of noises via random latent variables \mathbf{Z}_o , which controls the uncertainty of the graph-structured data. In order to generate the augmented graph-structured data, we perform a valid augmentation of the original data, which means we further assume that this augmentation only controls the structure of \mathbf{A} and is independent of the graph label y and the domain label d . There are a variety of methods to generate augmented graph data. In this paper, we adopt a strategy similar to DropEdge [44]. For each graph sample, we first generate a noise matrix \mathbf{A}_o by randomly dropping and adding edges in the original adjacency matrix \mathbf{A} . The details of $a_{ij} \in \mathbf{A}_o$ as shown in Equation(1).

$$a_{ij} = \begin{cases} 0 & \text{Constant between node } i \text{ and node } j. \\ 1 & \text{Add a edge between node } i \text{ and node } j. \\ -1 & \text{Drop a edge between node } i \text{ and node } j. \end{cases} \quad (1)$$

Specifically, the data augmentation process is formulated as follows:

$$\begin{aligned} \mathbf{A}' &= \mathbf{M}_{\text{drop}} \odot \mathbf{A} + \mathbf{M}_{\text{add}}, \\ m_{ij}^{\text{drop}} &\sim \text{Bernoulli}(p_{\text{drop}}), \\ m_{ij}^{\text{add}} &\sim \text{Bernoulli}(p_{\text{edge}} \cdot p_{\text{add}}), \\ \mathbf{A}_o &= \mathbf{A}' - \mathbf{A}, \end{aligned} \quad (2)$$

where p_{drop} , p_{add} , and p_{edge} refer to the dropping edge rate, adding edge rate, and the sparsity of \mathbf{A} , respectively. \mathbf{M}_{drop} and \mathbf{M}_{add} denote the mask matrix with the same shape as \mathbf{A} . m_{ij}^{drop} and m_{ij}^{add} denote the elements in two mask matrix, and they are sampled from the Bernoulli distribution independently.

Based on the aforementioned graph-structured data generative mechanism, we design a graph domain adaptation framework by first recovering the three types of latent variables via a VAE, and then disentangling them.

B. Latent Variables Reconstruction

The encoding and decoding modules in our framework are inspired by Variational Graph Auto-Encoders (VGAE) [12], a framework for unsupervised learning on graph-structured data based on the VAE. In VGAE, $q_{\phi}(\mathbf{Z}|\mathbf{A}, \mathbf{X})$ denotes the encoder with respect to the parameter ϕ to approximate the intractable true posterior distribution $p(\mathbf{Z}|\mathbf{A}, \mathbf{X})$, while $P_{\theta_{\phi}}(\mathbf{A}|\mathbf{Z})$ denotes the decoder with respect to the parameters

θ_r . $P(\mathbf{Z})$ is the prior distribution. The variational lower bound of the marginal likelihood is given as follow:

$$\mathcal{L}_{\text{ELBO}}(\phi, \theta_r) = \mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{A}, \mathbf{X})}[P_{\theta_r}(\mathbf{A}|\mathbf{Z})] - D_{KL}(q_\phi(\mathbf{Z}|\mathbf{A}, \mathbf{X})\|P(\mathbf{Z})), \quad (3)$$

where $D_{KL}(q(\cdot)\|P(\cdot))$ is the Kullback-Leibler divergence between distribution $q(\cdot)$ and $P(\cdot)$.

Unlike vanilla VGAE, we further decompose the latent variables \mathbf{Z} into \mathbf{Z}_d , \mathbf{Z}_y and \mathbf{Z}_o . As a result, Equation (3) can be derived as:

$$\begin{aligned} \mathcal{L}_{\text{ELBO}}(\phi_{f,d,y,o}, \theta_g) = & \mathbb{E}_{q_{\phi_{f,d,y,o}}(\mathbf{Z}_d, \mathbf{Z}_y, \mathbf{Z}_o|\mathbf{H})} [\log P_{\theta_g}(\mathbf{A}'|\mathbf{Z}_d, \mathbf{Z}_y, \mathbf{Z}_o)] \\ & - \sum_{k \in \{d,y,o\}} D_{KL}(q_{\phi_k}(\mathbf{Z}_k|\mathbf{A}', \mathbf{X})\|P(\mathbf{Z}_k)), \end{aligned} \quad (4)$$

, where $\mathbf{H} = F(\mathbf{A}', \mathbf{X}; \phi_f)$. F is the feature extractor that takes the node features and the adjacency matrix as input, and outputs the node representation $\mathbf{H} \in \mathbb{R}^{N \times K_H}$.

We further assume that the prior distributions on all the latent variables are normal distributions, i.e., $P(\mathbf{Z}_d), P(\mathbf{Z}_y), P(\mathbf{Z}_o) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Similar to VGAE, by applying a reparameterization trick, we use three Graph Convolutional Networks (GCNs) [19] based encoders E_d, E_y and E_o as the approximator of q , to encode the node representation \mathbf{H} into $\mathbf{Z}_d, \mathbf{Z}_y$ and \mathbf{Z}_o , respectively. Finally, we aggregate all the latent variables and employ an inner product decoder that maps each pair of node representations to a binary indicator of edge existence in \mathbf{A}' .

1) *Inference Model.*: We first extract the node representation via the GCNs. The GCNs comprises multiple stacked graph convolutional layers to extract features from multi-order neighbors. We take the node features \mathbf{X} and the augmented adjacency matrix \mathbf{A}' as the input of the GCNs, where each GCN layer can be formulated as follow:

$$\begin{aligned} \text{GCN}(\mathbf{A}', \mathbf{X}) &= \text{ReLU}(\hat{\mathbf{A}}\mathbf{X}\mathbf{W}), \\ \hat{\mathbf{A}} &= \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}, \end{aligned} \quad (5)$$

where $\tilde{\mathbf{A}} = \mathbf{A}' + \mathbf{I}_m$ is the augmented adjacency matrix with self-loop, $\tilde{\mathbf{D}}_{ii} = \text{diag}(\sum_j \tilde{\mathbf{A}}_{ij})$ is the diagonal degree matrix, \mathbf{W} is the trainable parameters and ReLU is the Rectified Linear Unit (ReLU) activation function. Specifically, in DGDA, we use a two-layer GCN as the feature extractor. For convenience, we let $\mathbf{H} = F(\mathbf{A}', \mathbf{X}; \phi_f)$ be the aforementioned process, ϕ_f are the trainable parameters.

In order to obtain three types of latent variables which are shown in Figure 1 (d), we further devise three encoders. These encoders are used to encode μ and $\log \sigma$ of $\mathbf{Z}_d, \mathbf{Z}_y$ and \mathbf{Z}_o , respectively. For the type $k \in \{d, y, o\}$ latent variables, the encoding model is formulated as follows:

$$\begin{aligned} q(\mathbf{Z}_k|\mathbf{A}', \mathbf{H}) &= \prod_{i=1}^N q(z_{ki}|\mathbf{A}', \mathbf{H}), \\ q(z_{ki}|\mathbf{A}', \mathbf{H}) &= \mathcal{N}(z_{ki}|\mu_{ki}, \text{diag}(\sigma_{ki}^2)), \\ \mu_k &= \text{GCN}_{\mu}(\mathbf{A}', \mathbf{H}) = \hat{\mathbf{A}}\mathbf{H}\mathbf{W}_{\mu k}, \\ \log \sigma_k &= \text{GCN}_{\sigma}(\mathbf{A}', \mathbf{H}) = \hat{\mathbf{A}}\mathbf{H}\mathbf{W}_{\sigma k}. \end{aligned} \quad (6)$$

Here, z_{ki} is the i th row of \mathbf{Z}_k , and the same for μ_{ki} and $\log \sigma_{ki}$. Let $\phi_k = \{\mathbf{W}_{\mu k}, \mathbf{W}_{\sigma k}\}$, so the encoder can be formulated as $E_k(\mathbf{A}', \mathbf{H}; \phi_k)$. By applying a reparameterization trick, we use three GCN based encoders $E_d(\mathbf{A}', \mathbf{H}; \phi_d)$, $E_y(\mathbf{A}', \mathbf{H}; \phi_y)$ and $E_o(\mathbf{A}', \mathbf{H}; \phi_o)$ to encode \mathbf{H} into three latent variables $\mathbf{Z}_d, \mathbf{Z}_y$ and \mathbf{Z}_o , respectively.

We also want to claim that our method is not restricted to the specific graph classification task. By replacing GCN with other graph neural networks like GAT [20], our method can be easily extended to different graph learning tasks.

2) *Generative Model.*: The generative model in VGAE is given by an inner product among latent variables. Since there are three types of latent variables in our graph decoder, we first employ an MLP layer to aggregate information from all latent variables and improve the expressiveness. Then we apply an inner product between each pair of node representations, which can be formulated as follows:

$$\begin{aligned} \mathbf{Z}_g &= \text{ReLU}([\mathbf{Z}_d \oplus \mathbf{Z}_y \oplus \mathbf{Z}_o]\mathbf{W}_{g0})\mathbf{W}_{g1}, \\ p(\mathbf{A}'|\mathbf{Z}_g) &= \prod_{i=1}^N \prod_{j=1}^N p(a'_{ij}|z_{gi}, z_{gj}) \end{aligned} \quad (7)$$

$$\text{with } p(a'_{ij} = 1|z_{gi}, z_{gj}) = \sigma(z_{gi}^\top z_{gj}),$$

where z_{gi} is the i th row of \mathbf{Z}_g . Please note that \oplus is the concatenate operation, a'_{ij} is the element of \mathbf{A}' and $\sigma(\cdot)$ is the logistic sigmoid function. Let $\theta_g = \{\mathbf{W}_{g0}, \mathbf{W}_{g1}\}$, so the decoder module can be formulated as $D_g(\mathbf{Z}_d, \mathbf{Z}_y, \mathbf{Z}_o; \theta_g)$.

C. Latent Variables Disentanglement

1) *Domain Variables Reconstruction.*: The disentanglement architecture of our framework is shown in the red dashed box in Figure 2, and it consists of three modules working together. The domain disentanglement module extracts the domain information by training a domain classifier C_d . The parameters θ_d are learned by minimizing the binary cross-entropy loss L_d . So the objective function of the domain disentanglement module is shown as follow:

$$\begin{aligned} \hat{d}_i &= C_d(\mathbf{Z}_d; \theta_d) \\ L(\hat{d}_i, d_i) &= d_i \log(\hat{d}_i) + (1 - d_i) \log(1 - \hat{d}_i) \\ \mathcal{L}_d(\theta_d) &= \frac{1}{n_S + n_T} \sum_{g_i \in (G_S, G_T)} L(\hat{d}_i, d_i), \end{aligned} \quad (8)$$

where d_i is the domain label.

2) *Semantic Variables Reconstruction.*: The $C_y(\mathbf{z}_y; \theta_y)$ is the label learning module that extracts the semantic information. This is done by training a label classifier C_y on the source domain since labeled data is unavailable in the target domain. As a result, the parameters θ_y in C_y are learned by minimizing the categorical cross-entropy loss \mathcal{L}_y . The objective function of the label disentanglement module is shown as follows:

$$\begin{aligned} \hat{y}_i &= C_y(\mathbf{Z}_y; \theta_y) \\ L(\hat{y}_i, y_i) &= y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \\ \mathcal{L}_y(\theta_y) &= \frac{1}{n_S} \sum_{(g_i, y_i) \in G_S} L(\hat{y}_i, y_i), \end{aligned} \quad (9)$$

where y_i is the graph label. Note that we are doing a graph classification task, the read out function is average over the output node features.

3) *Noise Variables Reconstruction.*: In this part, we give a description of noise variables reconstruction model $D_o(\mathbf{Z}_o; \theta_o)$. In order to disentangle the random latent variables \mathbf{Z}_o , we reconstruct the noise information introduced by data augmentation, i.e., the noise matrix \mathbf{A}_o , which means that we wish our model to infer \mathbf{A}_o from the augmented data \mathbf{A}' . The architecture of noise variables reconstructions module D_o is similar to the generative model D_g , which can be formulated as follow:

$$\begin{aligned} \mathbf{R} &= \text{ReLU}(\mathbf{Z}_o \mathbf{W}_{n0}) \mathbf{W}_{n1}, \\ p(\mathbf{A}_o | \mathbf{R}) &= \prod_{i=1}^N \prod_{j=1}^N p(a_{ij} | \mathbf{r}_i, \mathbf{r}_j) \quad (10) \\ \text{with } p(a_{ij} = 1 | \mathbf{r}_i, \mathbf{r}_j) &= \sigma(\mathbf{r}_i^\top \mathbf{r}_j), \end{aligned}$$

where \mathbf{r}_i is the i th row of \mathbf{R} . Finally, the objective function can be formulated as follows:

$$\mathcal{L}_o(\theta_o) = \mathbf{E}_{q_{\phi_o}(\mathbf{Z}_o | \mathbf{A}', \mathbf{H})} [\log P_{\theta_o}(\mathbf{A}_o | \mathbf{Z}_o)]. \quad (11)$$

4) *Latent Variables Regularization.*: In order to exclude the redundant information from the latent variables, we employ an element-wise maximum entropy loss \mathcal{L}_e restriction on each latent variables. The above three disentanglement terms ‘‘pull’’ the relevant information into each latent variables, while the regularizers ‘‘push’’ the irrelevant information away from the corresponding latent variables. Such regularization ensures the purity of the reconstructed latent variables. The details of the regularizers are shown as follows:

$$\begin{aligned} \mathcal{L}_e(\phi_{f,d,y,o}) &= \frac{1}{n_S + n_T} \sum_{g_i \in (G_S, G_T)} \sum_{k \in (d,y,o)} L_e(\mathbf{Z}_k), \\ L_e(\mathbf{Z}_k) &= \frac{1}{N \times K_k} \sum_{i=0}^N \sum_{j=0}^{K_k} \sigma(z_{ij}) \log \sigma(z_{ij}), \end{aligned} \quad (12)$$

where K_k is the dimension of type k latent variable.

D. Model Summary

We summarize our model as follows. The total loss of our proposed framework can be formulated as:

$$\mathcal{L}(\phi_{f,d,y,o}, \theta_{g,d,y,o}) = \mathcal{L}_{\text{ELBO}} + \gamma \mathcal{L}_d + \alpha \mathcal{L}_y + \omega \mathcal{L}_o + \delta \mathcal{L}_e, \quad (13)$$

where γ, α, ω , and δ are the hyperparameters that control the weight of losses. In our experiment, we let $\gamma = 1$, $\alpha = 1$, $\omega = 0.1$ and $\delta = 5$. Small weight for \mathcal{L}_o can effectively constrain the negative effect brought by dropping and adding edges. And a larger δ reduces more information from the latent variables. The training of our model can be divided into two steps. In each epoch, we first fix the noise variables reconstruction module (blue blocks in Figure 2) i.e., E_o and D_o , and jointly train the other modules (red blocks in Figure 2) with clean data on source and target domains. The training

algorithm is shown in Algorithm 1. The optimization process is shown as follow:

$$\left(\hat{\phi}_f, \hat{\phi}_d, \hat{\phi}_y, \hat{\theta}_g, \hat{\theta}_d, \hat{\theta}_y \right) = \arg \min_{\phi_{f,d,y}, \theta_{g,d,y}} \mathcal{L}(\phi_{f,d,y,o}, \theta_{g,d,y,o}). \quad (14)$$

In the second step, to prevent model collapse causing by a bad noise matrix, we train the noise variables reconstruction module independently with the augmented data, while the rest of the parameters remain fixed. The optimization process is shown as follow:

$$\left(\hat{\phi}_f, \hat{\phi}_o, \hat{\theta}_o \right) = \arg \min_{\phi_f, \phi_o, \theta_o} \mathcal{L}(\phi_{f,d,y,o}, \theta_{g,d,y,o}). \quad (15)$$

All parameters are optimized using the stochastic gradient descent algorithm. In the test phase, we predict the label as follow:

$$\hat{y} = C_y \left(E_y \left(F \left(\mathbf{A}, \mathbf{X}; \hat{\phi}_f \right); \hat{\phi}_y \right); \hat{\theta}_y \right) \quad (16)$$

Please note that we input the clean sample without the data augmentation process in the test phase.

Finally, we summarize the pseudo-code of training DGDA as follows:

Algorithm 1 DGDA training algorithm

- 1: **repeat**
 - 2: Randomly select g_s, y_s from G_S ;
 - 3: Randomly select g_t, y_t from G_T ;
 - 4: Forward propagation;
 - 5: Update $\phi_f, \phi_d, \phi_y, \theta_g, \theta_d, \theta_y$ based on Equation (14);
 - 6: Augment g_s, g_t based on Equation (2)
 - 7: Forward propagation;
 - 8: Update ϕ_f, ϕ_o, θ_o based on Equation (16);
 - 9: **until** Max Iteration or Early Stopping
-

TABLE II
STATISTICS OF IMDB&REDDIT.

Dataset	IMDB-Binary	Reddit-Binary
# of Nodes	19,773	859,254
# of Edges	96,531	995,508
# of Graphs	1,000	2,000
Avg Degree	4.88	1.16

TABLE III
STATISTICS OF EGO-NETWORK.

Dataset	OAG	Twitter	Weibo	Digg
# of Nodes	953,675	456,626	1,776,950	279,630
# of Edges	4,151,463	12,508,413	308,489,739	1,548,126
# of Graphs	499,848	499,160	779,164	244,128
Avg Degree	4.35	27.39	173.60	5.54

IV. EXPERIMENTS

In this section, We first describe the experimental setting from Section 4.1 to Section 4.3, and then compare DGDA against state-of-the-art baselines in Section 4.4. To verify the effectiveness of latent variables regularization and the graph data augmentation in our model, we perform ablation studies

in Section 4.5. In Section 4.6, we further give an analysis about the sensitivity of the hyper-parameters. In Section 4.7, we describe the implementation details including the model architecture and the hyper-parameter setting. Finally, we analyse the time complexity of our DGDA in Section 4.8.

TABLE IV
NODE FEATURES OF IMDB&REDDIT DATASET.

Feature Type	Feature Name
Vertex	Coreness [45]. Pagerank [46]. Hub score and authority score [47]. Eigenvector Centrality [48]. Clustering Coefficient [49]. Degree

TABLE V
NODE FEATURES OF EGO-NETWORK DATASET.

Feature Type	Feature Name
Vertex	Coreness [45]. Pagerank [46]. Hub score and authority score [47]. Eigenvector Centrality [48]. Clustering Coefficient [49]. Rarity (reciprocal of ego user’s degree) [50].
Embedding	64-dim DeepWalk [51] embedding.
Ego-net	The number/ratio of active neighbors [52]. Density of sub-network induced by active neighbors [53]. Number of Connected components formed by active neighbors [53].

A. Datasets

1) *IMDB&Reddit Dataset*: It is composed of two datasets: **IMDB-Binary** and **Reddit-Binary**. We take each dataset as a domain. The statistics of the IMDB&Reddit dataset is shown in II.

- **REDDIT-BINARY** [21] is a balanced dataset with each graph corresponding to an online discussion thread where each node corresponds to a user ². An edge is drawn between two nodes if at least one of them responds to another’s comment. This dataset contains top submissions from four popular communities, namely, *IAmA*, *AskReddit*, *TrollXChromosomes*, *atheism*. *IAmA* and *AskReddit* are two question/answer-based communities and *TrollXChromosomes* and *atheism* are two discussion-based communities. The task is to identify whether a given graph belongs to a question/answer-based community or a discussion-based community.
- **IMDB-BINARY** [21] is a movie collaboration dataset ³. Each graph corresponds to an ego-network for each actor/actress, where nodes correspond to actors/actresses, and an edge is drawn between two actors/actresses if they appear in the same movie. Each graph is derived from a pre-specified genre (*Romance* or *Action*) of movies, and the task is to classify the genre graph it is derived from. Note that a movie can belong to both genres at

the same time. Therefore, [3] discarded movies from *Romance* genre if the movie is already included in the *Action* genre.

Note that the original datasets do not contain any node features, so we generate other six different node features based on the graph structure. All used features are shown in Table IV.

2) *Ego-network Dataset*: This dataset is first used for social influence prediction in [54] ⁴ and we use it for unsupervised graph domain adaptation. It is collected from four different social network platforms: **OAG**, **Digg**, **Twitter** and **Weibo**. Similarly, we take each social network platform as a domain. The statistics of Ego-network dataset is shown in Table III. The details of each domain are introduced as follows:

- **OAG**: OAG (Open Academic Graph) is generated by combining two academic networks: Microsoft Academic Graph and AMiner. The social network is defined to be the co-author network, and the social action is defined to be citation behaviors, i.e., a researcher cites a paper. It is collected from 20 conferences focusing on data mining, information retrieval, machine learning, natural language processing, computer vision, and database.
- **Digg**: Digg is a news aggregator application that allows users to vote web content, a.k.a, story, up or down. The social network is defined to be the Digg friendship network, and the social action is defined as voting for the story.
- **Twitter**: This dataset was built after monitoring the spreading processes on Twitter before, during, and after the announcement of the discovery of a new particle (Higgs boson) in Jul. 4th, 2012. Similarly, the social network is defined to be the Twitter friendship network, and the social action is defined to be whether a user retweets “Higgs” related tweets.
- **Weibo**: Weibo is the most popular microblog application in China. The complete dataset contains the following networks and posting logs of 1,776,950 users between Sep. 28th, 2012 and Oct. 29th, 2012. Social action is defined as reposting (the retweeting behavior in Weibo).

We follow the same data preprocess as [54]. A fixed-size sub-network is sampled from the ego user’s neighborhood via random walking. More concretely, each sample consists of a 50-node graph labeled by the ego user’s action status (active or inactive), and the task is to predict the ego users’ status. All used features are shown in Table V.

B. Baselines

We not only compare with the baselines of traditional domain adaptation tasks but also compare with the state-of-the-art domain adaptation methods for graph-structured data. Moreover, we also consider the disentanglement-based domain adaptation methods like DIVA [42] and DSR [35]. We try three different random seeds on all the methods so we report the mean and variance.

Traditional Domain Adaptation Methods:

²<http://networkrepository.com/REDDIT-BINARY.zip>

³<http://networkrepository.com/IMDB-BINARY.zip>

⁴<https://github.com/xptree/DeepInf>

- Target: A 3-layer GCN that is trained and tested on labeled data from the target domain, and can serve as a performance upper bound in domain adaptation tasks.
- Domain Adversarial Neural Network (**DANN**) [55]: An adversarial representation learning framework in which one classifier aims to distinguish the domain labels, while the feature extractor coupled with a gradient reverse layer tries to confuse the domain classifier.
- Margin Disparity Discrepancy (**MDD**) [36]: A state-of-the-art unsupervised domain adaptation method for computer vision task.

Graph Neural Network based Methods:

- Domain Adaptive Network Embedding (**DANE**) [7]: It achieves domain adaptive network embedding via shared weight graph convolutional network and adversarial learning regularization.
- Unsupervised Domain Adaptive Graph Convolutional Networks (**UDA-GCN**) [9]: It applies an attention mechanism to integrate global and local consistency and extract cross-domain node embedding with the help of Gradient Reversal Layer [55].

Disentanglement based Methods:

- Learning Disentangled Semantic Representation for Domain Adaptation (**DSR**) [35]: An unsupervised domain adaptation framework for computer vision task. It disentangles the input image into semantic latent variables and domain latent variables.
- Domain Invariant Variational Autoencoders (**DIVA**) [42]: A disentanglement based domain generalization method for medical image classification tasks.

C. Evaluation Metrics

We use F1-Score, the harmonic mean of the precision and recall, to quantitatively evaluate all methods, which is also widely used in other related works. It can be formulated as:

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} \quad (17)$$

D. Experimental Results

1) *IMDB&Reddit Result.*: The experiment results of the IMDB&Reddit dataset are shown in Table VI. Our method not only outperforms the conventional unsupervised domain adaptation methods but also outperforms the state-of-the-art graph domain adaptation methods on all the domain adaptation tasks. It’s remarkable that our method overpasses the comparison methods on task $I \rightarrow R$. However, the improvement of our method on $R \rightarrow I$ is not so remarkable. This is because the dataset size of Reddit-B is larger than that of IMDB-B. Furthermore, we can find that the performance of our method is very closed to that trained on the target domain (76.7), which is also another reason that leads to the indistinction of our method.

2) *Ego-network Result.*: The experiment results are shown in Table VII. Our method also outperforms all the other baselines on most tasks. We can find that our method achieves significant performance between the analogical domains, e.g., $T \rightarrow W$ and $W \rightarrow T$, which are both micro-blog platforms. What’s more, our method also performs well between the domains with large domain discrepancy like $O \rightarrow T$ and $O \rightarrow D$. In some tasks, the improvement is even larger than 4 points. However, the performance of our model on $T \rightarrow O$ is lower than that of UDAGCN and DIVA. According to Table III, domain OAG (O) is sparser than the other domains, so all the compared models perform badly when the target domain is OAG. We also conduct the Wilcoxon signed-rank test on the experiment result with five different random seeds. The p-value is 0.0216, which shows that our method significantly outperforms the baselines.

E. Ablation Study

To investigate the effectiveness of the latent variables regularization and the random latent variables, we further devise the following three ablated models.

- **DGDA-D**: We remove the maximum entropy loss on all latent variables.
- **DGDA-O**: We replace the augmented data with the original data.
- **DGDA-M**: We remove the noise variables reconstruction module.

The experiment results are shown Table VIII.

1) *Study of Latent Variables Regularization.*: To study the effectiveness of the latent variables regularization, we devise **DGDA-D** in which the regularization loss is removed. The experiment results are shown in Table VIII. According to the experiment results, we can find that the ablated model still achieves a comparable performance compared with the state-of-the-art baseline. But the performance drops in all the transfer learning tasks, which shows that the latent variables regularization can restrict these three types of latent variables and contribute to disentanglement.

2) *Study of Noise Variables Reconstruction.*: To study the effectiveness of manipulation latent variables, we devise **DGDA-O** and **DGDA-M**. Comparing with the result of **DGDA** and **DGDA-O**, we find that **DGDA** achieves a better result, which testifies that the data augmentation process is beneficial to learning the noise latent variables. Comparing the result of **DGDA** and **DGDA-M**, we also find **DGDA** performs better than **DGDA-M**, which testifies that the noise variables reconstruction module improves the model robustness and the transferring performance.

3) *Study of Graph Data Augmentation on Baselines.*: For a fair comparison, we applied the graph data augmentation to the baseline methods. The experiment results on the Ego-networks dataset are shown in Table IX.

According to the experiment results, most of the compared methods benefit from the augmented training data in different degrees. Especially, DSR gains an average improvement of more than 1 point and achieves the best performance in task $T \rightarrow O$. Even so, our proposed method DGDA still

TABLE VI
IMDB&REDDIT DATASET TEST SET F1-SCORE(%).

task	Target	DANN	MDD	DANE	UDAGCN	DSR	DIVA	DGDA(Ours)
I→R	90.0±1.6	74.4±0.6	67.3±2.0	73.9±1.1	73.7±0.9	75.7±1.4	75.9±0.7	78.1±1.0
R→I	76.7±2.6	73.8±1.4	71.0±0.8	72.1±1.6	73.9±0.7	73.2±1.1	74.0±1.3	74.7±0.9

TABLE VII
EGO-NETWORK DATASETS TEST SET F1-SCORE(%).

Methods	O→T	O→W	O→D	T→O	T→W	T→D	W→O	W→T	W→D	D→O	D→T	D→W	Avg
Target	53.3	53.2	62.9	43.0	53.2	62.9	43.0	53.3	62.9	43.0	53.3	53.2	53.1
DANN	46.7±0.5	41.4±0.4	42.9±1.4	40.2±0.6	42.5±0.8	41.1±0.4	40.2±0.2	47.1±0.3	41.8±0.6	40.2±0.1	45.0±0.3	41.6±0.2	42.6
MDD	47.1±0.7	40.6±0.5	39.8±1.2	40.1±0.5	43.1±0.8	40.6±0.9	40.0±0.1	47.5±0.2	40.7±0.7	40.3±0.1	45.1±0.2	40.2±0.1	42.2
DANE	42.2±0.6	40.1±0.6	46.0±1.9	40.2±0.2	42.3±0.6	40.8±1.9	40.2±0.2	47.7±0.5	42.6±1.4	39.9±0.1	43.3±2.8	40.2±0.4	42.1
UDAGCN	40.2±0.6	42.1±0.5	45.4±1.2	40.4±0.2	42.9±0.9	40.8±0.8	40.3±0.1	48.3±0.3	41.1±0.7	40.0±0.6	43.4±0.2	40.3±0.2	42.1
DSR	47.7±0.5	41.0±0.6	43.4±1.8	40.2±0.1	41.1±0.5	41.4±2.0	40.3±0.2	46.8±0.5	43.4±1.4	40.2±0.1	45.5±2.8	41.2±0.5	42.7
DIVA	47.1±0.5	42.0±0.1	44.6±0.3	40.7±0.4	42.3±0.5	45.1±1.0	40.8±0.2	48.5±0.1	43.7±2.2	40.9±0.4	46.4±0.7	41.7±0.1	43.6
DGDA(Ours)	49.2±0.1	42.9±0.4	49.7±1.4	40.3±0.4	44.7±0.6	46.7±1.7	41.5±0.3	49.9±0.3	48.3±1.3	41.1±0.4	48.4±0.6	42.7±0.8	45.5

TABLE VIII
THE ABLATION STUDY RESULTS ON EGO-NETWORK DATASETS.

Methods	O→T	O→W	O→D	T→O	T→W	T→D	W→O	W→T	W→D	D→O	D→T	D→W	Avg
Target	53.3	53.2	62.9	43.0	53.2	62.9	43.0	53.3	62.9	43.0	53.3	53.2	53.1
DGDA-D	48.7±0.3	42.3±0.6	49.0±1.8	40.3±0.1	42.2±0.3	45.1±1.9	41.0±0.4	48.3±0.4	47.6±0.9	41.0±0.4	47.7±0.8	42.0±0.6	44.6
DGDA-O	48.6±0.5	41.2±0.5	46.3±1.1	40.2±0.0	41.7±0.9	42.7±1.5	40.3±0.1	49.1±0.3	44.5±1.8	40.3±0.0	48.0±0.7	40.5±0.2	43.6
DGDA-M	48.3±0.1	41.0±0.6	47.1±2.6	40.2±0.0	42.1±0.7	43.0±2.3	40.4±0.4	48.8±0.2	44.7±1.6	40.2±0.0	47.6±0.4	41.0±0.6	43.7
DGDA	49.2±0.1	42.9±0.4	49.7±1.4	40.3±0.4	44.7±0.6	46.7±1.7	41.5±0.3	49.9±0.3	48.3±1.3	41.1±0.4	48.4±0.6	42.7±0.8	45.5

TABLE IX
THE ABLATION STUDY RESULTS OF GRAPH DATA AUGMENTATION PROCESS. ALL THE BASELINES ARE TRAINED WITH BOTH ORIGINAL AND AUGMENTED DATA.

Methods	O→T	O→W	O→D	T→O	T→W	T→D	W→O	W→T	W→D	D→O	D→T	D→W	Avg
Target	53.3	53.2	62.9	43.0	53.2	62.9	43.0	53.3	62.9	43.0	53.3	53.2	53.1
DANN+Augmented	46.8±0.5	41.0±0.2	43.3±0.1	40.2±0.1	42.4±0.3	41.6±0.1	40.2±0.0	47.1±0.3	43.3±1.5	40.1±0.0	46.5±0.3	41.6±0.1	42.8
MDD+Augmented	46.4±0.7	42.0±0.1	41.8±0.3	40.3±0.2	43.2±0.4	43.0±1.5	40.1±0.0	47.4±0.3	41.3±0.9	40.2±0.0	43.1±0.6	40.7±0.3	42.4
DANE+Augmented	42.7±0.5	40.2±0.6	46.1±1.6	40.1±0.7	41.9±1.7	41.4±1.2	40.0±0.1	48.0±0.5	43.5±1.2	40.3±0.2	43.5±2.0	41.0±0.3	42.3
UDAGCN+Augmented	40.3±0.7	42.5±0.8	47.0±1.2	40.2±0.5	42.9±0.5	40.6±0.4	40.3±0.1	48.2±0.3	42.0±0.4	40.1±0.1	43.1±0.4	40.7±0.6	42.3
DSR+Augmented	47.6±0.1	42.2±0.4	46.4±1.0	41.1±0.2	42.4±0.3	43.8±1.2	41.3±0.2	48.5±0.4	46.6±0.4	40.4±0.4	47.4±0.6	42.4±0.1	44.2
DIVA+Augmented	47.8±0.5	42.0±0.1	43.5±0.9	40.2±0.3	42.5±0.6	43.7±1.1	40.8±0.1	48.0±0.1	45.2±2.2	41.0±0.7	46.9±0.3	41.6±0.1	43.6
DGDA(Ours)	49.2±0.1	42.9±0.4	49.7±1.4	40.3±0.4	44.7±0.6	46.7±1.7	41.5±0.3	49.9±0.3	48.3±1.3	41.1±0.4	48.4±0.6	42.7±0.8	45.5

outperforms most of the baselines. This ablation study again validates the superior performance of our method for the graph unsupervised domain adaptation task.

F. Sensitivity Analysis

In this section, we investigate the sensitivity of hyperparameters. More specifically, we evaluate how the different dimensions of the semantic latent variables Z_y , and the latent variables regularization weight δ , as well as the different edge drop rates, affect the performance of **DGDA**. We report the test set F1 scores on the $W \rightarrow T$ and $T \rightarrow W$ of the ego-network datasets in Figure 3.

First, As shown in Figure 3(a), we try different dimensions of semantic latent variables Z_y from 32 to 300 while fixing the dimension of Z_d and Z_o . We observed that 256-dimension Z_y achieves the best performance on both two tasks. Furthermore, when the dimension is larger than 64, only slight differences can be observed with different numbers of dimensions. These results show that **DGDA** is stable with the dimension of Z_y .

Second, as shown in Figure 3(b) parameter δ denotes the weight of maximum entropy loss, which controls how

much information we exclude from the latent variables. We try different δ values from 0.1 to 10.0. We find that the performance is stable in both $W \rightarrow T$ and $T \rightarrow W$ when δ is larger than 1.0. In our experiment, the value of δ is 5.0.

Finally, Figure 3 (c)-(e) show the performance of **DGDA** in the condition of different edge perturbation ratio. We vary the perturbation rate from 0.0 to 0.6 and find that: (1) the F1 score increases in both tasks as the increment of the edge perturbation ratio and achieves the best result with a 0.1 to 0.3 edge perturbation ratio, which shows that considering the uncertainty of graph data is beneficial to the robustness of the graph neural networks. (2) the F1 score gradually drops when the edge perturbation ratio becomes larger. This is because a too large edge perturbation ratio leads to the destruction of semantic-related structures and results in the degeneration of the performance. As a result, we set $p_{\text{drop}} = 0.1$ and $p_{\text{add}} = 0.1$ in our experiments.

We also investigate the sensitivity of the loss weights γ , α , and ω in Equation (12) of the main manuscript. They correspond to the weight of \mathcal{L}_d , \mathcal{L}_y , and \mathcal{L}_o , respectively. The results are shown in Figure 4.

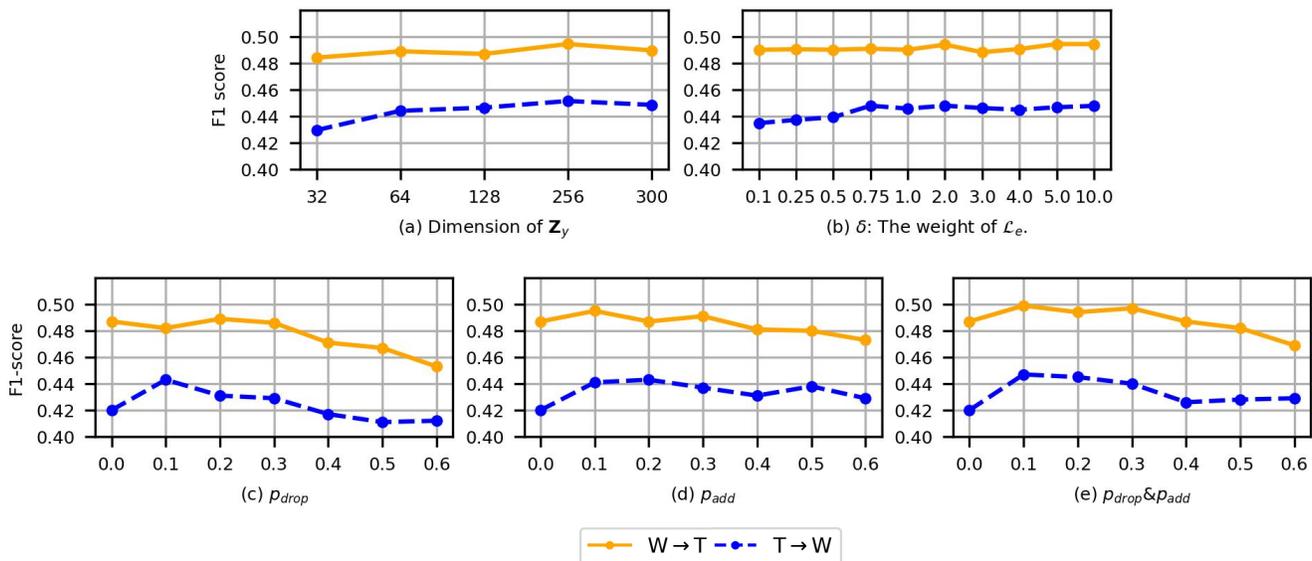


Fig. 3. Sensitivities of hyperparameters on the performance of DGDA.

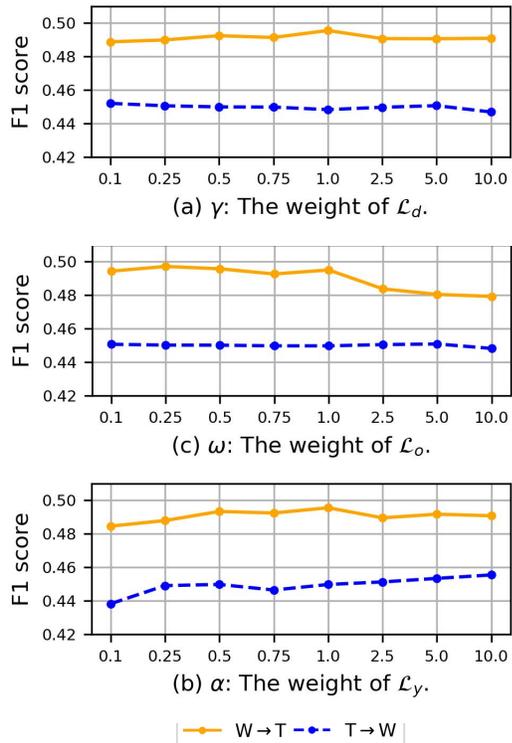


Fig. 4. Sensitivities of the loss weights on the performance of DGDA.

Firstly, as shown in Figure 4 (a), our DGDA is stable with γ . This is because it is easy to determine the domain label of the graph representation. We set $\gamma = 1$ in our experiment. Secondly, according to Figure 4 (b), we find that the performance drops around 0.01 on both tasks when the value α is small (≤ 0.25). As α increases, the performance increases and becomes more stable. We set $\alpha = 1$ in our experiment. Thirdly, as shown in Figure 4 (c), we find that the performance is stable when the value of ω is small (≤ 1.0)

and the performance may drop when ω becomes larger. This is because a large weight on \mathcal{L}_o submerges the gradient from \mathcal{L}_y and causes model collapse. We set $\omega = 0.1$ in our experiment.

G. Implementation Details

For fairness, we use the same architectures in our DGDA and all compared baselines, which is shown in Table X. The other hyper-parameters are introduced in Table XI.

TABLE X
THE COMPONENTS OF OUR DGDA AND ALL COMPARED BASELINES.

Component	Description
Feature extractor/Encoder	3-layer GCN.
Class classifier	2-layer MLP.
Domain classifier	One linear layer.
Decoder	2-layer MLP.
Readout function	Mean [56].
Activation function	ReLU [57].

TABLE XI
THE HYPER-PARAMETERS OF DGDA.

Parameter	IMDB&Reddit	Ego-network
Batch size	64	1024
Learning rate	0.001	0.01
Encoder hidden size	256	256
Dimension of Z_d	256	64
Dimension of Z_y	256	256
Dimension of Z_o	128	128
Decoder hidden size	64	64
Dropout rate	0.2	0.5
Weight decay	0.0005	0.0005
p_{drop}	0.1	0.1
p_{add}	0.1	0.1

H. Time Complexity Analysis

In this section, we briefly analyse the time complexity our DGDA. We denote the number of nodes in a graph

as N , the number of input, hidden, and output dimension as D_{in} , D_h , and D_{out} . Assuming that $N > D_{in}, D_h, D_{out}$, the computational complexity of one GCN layer is $O(N^3 + N^2D_{out} + ND_{in}D_{out}) = O(N^3)$, and one linear layer takes $O(ND_{in}D_{out})$. The time complexity of each block is shown as follows:

- Encoder Block: Assuming we employ L layer GCNs, its computational complexity is $O(LN^3)$.
- Decoder Block: it consists of a two-layer MLP and an inner product, so it costs $O(N^2D_{out} + ND_{in}D_h + ND_hD_{out}) = O(N^2D_{out})$.
- Disentanglement Block: Similarly, it costs $O(N^2D_{out})$.

As a summary, the time complexity of DGDA is $O(LN^3 + N^2D_{out} + N^2D_{out}) = O(LN^3)$, which is identical to the complexity of stacking L layer GCNs.

V. CONCLUSION

In this paper, we present a disentanglement-based framework for the graph unsupervised domain adaptation task. Starting from the generation process of graph-structured data, our approach extracts the disentangled semantic information on the recovered latent space. Specifically, we apply a variational graph auto-encoder architecture with three types of disentanglement modules. Experimental results on two sets of graph classification datasets show that our method outperforms existing methods for graph domain adaptation tasks, which further implies our proposal provides an effective solution for the graph data with high uncertainty and complexity. Following this principle, one promising direction is to go beyond GCN, to investigate more powerful graph neural network architectures for domain adaptation tasks. It would also be interesting to understand and visualize actually what substructure is transferable between graphs from different domains, to provide more interpretability.

REFERENCES

- [1] R. Cai, Z. Zhang, Z. Hao, and M. Winslett, "Understanding social causalities behind human action sequences," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 28, no. 8, pp. 1801–1813, 2017.
- [2] W. Chen, R. Cai, Z. Hao, C. Yuan, and F. Xie, "Mining hidden non-redundant causal relationships in online social networks," *Neural Computing and Applications*, vol. 32, no. 11, pp. 6913–6923, 2020.
- [3] X. Wang, L. Li, W. Ye, M. Long, and J. Wang, "Transferable attention for domain adaptation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 5345–5352.
- [4] B. Y. Lin and W. Lu, "Neural adaptation layers for cross-domain named entity recognition," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 2012–2022.
- [5] Z. Hao, D. Lv, Z. Li, R. Cai, W. Wen, and B. Xu, "Semi-supervised disentangled framework for transferable named entity recognition," *Neural Networks*, vol. 135, pp. 127–138, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0893608020304159>
- [6] R. Cai, J. Chen, Z. Li, W. Chen, K. Zhang, J. Ye, Z. Li, X. Yang, and Z. Zhang, "Time series domain adaptation via sparse associative structure alignment," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 8, pp. 6859–6867, May 2021. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/16846>
- [7] Y. Zhang, G. Song, L. Du, S. Yang, and Y. Jin, "Dane: domain adaptive network embedding," in *Proceedings of the 28th International Joint Conference on Artificial Intelligence*. AAAI Press, 2019, pp. 4362–4368.
- [8] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," *Advances in neural information processing systems*, vol. 27, pp. 2672–2680, 2014.
- [9] M. Wu, S. Pan, C. Zhou, X. Chang, and X. Zhu, "Unsupervised domain adaptive graph convolutional networks," in *Proceedings of The Web Conference 2020*, 2020, pp. 1457–1467.
- [10] Y. Ganin and V. Lempitsky, "Unsupervised domain adaptation by backpropagation," in *Proceedings of the 32nd International Conference on International Conference on Machine Learning-Volume 37*, 2015, pp. 1180–1189.
- [11] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.
- [12] T. N. Kipf and M. Welling, "Variational graph auto-encoders," *arXiv preprint arXiv:1611.07308*, 2016.
- [13] D. Wang, P. Cui, and W. Zhu, "Structural deep network embedding," in *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, 2016, pp. 1225–1234.
- [14] W. L. Hamilton, R. Ying, and J. Leskovec, "Inductive representation learning on large graphs," *arXiv preprint arXiv:1706.02216*, 2017.
- [15] R. Ying, R. He, K. Chen, P. Eksombatchai, W. L. Hamilton, and J. Leskovec, "Graph convolutional neural networks for web-scale recommender systems," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 974–983.
- [16] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, and Q. Mei, "Line: Large-scale information network embedding," in *Proceedings of the 24th international conference on world wide web*, 2015, pp. 1067–1077.
- [17] Y. Dong, N. V. Chawla, and A. Swami, "metapath2vec: Scalable representation learning for heterogeneous networks," in *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, 2017, pp. 135–144.
- [18] J. Qiu, Y. Dong, H. Ma, J. Li, K. Wang, and J. Tang, "Network embedding as matrix factorization: Unifying deepwalk, line, pte, and node2vec," in *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, 2018, pp. 459–467.
- [19] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.
- [20] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.
- [21] P. Yanardag and S. Vishwanathan, "Deep graph kernels," in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2015, pp. 1365–1374.
- [22] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How powerful are graph neural networks?" *arXiv preprint arXiv:1810.00826*, 2018.
- [23] O. Vinyals, S. Bengio, and M. Kudlur, "Order matters: Sequence to sequence for sets," *arXiv preprint arXiv:1511.06391*, 2015.
- [24] M. Zhang, Z. Cui, M. Neumann, and Y. Chen, "An end-to-end deep learning architecture for graph classification," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [25] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," *arXiv preprint arXiv:1606.09375*, 2016.
- [26] Z. Ying, J. You, C. Morris, X. Ren, W. Hamilton, and J. Leskovec, "Hierarchical graph representation learning with differentiable pooling," in *Advances in neural information processing systems*, 2018, pp. 4800–4810.
- [27] Y. Ma, S. Wang, C. C. Aggarwal, and J. Tang, "Graph convolutional networks with eigenpooling," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 723–731.
- [28] J. Lee, I. Lee, and J. Kang, "Self-attention graph pooling," in *International Conference on Machine Learning*. PMLR, 2019, pp. 3734–3743.
- [29] A. H. Khasahmadi, K. Hassani, P. Moradi, L. Lee, and Q. Morris, "Memory-based graph networks," *arXiv preprint arXiv:2002.09518*, 2020.
- [30] J. Huang, Z. Li, N. Li, S. Liu, and G. Li, "Attpool: Towards hierarchical feature representation in graph convolutional networks via attention mechanism," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6480–6489.
- [31] H. Gao and S. Ji, "Graph u-nets," in *international conference on machine learning*. PMLR, 2019, pp. 2083–2092.
- [32] D. Mesquita, A. H. Souza, and S. Kaski, "Rethinking pooling in graph neural networks," *arXiv preprint arXiv:2010.11418*, 2020.

- [33] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell, "Deep domain confusion: Maximizing for domain invariance," arXiv preprint arXiv:1412.3474, 2014.
- [34] M. Long, Y. Cao, J. Wang, and M. Jordan, "Learning transferable features with deep adaptation networks," in International conference on machine learning. PMLR, 2015, pp. 97–105.
- [35] R. Cai, Z. Li, P. Wei, J. Qiao, K. Zhang, and Z. Hao, "Learning disentangled semantic representation for domain adaptation," in Proceedings of the 28th International Joint Conference on Artificial Intelligence. AAAI Press, 2019, pp. 2060–2066.
- [36] Y. Zhang, T. Liu, M. Long, and M. Jordan, "Bridging theory and algorithm for domain adaptation," in International Conference on Machine Learning, 2019, pp. 7404–7413.
- [37] C. Shui, Z. Li, J. Li, C. Gagné, C. Ling, and B. Wang, "Aggregating from multiple target-shifted sources," arXiv preprint arXiv:2105.04051, 2021.
- [38] Z. Li, R. Cai, H. W. Ng, M. Winslett, T. Z. Fu, B. Xu, X. Yang, and Z. Zhang, "Causal mechanism transfer network for time series domain adaptation in mechanical systems," ACM Transactions on Intelligent Systems and Technology (TIST), vol. 12, no. 2, pp. 1–21, 2021.
- [39] X. Shen, Q. Dai, F.-I. Chung, W. Lu, and K.-S. Choi, "Adversarial deep network embedding for cross-network node classification," arXiv preprint arXiv:2002.07366, 2020.
- [40] M. Pilanci and E. Vural, "Domain adaptation on graphs by learning aligned graph bases," IEEE Transactions on Knowledge and Data Engineering, 2020.
- [41] E. Vural, "Domain adaptation on graphs by learning graph topologies: theoretical analysis and an algorithm," Turkish Journal of Electrical Engineering & Computer Sciences, vol. 27, no. 3, pp. 1619–1635, 2019.
- [42] M. Ilse, J. M. Tomczak, C. Louizos, and M. Welling, "Diva: Domain invariant variational autoencoders," in Medical Imaging with Deep Learning. PMLR, 2020, pp. 322–348.
- [43] C. Zhang, K. Zhang, and Y. Li, "A causal view on robustness of neural networks," arXiv preprint arXiv:2005.01095, 2020.
- [44] Y. Rong, W. Huang, T. Xu, and J. Huang, "Dropedge: Towards deep graph convolutional networks on node classification," in International Conference on Learning Representations, 2019.
- [45] V. Batagelj and M. Zaversnik, "An $o(m)$ algorithm for cores decomposition of networks," arXiv preprint cs/0310049, 2003.
- [46] L. Page, S. Brin, R. Motwani, and T. Winograd, "The pagerank citation ranking: Bringing order to the web." Stanford InfoLab, Tech. Rep., 1999.
- [47] S. Chakrabarti, B. E. Dom, S. R. Kumar, P. Raghavan, S. Rajagopalan, A. Tomkins, D. Gibson, and J. Kleinberg, "Mining the web's link structure," Computer, vol. 32, no. 8, pp. 60–67, 1999.
- [48] P. Bonacich, "Power and centrality: A family of measures," American journal of sociology, vol. 92, no. 5, pp. 1170–1182, 1987.
- [49] D. J. Watts and S. H. Strogatz, "Collective dynamics of 'small-world' networks," nature, vol. 393, no. 6684, pp. 440–442, 1998.
- [50] L. A. Adamic and E. Adar, "Friends and neighbors on the web," Social networks, vol. 25, no. 3, pp. 211–230, 2003.
- [51] B. Perozzi, R. Al-Rfou, and S. Skiena, "Deepwalk: Online learning of social representations," in Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining, 2014, pp. 701–710.
- [52] L. Backstrom, D. Huttenlocher, J. Kleinberg, and X. Lan, "Group formation in large social networks: membership, growth, and evolution," in Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining, 2006, pp. 44–54.
- [53] J. Ugander, L. Backstrom, C. Marlow, and J. Kleinberg, "Structural diversity in social contagion," Proceedings of the National Academy of Sciences, vol. 109, no. 16, pp. 5962–5966, 2012.
- [54] J. Qiu, J. Tang, H. Ma, Y. Dong, K. Wang, and J. Tang, "Deepinf: Social influence prediction with deep learning," in Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2018, pp. 2110–2119.
- [55] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," Journal of Machine Learning Research, vol. 17, pp. 1–35, 2016.
- [56] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Neural message passing for quantum chemistry," in ICML, 2017.
- [57] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in Icml, 2010.