Heterogeneous Information Network Embedding with Adversarial Disentangler

Ruijia Wang*, Chuan Shi*[†], Member, IEEE, Tianyu Zhao, Xiao Wang, Yanfang Ye

Abstract—Heterogeneous information network (HIN) embedding has gained considerable attention in recent years, which learns lowdimensional representation of nodes while preserving the semantic and structural correlations in HINs. Many of existing methods which exploit meta-path guided strategy have shown promising results. However, the learned node representations could be highly entangled for downstream tasks; for example, an author's publications in multidisciplinary venues may make the prediction of his/her research interests difficult. To address this issue, we develop a novel framework named HEAD (*i.e.*, HIN Embedding with Adversarial Disentangler) to separate the distinct, informative factors of variations in node semantics formulated by meta-paths. More specifically, in HEAD, we first propose the meta-path disentangler to separate node embeddings from various meta-paths into intrinsic and specific spaces; then with meta-path schemes as self-supervised information, we design two adversarial learners (*i.e.*, meta-path and semantic discriminators) to make the intrinsic embedding more independent from the designed meta-paths while the specific embedding more meta-path dependent. To comprehensively evaluate the performance of HEAD, we perform a set of experiments on four real-world datasets. Compared to the state-of-the-art baselines, the maximum 15% improvement of performance demonstrates the effectiveness of HEAD and the benefits of the learned disentangled representations.

Index Terms—Heterogeneous Information Networks, Network Embedding, Representation Disentanglement

1 INTRODUCTION

ETEROGENEOUS information networks (HINs) [1], involving different types of nodes and relations, are ubiquitous in the real world, ranging from bibliographic and social networks to transportation and telecommunication systems [2]. With heterogeneous types of nodes and relations, HINs are able to model complex interactions and immensely rich semantics in real-world scenarios. Thus, HIN analysis has emerged as a promising direction for many data mining tasks [3]. With the surge of network embedding [4], [5], numerous recent research has shifted towards HIN embedding [6], [7], [8], [9], [10], [11], which aims to project the nodes into a low-dimensional space whilst preserving the structural and semantic properties of HINs. The learned low-dimensional embedding has been a de facto solution towards fundamental problems in various applications, such as node classification [12], [13], link prediction [14], [15] and recommendation [16], [17].

HIN embedding has recently attracted considerable attention, the major line of which [7], [8], [16], [18] adopts meta-path [19] to retain both structural and semantic correlations between nodes. Specifically, meta-path is a typeconstrained relation sequence connecting two nodes, which is regarded as a basic tool in HINs to extract sub-structure and semantic. As illustrated in Figure 1 (a) and (b) for the bibliographic data, the toy HIN consists of multiple types of

• * Both authors contributed equally to this research.

 [†] indicates corresponding author. Ruijia Wang, Chuan Shi, Tianyu Zhao and Xiao Wang are with the Beijing Key Lab of Intelligent Telecommunications Software and Multimedia, Beijing University of Posts and Telecommunications, China. E-mail: {wangruijia,shichuan,tyzhao,xiaowang}@bupt.edu.cn

• Yanfang Ye is with the Department of Computer and Data Sciences, Case Western Reserve Univerity, USA. E-mail: yanfang.ye@case.edu.



Fig. 1: A toy example of heterogeneous information network and the two-step framework of meta-path based HIN embedding methods.

nodes and relations; meta-path APA and APCPA describe the co-author and co-conference relationship between two authors, respectively. These meta-path based HIN embedding methods either use only one meta-path inevitably resulting in information loss, or rely on certain strategies to merge multiple meta-paths, which often boil down to a twostep framework as displayed in Figure 1 (c). Firstly, they extract node sequences along each meta-path and embed them with optimization models, such as skip-gram [20]. Secondly, the node embeddings learned by multiple meta-paths are fused with some strategies (*e.g.*, concatenation or weighted combination [7], [18]) to obtain the final representations. Through the above operations, they produce structure- and semantic-preserving embeddings.

Although these meta-path based HIN embedding methods have demonstrated great success in learning representations, they may fail to uncover the latent explanatory factors hidden in the node embeddings under various meta-paths. Taking meta-paths APA and APCPA in Figure 1 as instances,

an author's publications with interdisciplinary authors or in multidisciplinary venues may make the prediction of his/her research interests difficult. To address this issue, we review the role of meta-path in HIN analysis, and have the following two insights. (1) Meta-paths are highly correlated, and their common part exposes intrinsic factors of nodes which absorb informative characteristic from all meta-paths and could more essentially reflect the characteristics. Metapaths APCPA and APA depict the micro co-author and macro co-conference relations respectively, but they actually capture the intrinsic correlation among authors, *i.e.*, the same research interests, from different perspectives. (2) Meta-paths describe distinguished sub-structures from various aspects, and their specific parts sometimes interfere with the analysis of nodes once connections along metapath are sparse or noisy [8], [16]. For instance, a broadsubject conference in APCPA may provide misleading information for predicting author's research interest, so does a broad-interest author in APA. Therefore, the mentioned highly-correlated but distinguished sub-structures explored by meta-paths provide an impetus to purify out intrinsic factors during meta-path fusion.

Nevertheless, the type sequences among meta-paths are overlapping somehow but different, revealing a high degree of semantic entanglement. Although aware of the necessity of extracting intrinsic factors, its materialization is nontrivial. (1) How to determine intrinsic and specific factors of nodes extracted by multiple meta-paths? From the perceptual understanding of intrinsic and specific factors, the intrinsic factors should be invariant to meta-paths and correspond to the essential feature of nodes, while the specific factors should depend on the meta-path specific semantic. For example, the author's research interests determine which conferences he will publish and which researchers he will collaborate with, thus should be invariant to meta-paths APCPA and APA. But misleading information brought by specific conferences and collaborators is strongly meta-path dependent. (2) How to disentangle intrinsic and specific factors without explicit supervised information? There are no supervised labels about intrinsic or specific factors. The author's research interests are sometimes expensive to obtain, let alone meta-path dependent misleading information. Hence, it is imperative to design a mechanism that utilizes the applicable selfsupervised information to guarantee the disentanglement.

In this paper, we make the first attempt to employ representation disentanglement on HINs to disentangle intrinsic and specific node embedding from multiple meta-paths. We propose HEAD, a novel HIN Embedding framework with Adversarial Disentangler. To detach intrinsic and specific embedding, a meta-path disentangler is designed to encode node embedding generated by each meta-path into intrinsic and specific space. Furthermore, utilizing metapath schemes as self-supervised information, two adversarial learners (i.e., meta-path and semantic discriminator) are deployed to try the best to separate intrinsic and specific embedding. Specifically, the meta-path discriminator learns a more intrinsic embedding while the semantic discriminator enhances the correlation between specific embedding and meta-path semantic. The major contributions of our work are summarized as followings.

To the best of our knowledge, we are the first to investi-

gate intrinsic and specific factors of nodes in HIN embedding and propose a disentangled representation solution to purify out intrinsic embedding. Purifying HIN embedding is important for HIN analysis, because it alleviates meta-path dependence and reduces noise interference.

- We devise a novel framework HEAD that disentangles intrinsic embedding from node embeddings generated by various meta-paths with adversarial disentangler, which boosts the robustness of representation. Besides, HEAD is an unsupervised learner to effectively absorb intrinsic factors from various meta-paths, thus enhancing existing meta-path based embedding methods.
- We perform experiments on four public datasets. Compared to state-of-the-art baselines, the maximum 15% improvement demonstrates the effectiveness of HEAD and the benefit of intrinsic representations.

2 RELATED WORK

We review the most related works in network embedding, HIN embedding and representation disentanglement.

2.1 Network embedding

Network embedding, *i.e.*, network representation learning (NRL), is proposed to learn structure-preserving node representations which can be applied to downstream network tasks [4], [5], [21], [22].

Contemporary methods usually explore network topology as context information. DeepWalk [23] and node2vec [24] construct node sequences by randomly walking on the network, then utilize skip-gram based models to learn node embeddings. LINE [25] and SDNE [26] characterize firstorder and second-order proximity to preserve neighborhood information for nodes. Furthermore, GraRep [27] and HOPE [28] are both designed to model the high-order proximity between nodes. Besides learning network embedding only from topology, there are also some works leveraging node content information [29], [30] or temporal information [31], [32], [33] for a more robust representation. Recently, graph neural networks [34], [35], [36] have demonstrated their remarkable ability in network representation learning, which learn a function that generates embeddings by aggregating feature from local neighborhood. For example, Graph Convolutional Networks [35] performs localized first-order approximation of spectral graph convolutions. Inspired by attention mechanism [37], [38], Graph Attention Networks [39] is proposed to distinguish the importance between neighbors and fuse them to obtain final node embeddings.

Nonetheless, these methods mainly deal with homogeneous networks, and they cannot be directly applied to HINs which contain multiple types of nodes and relations.

2.2 HIN embedding

In the past decade, heterogeneous information networks (HINs) [1] have been proposed to model complex entities and their rich relations in various applications [17], [40].

One major line of work leverages meta-path to learn semantic-preserving embedding for HINs. ESim [7] accepts predefined meta-paths as guidance to learn node embedding for similarity search. Even though ESim utilize multiple meta-paths, it needs to conduct grid search to find the optimal weights of meta-paths. Metapath2vec [8] and HERec [16] design meta-path based random walk and utilize skip-gram to perform HIN embedding. HIN2Vec [41] carries out multi-task classification which learns the latent embeddings of nodes and meta-paths simultaneously. Based on the attention mechanism, HAN [18] proposes node- and semantic-level attention to learn the importance of nodes and meta-paths, respectively. In summary, these aforementioned meta-path based algorithms do not distinguish the intrinsic and specific factors of nodes when fusing multiple meta-paths, and may be susceptible to noise.

On another line, there exist some methods utilizing relation type rather than meta-path to perform HIN embedding. PME [42] projects different types of nodes into the same relation space and conducts heterogeneous link prediction. Based on the adversarial principle, HeGAN [43] designs relation-aware discriminator and generator to learn node distribution for better negative samples. HetGNN [44] introduces a random walk with restart strategy to sample a fixed size of neighbors, and aggregate feature information from them based upon node types. It is worth noting that meta-path has ability to accurately capture corresponding semantic, but these meta-path free methods may lose the advantage of meta-path, making the semantic fuzzy, which may harm the performance of specific relation prediction.

2.3 Representation disentanglement

Recently, representation disentanglement has gained considerable attention, particularly in the field of image representation learning [45], [46], [47], [48]. It aims to disentangle latent factors from image variants, which leads to the understanding of observed data. For instance, [49] learns invertible graphic codes for 3D image rendering in a fully supervised setting. InfoGAN [45] decomposes representation by maximizing the mutual information between latent factors and synthesized images. Such disentangled representations are demonstrated to bring enhanced generalization ability and robustness in downstream tasks [50], [51].

However, the complex network structure makes representation disentanglement on networks rarely explored. The only studies [52], [53] focus on homogeneous networks to identify the latent factor that may cause the link between the node pair. These methods ignore the heterogeneity of networks, and depart from our goal to disentangle intrinsic and specific factors for semantic-rich HINs.

3 PRELIMINARY

In this section, we formalize the problem of HIN embedding and introduce the background on VAEs and GANs.

HIN embedding. Heterogenous information network is a special kind of information network, which contains either multiple types of nodes or relations. It can be formally defined as follows.

Definition 1. Heterogeneous Information Network (HIN) [1]. An HIN $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{A}, \mathcal{R}, \phi, \varphi)$ is a form of information network, where \mathcal{V} and \mathcal{E} denote sets of nodes and relations, respectively. It is also associated with a node type mapping function $\phi : \mathcal{V} \to \mathcal{A}$ and an relation type mapping function $\varphi : \mathcal{E} \to \mathcal{R}$, where \mathcal{A} and \mathcal{R} denote sets of node and relation types such that $|\mathcal{A}| + |\mathcal{R}| > 2$.

A toy example of HIN is illustrated in Figure 1 for bibliographic data. It is observed that it consists of three types of nodes (*i.e.*, author, paper and conference) and their semantic information based on meta-paths [19] (*e.g.*, author-paper-author).

The goal of HIN embedding is to embed each node $v \in \mathcal{V}$ to a low-dimensional space \mathbb{R}^d ($d \ll |\mathcal{V}|$) while preserving structure and rich semantic on the HINs so that the learned embeddings can be applied to downstream tasks.

Variational Autoencoders. VAEs [54], [55] are deep generative models that simultaneously train both a probabilistic encoder and generator for a dataset $\mathcal{D} = \{\mathbf{x}^1, \dots, \mathbf{x}^N\}$. The central analogy is that an encoding \mathbf{z} can be considered as a latent variable, casting the generator as a conditional probability density $p_{\theta}(\mathbf{x}|\mathbf{z})$, where θ denotes parameters of the generator. By placing a weak prior $p(\mathbf{z})$ over \mathbf{z} , the generator defines a posterior and joint distribution $p(\mathbf{z}|\mathbf{x}) \propto$ $p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$. And inference in VAEs can be performed using a variational method that approximates the posterior distribution $p(\mathbf{z}|\mathbf{x})$ using an encoder $q_{\phi}(\mathbf{z}|\mathbf{x})$, where ϕ denotes parameters of the encoder. The encoder and the generator are trained jointly by performing stochastic gradient ascent [54] on the *evidence lower bound* $\mathcal{L}(\phi, \theta; \mathcal{D})$ (ELBO),

$$\mathcal{L}(\phi, \theta; \mathcal{D}) = \sum_{n=1}^{N} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^{n})} \left[\log p_{\theta}\left(\mathbf{x}^{n}|\mathbf{z}\right) + \log p(\mathbf{z}) - \log q_{\phi}\left(\mathbf{z}|\mathbf{x}^{n}\right) \right]$$

Typically, the first term $\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^n)} [\log p_{\theta}(\mathbf{x}^n|\mathbf{z})]$ is approximated by a Monte Carlo estimate and the remaining two terms are expressed as a divergence $-\text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^n) || p(\mathbf{z}))$, which can be computed analytically when the encoder model and the prior are *Gaussian*.

Generative Adversarial Networks. Our model is based on adversarial learning, while the typical application of adversarial learning is GANs, which can be viewed as a minimax game between two players, *i.e.*, generator G and discriminator D, in the following manner.

$$\min_{\theta^{G}} \max_{\theta^{D}} \qquad \mathbb{E}_{x \sim P_{data}} \left[\log D(x; \theta^{D}) \right] \\ + \qquad \mathbb{E}_{z \sim P_{Z}} \left[\log \left(1 - D(G(z; \theta^{G}); \theta^{D}) \right) \right]$$

The generator G tries to generate fake samples as close to true data as possible with the noise z from a predefined distribution P_Z , where θ^G denotes parameters of the generator. On the contrary, the discriminator D aims to distinguish real data from fake samples, where θ^D represents parameters of the discriminator. In practice, GANs has been found to work better if the generator minimizes $-\log D(G(\cdot; \theta^G); \theta^D)$ instead of $\log(1 - D(G(\cdot; \theta^G); \theta^D))$ [56].

4 HEAD: THE PROPOSED MODEL

In this section, we present the proposed model HEAD, a novel HIN embedding framework with adversarial disentangler. We begin with the overview of HEAD. Subsequently we zoom into the meta-path disentangler, followed by elaborations on adversarial learners that guarantee disentanglement. Lastly, we discuss the optimization and characteristics



Fig. 2: Overview of our proposed model HEAD for disentangled representation learning on HINs, where takes two meta-paths APA and APCPA as an example. Note that the diagrams of encoder refer to the same one, as well as the generator.

Notion	Explanation
D_{mp}, D_{sem}	meta-path and semantic discriminator, resp.
$\mathbf{x}^1, \mathbf{\hat{x}}^1, \mathbf{\bar{x}}^1$	input, reconstructed and synthesized node presentation of meta-path APA, resp.
$\mathbf{x}^2, \mathbf{\hat{x}}^2, \mathbf{\bar{x}}^2$	input, reconstructed and synthesized node presentation of meta-path APCPA, resp.
$\mathbf{e}_{in}^1, \mathbf{e}_{sp}^1$	intrinsic and specific embedding under APA, resp.
$\mathbf{e}_{in}^2, \mathbf{e}_{sp}^2$	intrinsic and specific embedding under APCPA, resp.
$\hat{l}_{mp}, \hat{l}_{sem}$	meta-path type predicted by meta-path classifier and sem- antic classifier, resp.
\bar{l}_{mp}	meta-path type of the randomly selected specific embedding
$\theta_E, \theta_{lkp}, \theta_G$	parameters of encoder, embedding lookup layer and gener- ator, resp.
$\theta_{mp}, \theta_{sem}$	parameters of meta-path and semantic discriminator, resp.

TABLE 1: Summary of notations.

of our model. Notations we use throughout the paper are summarized in Table 1.

4.1 Overview

In order to detach intrinsic and specific factors from multiple meta-paths, we propose a disentangled representation solution which firstly coarsely detaches intrinsic and specific embeddings under each meta-path with a meta-path disentangler, then refine them with adversarial learners to make intrinsic embedding more meta-path independent and specific embedding more meta-path dependent. Figure 2 shows the overall framework.

Given an HIN and a set of meta-paths, HEAD first obtains node embedding for each meta-path with a meta-path based embedding method (*e.g.*, Metapath2vec or HERec) as input. For the node embedding \mathbf{x}^{ϕ} of a meta-path ϕ , the meta-path disentangler employs a modified VAE architecture to detach the intrinsic embedding \mathbf{e}_{in}^{ϕ} and specific embedding \mathbf{e}_{sp}^{ϕ} from \mathbf{x}^{ϕ} , which concatenates specific embedding \mathbf{e}_{sp}^{ϕ} in the middle layer of VAE to reconstruct the original node embedding \mathbf{x}^{ϕ} . As exemplified in Figure 2, \mathbf{x}^{1} is reconstructed from $\mathbf{e}_{in}^{1} || \mathbf{e}_{sp}^{1}$.

In order to refine the intrinsic and specific embedding, we design two adversarial learners. Specifically, a meta-path discriminator D_{mp} makes the intrinsic embedding more meta-path independent and a semantic discriminator D_{sem} makes the specific embedding more meta-path dependent. Since there is no label information to supervise the disentangled learning in our problem setting, we leverage metapath schemes as self-supervised information to distinguish intrinsic and specific embedding. More precisely, to purify the intrinsic embedding \mathbf{e}_{in}^{ϕ} , the meta-path discriminator D_{mp} makes \mathbf{e}_{in}^{ϕ} more difficult to be distinguished which meta-path it belongs to. To constrain specific embedding \mathbf{e}_{sp}^{ϕ} to be more meta-path dependent, the semantic discriminator D_{sem} is designed to distinguish a synthesized semantic, generating by the concatenation of intrinsic embedding and specific embedding under different meta-paths (e.g., $\bar{\mathbf{x}}^2$ from $\mathbf{e}_{in}^1 || \mathbf{e}_{sn}^2$). There are two classifiers inside semantic discriminator \hat{D}_{sem} : (1) the semantic classifier forces the synthesized semantic belonging to the meta-path type of specific embedding (e.g., $\bar{\mathbf{x}}^2$ belongs to APCPA), which makes the specific embedding more meta-path dependent. (2) The real/fake classifier guarantees the quality of synthesized semantic by adversarial learning, since bad synthesized semantic can be easily tested by the semantic classifier. Note that metapath disentangler and adversarial learners are shared by all meta-paths. The parameters of these learners are iteratively training as other adversarial models.

4.2 Meta-path Disentangler

Given an HIN \mathcal{G} and a series of meta-paths $(\phi^1, \phi^2, \dots, \phi^K)$, we first learn node embeddings under each meta-path with existing meta-path based embedding methods (*e.g.*, Mp2vec and HERec) as input of HEAD. Note that we employ Mp2vec in the experiments to obtain the input node embeddings $(\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^K)$ for meta-path set $(\phi^1, \phi^2, \dots, \phi^K)$.

5

Similar experimental results can also be obtained for other meta-path based HIN embedding methods (*e.g.*, HERec).

With node embeddings of various meta-paths as input, the meta-path disentangler coarsely detaches intrinsic and specific embedding for each meta-path. Specifically, taking the node embedding \mathbf{x}^{ϕ} of a meta-path ϕ as input, the encoder derives the intrinsic embedding \mathbf{e}_{in}^{ϕ} of \mathbf{x}^{ϕ} , then concatenates it with the specific embedding \mathbf{e}_{sp}^{ϕ} from the embedding lookup layer to reconstruct the original node embedding \mathbf{x}^{ϕ} via generator.

It is worth noting that there exists uncertainty in the intrinsic embedding. For example, we cannot pinpoint the research interest of a broad-interest author. To model uncertainty, we are inspired by the development of variational autoencoders (VAEs) [54], [55], [57] to encode intrinsic embedding \mathbf{e}_{in}^{ϕ} into random variables with variance, rather than fixed values. Furthermore, we hypothesize that the intrinsic embeddings under various meta-paths have similar form inside the input embeddings and could be captured by the same encoder, but the specific embeddings are so diverse that they could be generated from the embedding lookup layer with random initialization. Thus we modify the VAE architecture as meta-path disentangler, which concatenates specific embedding \mathbf{e}_{sp}^{ϕ} from the middle layer, and the encoder and generator are implemented using multilayer perceptron (MLP). The objective function of meta-path disentangler is defined as:

$$\mathcal{L}_{vae} = \|\hat{\mathbf{x}}^{\phi} - \mathbf{x}^{\phi}\|_{F}^{2} + KL\left(q(\mathbf{e}_{in}^{\phi}|\mathbf{x}^{\phi})\|p(\mathbf{e}_{in}^{\phi})\right), \qquad (1)$$

where the first term aims at recovering the original node embedding \mathbf{x}^{ϕ} , and the second term calculates *Kullback-Leibler divergence* which penalizes the deviation of intrinsic embedding \mathbf{e}_{in}^{ϕ} from the prior distribution $p(\mathbf{e}_{in}^{\phi})$ (as $\mathbf{e}_{in}^{\phi} \sim \mathcal{N}(0, \sigma)$). The loss function in Eq. (1) retains the reconstruction ability and aligns the latent intrinsic embedding to a Gaussian distribution. However, this property cannot guarantee that intrinsic embedding \mathbf{e}_{in}^{ϕ} is well disentangled from the specific embedding \mathbf{e}_{sp}^{ϕ} .

4.3 Adversarial Learners

In order to better detach intrinsic and specific embeddings, we need to try the best to differentiate them with some guidance. However, there is no supervised information in our problem setting. Considering that embeddings are learned from diverse meta-paths, we utilize the meta-path type (denoted as l_{mp}) as self-supervised information to guide the learning. Moreover, enlightened by the adversarial learning [45], we design two adversarial learners to refine the disentanglement, including a meta-path discriminator purifying intrinsic embedding and a semantic discriminator specializing specific embedding.

4.3.1 Meta-path Discriminator

Intrinsic embedding does not depend on meta-paths and should be invariant to meta-path types. Thus to further refine the intrinsic embedding, we deploy meta-path discriminator D_{mp} in intrinsic space, which makes it difficult to tell which meta-path the intrinsic embedding comes from.

More precisely, the meta-path classifier inside D_{mp} takes intrinsic embedding \mathbf{e}_{in}^{ϕ} as input and predict the meta-path type l_{mp} . In the adversarial setting, the encoder attempts to purify out intrinsic embedding \mathbf{e}_{in}^{ϕ} that invariant to metapath type l_{mp} , whereas the meta-path discriminator tries to distinguish meta-path type l_{mp} from encoded intrinsic embedding \mathbf{e}_{in}^{ϕ} . The better trained meta-path discriminator would then force the encoder to produce more intrinsic embedding, and the process is repeated. During such iterations, both the encoder and meta-path discriminator receive mutual, positive reinforcement. Therefore, the objective functions of meta-path discriminator \mathcal{L}_{mp} and encoder \mathcal{L}_{E}^{adv} are derived as follows:

$$\mathcal{L}_{mp} = \mathbb{H}\left[\log P\left(\hat{l}_{mp} = l_{mp} | \mathbf{e}_{in}^{\phi}\right)\right], \quad (2)$$

$$\mathcal{L}_{E}^{adv} = -\mathcal{L}_{mp} = -\mathbb{H}\left[\log P\left(\hat{l}_{mp} = l_{mp}|\mathbf{e}_{in}^{\phi}\right)\right], \quad (3)$$

where *P* is probability distribution over meta-paths, which is produced by meta-path classifier, l_{mp} is the ground-truth meta-path type of \mathbf{e}_{in}^{ϕ} , and \mathbb{H} represents the entropy.

4.3.2 Semantic Discriminator

To force specific embedding more meta-path dependent, we introduce semantic discriminator D_{sem} , including a semantic classifier and a real/fake classifier, to distinguish synthesized semantics generated by the concatenation of intrinsic embedding and specific embedding under different meta-paths.

Semantic Classifier. For meta-paths $(\phi^1, \phi^2, \cdots, \phi^K)$, we shuffle $(\mathbf{e}_{sp}^1, \mathbf{e}_{sp}^2, \cdots, \mathbf{e}_{sp}^K)$ to $(\mathbf{e}_{sp}^{i_1}, \mathbf{e}_{sp}^{i_2}, \cdots, \mathbf{e}_{sp}^{i_K})$, and concatenate them with intrinsic embeddings $(\mathbf{e}_{in}^1, \mathbf{e}_{in}^2, \cdots, \mathbf{e}_{in}^K)$, which forms the concatenation of each pair $\mathbf{e}_{in}^{k} || \mathbf{e}_{sp}^{i_{k}}$ to generate synthesized semantics. Although generated by intrinsic embedding and specific embedding under different metapaths, a good synthesized semantic should be accurately distinguished as the meta-path type of specific embedding, because a good specific embedding is meta-path dependent while a good intrinsic embedding is meta-path independent. Thus the semantic classifier forces synthesized semantic belonging to the meta-path type of specific embedding. Without loss of generalization, we take two meta-paths in Figure 2 as an example to illustrate the semantic classifier. The synthesized semantic $\bar{\mathbf{x}}^2$ is generated from $\mathbf{e}_{in}^1 || \mathbf{e}_{sp'}^2$ and the semantic classifier forces it to be APCPA type, the objective function of which is derived as:

$$\mathcal{L}_{sem}^{clf} = \mathbb{H} \left[\log P \left(\hat{l}_{sem} = l_{mp} | \mathbf{x}^2 \right) \right] \\ + \mathbb{H} \left[\log P \left(\hat{l}_{sem} = \bar{l}_{mp} | \bar{\mathbf{x}}^2 \right) \right], \tag{4}$$

where \bar{l}_{mp} is the meta-path type of specific embedding (*i.e.*, APCPA) and \hat{l}_{sem} is the meta-path type prediction of semantic classifier.

Real/Fake Classifier. We introduce a real/fake classifier to distinguish real and synthesized semantics, which compels the generator to mimic the real semantic, so as to reduce the impact of bad synthesized semantic on classification. The objective functions $\mathcal{L}_{sem}^{r/f}$ and \mathcal{L}_{G}^{adv} of classifier and generator are defined as:

$$\mathcal{L}_{sem}^{r/f} = \mathbb{H}\left[\log\left(C_{r/f}(\bar{\mathbf{x}}^2)\right)\right] + \mathbb{H}\left[\log\left(1 - C_{r/f}(\mathbf{x}^2)\right)\right], \quad (5)$$
$$\mathcal{L}_{G}^{adv} = -\mathbb{H}\left[\log\left(C_{r/f}(\bar{\mathbf{x}}^2)\right)\right], \quad (6)$$

Algorithm 1: Model training for HEAD

Input : HIN \mathcal{G} , node embeddings $(\mathbf{x}^1, \mathbf{x}^2, \cdots, \mathbf{x}^K)$ under meta-path set $(\phi^1, \phi^2, \cdots, \phi^K)$, prior distribution $\mathcal{N}(0, \sigma)$ Output: intrinsic embeddings $(\mathbf{e}_{in}^1, \mathbf{e}_{in}^2, \cdots, \mathbf{e}_{in}^K)$, specific embeddings $(\mathbf{e}_{sp}^1, \mathbf{e}_{sp}^2, \cdots, \mathbf{e}_{sp}^K)$

1 Initialize θ_E , θ_{lkp} , θ_G , θ_{mp} and θ_{sem} ; 2 while not converge do // Meta-path discriminator training 3 Encode $(\mathbf{x}^1, \mathbf{x}^2, \cdots, \mathbf{x}^K)$ into $(\mathbf{e}_{in}^1, \mathbf{e}_{in}^2, \cdots, \mathbf{e}_{in}^K)$; 4 Update θ_{mp} with Eq. 7; 5 // Semantic discriminator training 6 Lookup $(\mathbf{e}_{sp}^1, \mathbf{e}_{sp}^2, \cdots, \mathbf{e}_{sp}^K)$; Shuffle $(\mathbf{e}_{sp}^1, \mathbf{e}_{sp}^2, \cdots, \mathbf{e}_{sp}^K)$ to $(\mathbf{e}_{sp}^{i_1}, \mathbf{e}_{sp}^{i_2}, \cdots, \mathbf{e}_{sp}^{i_K})$; Generate $(\bar{\mathbf{x}}^1, \bar{\mathbf{x}}^2, \cdots, \bar{\mathbf{x}}^K)$ for each pair $\mathbf{e}_{in}^k || \mathbf{e}_{sp}^{i_k}$; 7 8 9 Update θ_{sem} with Eq. 8; 10 // Meta-path disentangler training 11 Reconstruct $(\hat{\mathbf{x}}^1, \hat{\mathbf{x}}^2, \cdots, \hat{\mathbf{x}}^K)$; 12 Update θ_E , θ_{lkp} and θ_G with Eq. 9, Eq. 10 and 13 Eq. 11, respectively; 14 return $(\mathbf{e}_{in}^1, \mathbf{e}_{in}^2, \cdots, \mathbf{e}_{in}^K)$ and $(\mathbf{e}_{sp}^1, \mathbf{e}_{sp}^2, \cdots, \mathbf{e}_{sp}^K)$;

where $C_{r/f}()$ denotes the prediction result of real/fake classifier.

4.4 Optimization

Following the optimization technique in adversarial learning, we iteratively optimize three main modules to train HEAD. In each iteration, we alternately update meta-path discriminator D_{mp} , semantic discriminator D_{sem} and metapath disentangler. For meta-path discriminator D_{mp} , it needs to accurately classify intrinsic embedding, so the gradient is:

$$\theta_{mp} \stackrel{+}{\leftarrow} -\Delta_{\theta_{mp}} \left(\mathcal{L}_{mp} \right). \tag{7}$$

With regard to semantic discriminator D_{sem} , its purpose is to identify the meta-path type of synthesized semantic, as well as the real/fake label:

$$\theta_{sem} \stackrel{+}{\leftarrow} -\Delta_{\theta_{sem}} \left(\mathcal{L}_{sem}^{clf} + \mathcal{L}_{sem}^{r/f} \right). \tag{8}$$

Meta-path disentangler includes encoder, embedding lookup layer and decoder. First of all, their common goal is to reconstruct input node embedding. Then, the encoder also makes intrinsic embedding difficult to distinguish the meta-path type, while the embedding lookup layer makes specific embedding richer with meta-path specific information and the generator ensures the quality of synthesized semantics, so their gradients are as follows:

$$\theta_E \stackrel{+}{\leftarrow} -\Delta_{\theta_E} \left(\mathcal{L}_{vae} + \mathcal{L}_E^{adv} \right), \tag{9}$$

$$\theta_{lkp} \stackrel{+}{\leftarrow} -\Delta_{\theta_E} \left(\mathcal{L}_{vae} + \mathcal{L}_{sem}^{clf} \right), \tag{10}$$

$$\theta_G \stackrel{+}{\leftarrow} -\Delta_{\theta_G} \left(\mathcal{L}_{vae} + \mathcal{L}_G^{adv} + \mathcal{L}_{sem}^{clf} \right), \tag{11}$$

where θ denotes the parameters of corresponding modules. The training framework of HEAD is outlined in Algorithm

TABLE 2: Description of datasets.

Dataset	Relations (A-B)	#A	#B	#A-B
Yelp	Business-User	2, 614	1,286	30, 838
	Business-Service	2, 614	2	2, 614
	Business-Reservation	2, 614	2	2, 614
ACM	Paper-Author	4, 019	7, 167	13, 407
	Paper-Subject	4, 019	60	4, 019
DBLP	Paper-Author	14, 376	14, 475	41, 794
	Paper-Conference	14, 376	20	14, 376
	Paper-Term	14, 376	8, 920	114, 624
AMiner	Paper-Author	28, 722	46, 583	78, 767
	Paper-Reference	28, 722	95, 970	250, 012

1. Note that we concatenate the intrinsic embeddings of all meta-paths as the final intrinsic embedding for HEAD.

4.5 Discussion

We now conduct a complexity analysis. The complexity of HEAD lies in the iterative updating process. Each iteration mainly involves the updating of projection matrices in the meta-path disentangler and discriminators. It is worth noting that the meta-path disentangler and discriminators are shared by all meta-paths. Therefore, the time complexity is $O(K \cdot |\mathcal{V}| \cdot d^2)$, where *K* is the number of input meta-paths, $|\mathcal{V}|$ is the number of target-type nodes and *d* is the embedding dimension. Although HEAD adds complexity in iterative updating, this updating process is still efficient, which is linear with the number of target-type nodes.

Our HEAD provides a flexible framework to absorb multiple meta-paths for robust node embedding and can be seen as a fine-tuning process of existing meta-path based methods, which arms them with the ability to purify out intrinsic factors from each meta-path. Thus, HEAD is capable to enhance existing meta-path based methods.

5 EXPERIMENTS

In this section, we evaluate the effectiveness of HEAD on four tasks, including node classification, node clustering, relation prediction and network visualization. We further analyze the underlying mechanism inside HEAD. Lastly, we investigate the hyper-parameter sensitivity.

5.1 Experimental Setup

5.1.1 Datasets

We conduct extensive experiments on four heterogeneous information networks, as summarized in Table 2.

- Yelp [43]. This dataset consists of four types of nodes, *i.e.*, business (B), user (U), service (S) and reservation (R), where businesses are divided into three classes. The metapaths we use are {*UBU*, *UBSBU*, *UBRBU*}.
- ACM [18]. This dataset contains three types of nodes, *i.e.*, author (A), paper (P) and subject (S), where papers are divided into three classes according to the conferences they published on. Here we employ the meta-path set {*PAP*, *PSP*} to perform experiments.
- **DBLP** [41]. This dataset involves four types of nodes, *i.e.*, author (A), paper (P), conference (C) and term (T),

TABLE 3: Quantitative results ($\% \pm \sigma$) on node classification. (bold: best; underline: runner-up)

Datasets	Metrics	Training	DeepWalk	LINE-1st	LINE-2nd	ESim	HERec	Mp2vec	HIN2Vec	HetGNN	HeGAN	HEAD
Yelp	Macro-F1	20% 40% 60% 80%	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} \underline{75.55} \pm 1.23 \\ \overline{76.60} \pm 0.95 \\ \overline{77.60} \pm 0.82 \\ \overline{77.61} \pm 1.33 \end{array}$	$\begin{array}{c} 69.79{\pm}0.90\\ 72.57{\pm}0.82\\ 73.57{\pm}1.04\\ 74.50{\pm}1.58\end{array}$	$\begin{array}{c} 74.83{\pm}1.28\\ \underline{76.72}{\pm}1.19\\ \underline{77.86}{\pm}1.41\\ \underline{77.77}{\pm}1.99\end{array}$	$70.68{\pm}1.30\\71.56{\pm}1.41\\72.94{\pm}1.72\\74.02{\pm}1.87$	$71.42{\pm}1.38 \\ 73.15{\pm}1.37 \\ 74.36{\pm}1.32 \\ 75.21{\pm}1.19$	$74.27{\pm}1.08 \\ 76.15{\pm}0.77 \\ 77.17{\pm}1.21 \\ 76.76{\pm}1.99$	$\begin{array}{c} 68.92{\pm}1.06\\ 70.54{\pm}1.59\\ 71.76{\pm}1.55\\ 73.14{\pm}1.88\end{array}$	$70.51{\pm}1.12 \\72.01{\pm}1.32 \\73.41{\pm}1.44 \\73.56{\pm}1.96$	$\begin{array}{c} \textbf{75.69}{\pm}1.25\\ \textbf{78.16}{\pm}1.17^{*}\\ \textbf{78.61}{\pm}1.50^{*}\\ \textbf{79.23}{\pm}2.01^{*} \end{array}$
	Micro-F1	20% 40% 60% 80%	$\begin{array}{c} 77.19 {\pm} 0.88 \\ 79.18 {\pm} 1.24 \\ 80.03 {\pm} 1.51 \\ 80.13 {\pm} 2.25 \end{array}$	$\begin{array}{c} 79.12{\pm}0.93\\ 79.92{\pm}0.76\\ 80.64{\pm}0.84\\ 80.82{\pm}1.43\end{array}$	$73.78 {\pm} 0.94 \\76.63 {\pm} 0.65 \\77.67 {\pm} 0.73 \\78.55 {\pm} 1.55$	$\begin{array}{c} \underline{79.28}{\pm}0.76\\ \underline{80.57}{\pm}0.97\\ \underline{81.36}{\pm}1.22\\ \underline{81.51}{\pm}2.06 \end{array}$	$\begin{array}{c} 76.10 {\pm} 0.83 \\ 77.07 {\pm} 1.09 \\ 77.95 {\pm} 1.54 \\ 79.06 {\pm} 1.55 \end{array}$	$75.64{\pm}0.98 \\ 76.89{\pm}1.35 \\ 78.16{\pm}1.42 \\ 79.18{\pm}1.21$	$78.75 {\pm} 0.77 \\80.14 {\pm} 0.52 \\80.89 {\pm} 1.15 \\80.75 {\pm} 1.92$	$74.48 {\pm} 0.94 \\75.98 {\pm} 1.00 \\77.22 {\pm} 1.25 \\78.39 {\pm} 1.63$	$75.42{\pm}0.68 \\ 76.88{\pm}1.02 \\ 78.19{\pm}1.16 \\ 78.45{\pm}1.69$	79.62 ±0.89 80.92 ±0.96* 82.07 ±1.20* 83.31 ±1.91*
АСМ	Macro-F1	20% 40% 60% 80%	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 80.34{\pm}0.74\\ 81.14{\pm}0.81\\ 81.69{\pm}0.76\\ 81.49{\pm}1.56\end{array}$	$\begin{array}{c} 80.62{\pm}0.52\\ 81.20{\pm}0.78\\ 81.48{\pm}0.94\\ 81.57{\pm}1.34\end{array}$	$78.94{\pm}0.69\\80.41{\pm}0.51\\80.76{\pm}0.87\\80.74{\pm}1.08$	$\begin{array}{c} 81.45{\pm}0.81\\ 82.79{\pm}0.49\\ 83.07{\pm}1.02\\ 83.66{\pm}1.16\end{array}$	$\begin{array}{c} \underline{82.94}{\pm 0.38}\\ \underline{83.66}{\pm 0.55}\\ \underline{84.04}{\pm 0.68}\\ \underline{84.60}{\pm 1.62}\end{array}$	$\begin{array}{c} 82.09{\pm}0.62\\ 83.00{\pm}0.68\\ 82.61{\pm}0.73\\ 82.91{\pm}1.04\end{array}$	$79.97{\pm}0.60\\81.03{\pm}0.74\\81.70{\pm}0.82\\82.10{\pm}1.18$	$\begin{array}{c} 80.36{\pm}0.63\\ 80.93{\pm}0.60\\ 81.02{\pm}0.80\\ 81.76{\pm}1.18\end{array}$	$\begin{array}{c} \textbf{84.14}{\pm}0.67^{*} \\ \textbf{85.94}{\pm}0.48^{*} \\ \textbf{86.45}{\pm}0.59^{*} \\ \textbf{87.03}{\pm}1.52^{*} \end{array}$
	Micro-F1	20% 40% 60% 80%	$ \begin{vmatrix} 80.05 \pm 0.62 \\ 80.81 \pm 0.63 \\ 81.29 \pm 0.93 \\ 81.60 \pm 1.14 \end{vmatrix} $	$\begin{array}{c} 80.81{\pm}0.70\\ 81.48{\pm}0.82\\ 82.01{\pm}0.85\\ 81.82{\pm}1.61\end{array}$	$\begin{array}{c} 81.07{\pm}0.55\\ 81.75{\pm}0.83\\ 82.01{\pm}1.01\\ 82.20{\pm}1.36\end{array}$	$79.93 {\pm} 0.59 \\81.20 {\pm} 0.49 \\81.43 {\pm} 0.93 \\81.39 {\pm} 1.14$	81.87±0.74 83.12±0.48 83.40±1.04 84.00±1.14	$\frac{\underline{83.19} \pm 0.41}{\underline{83.83} \pm 0.58} \\ \underline{\underline{84.20}} \pm 0.79 \\ \underline{\underline{84.79}} \pm 1.50 \\ \end{array}$	$\begin{array}{c} 82.56 {\pm} 0.55 \\ 83.38 {\pm} 0.64 \\ 82.95 {\pm} 0.72 \\ 83.26 {\pm} 0.98 \end{array}$	$\begin{array}{c} 80.67{\pm}0.67\\ 81.70{\pm}0.73\\ 82.29{\pm}0.83\\ 82.79{\pm}1.15\end{array}$	$\begin{array}{c} 81.19{\pm}0.69\\ 81.66{\pm}0.56\\ 81.73{\pm}0.79\\ 82.55{\pm}1.20\end{array}$	$\begin{array}{c} \textbf{84.41}{\pm}0.62^{*}\\ \textbf{86.10}{\pm}0.50^{*}\\ \textbf{86.59}{\pm}0.61^{*}\\ \textbf{87.20}{\pm}1.47^{*} \end{array}$
DBLP	Macro-F1	20% 40% 60% 80%	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 89.50 {\pm} 0.35 \\ 90.68 {\pm} 0.64 \\ 91.39 {\pm} 0.58 \\ 91.61 {\pm} 0.67 \end{array}$	$\begin{array}{c} 90.24{\pm}0.43\\ 91.01{\pm}0.59\\ 91.58{\pm}0.60\\ 91.80{\pm}0.76\end{array}$	$\begin{array}{c} 92.65{\pm}0.36\\ 93.19{\pm}0.45\\ 93.71{\pm}0.44\\ 93.57{\pm}0.75\end{array}$	$\begin{array}{c} 90.86{\pm}0.27\\ 91.66{\pm}0.48\\ 92.33{\pm}0.50\\ 92.70{\pm}0.37\end{array}$	$\begin{array}{c} 92.55{\pm}0.43\\ 92.80{\pm}0.56\\ 92.80{\pm}0.52\\ 92.63{\pm}0.75\end{array}$	$\begin{array}{c} 84.63{\pm}1.90\\ 88.95{\pm}1.63\\ 91.13{\pm}0.65\\ 92.23{\pm}0.78\end{array}$	$\begin{array}{c} 93.11{\pm}0.22\\ 93.34{\pm}0.34\\ 93.45{\pm}0.42\\ 93.41{\pm}0.73\end{array}$	$\begin{array}{r} \underline{93.42}{\pm 0.28}\\ \underline{93.67}{\pm 0.61}\\ \underline{93.76}{\pm 0.67}\\ \underline{93.64}{\pm 0.76}\end{array}$	$\begin{array}{c} \textbf{94.17}{\pm}0.23^{*} \\ \textbf{94.55}{\pm}0.24^{*} \\ \textbf{94.59}{\pm}0.46^{*} \\ \textbf{94.92}{\pm}0.75^{*} \end{array}$
	Micro-F1	20% 40% 60% 80%	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 90.00{\pm}0.32\\ 91.13{\pm}0.61\\ 91.82{\pm}0.49\\ 92.04{\pm}0.54\end{array}$	$\begin{array}{c} 90.72{\pm}0.41\\ 91.46{\pm}0.56\\ 91.99{\pm}0.55\\ 92.22{\pm}0.66\end{array}$	$\begin{array}{c} 93.15{\pm}0.32\\ 93.63{\pm}0.40\\ 94.12{\pm}0.35\\ 94.05{\pm}0.67\end{array}$	$\begin{array}{c} 91.46{\pm}0.27\\ 92.21{\pm}0.45\\ 92.87{\pm}0.41\\ 93.24{\pm}0.34\end{array}$	$\begin{array}{c} 93.11{\pm}0.38\\ 93.31{\pm}0.52\\ 93.31{\pm}0.49\\ 93.18{\pm}0.73\end{array}$	$\begin{array}{c} 87.58{\pm}1.28\\ 90.52{\pm}1.21\\ 92.16{\pm}0.56\\ 93.07{\pm}0.68\end{array}$	$\begin{array}{c} 93.56 {\pm} 0.19 \\ 93.77 {\pm} 0.30 \\ 93.88 {\pm} 0.34 \\ 93.90 {\pm} 0.61 \end{array}$	$\begin{array}{c} \underline{93.87}{\pm 0.27} \\ \underline{94.11}{\pm 0.56} \\ \underline{94.20}{\pm 0.57} \\ \underline{94.14}{\pm 0.69} \end{array}$	94.56±0.21* 94.90±0.21* 94.96±0.40* 95.30±0.65*
AMiner	Macro-F1	20% 40% 60% 80%	$\begin{array}{c} 57.46 {\pm} 0.60 \\ 58.15 {\pm} 0.33 \\ 58.23 {\pm} 0.35 \\ 58.16 {\pm} 0.97 \end{array}$	$\begin{array}{c} \underline{63.11}{\pm 0.35} \\ \underline{63.45}{\pm 0.30} \\ \underline{63.61}{\pm 0.46} \\ \underline{63.35}{\pm 1.03} \end{array}$	$\begin{array}{c} 53.78 {\pm} 0.31 \\ 54.39 {\pm} 0.45 \\ 54.43 {\pm} 0.50 \\ 54.39 {\pm} 0.96 \end{array}$	$\begin{array}{c} 58.41{\pm}0.33\\ 58.72{\pm}0.31\\ 58.90{\pm}0.48\\ 58.93{\pm}0.99\end{array}$	$\begin{array}{c} 59.33 {\pm} 0.36 \\ 60.18 {\pm} 0.43 \\ 60.33 {\pm} 0.41 \\ 60.53 {\pm} 0.86 \end{array}$	$\begin{array}{c} 61.63{\pm}0.39\\ 62.43{\pm}0.30\\ 62.70{\pm}0.30\\ 63.01{\pm}0.50\end{array}$	$\begin{array}{c} 62.78 {\pm} 0.31 \\ 63.16 {\pm} 0.31 \\ 63.19 {\pm} 0.41 \\ \underline{63.45} {\pm} 0.93 \end{array}$	$\begin{array}{c} 61.86{\pm}0.55\\ 62.34{\pm}0.41\\ 62.46{\pm}0.39\\ 62.68{\pm}0.64\end{array}$	$\begin{array}{c} 62.00{\pm}0.67\\ 62.72{\pm}0.39\\ 63.18{\pm}0.68\\ 62.94{\pm}1.39\end{array}$	$\begin{array}{c} 69.08 {\pm} 0.67^{*} \\ 71.46 {\pm} 0.65^{*} \\ 72.37 {\pm} 0.80^{*} \\ 72.61 {\pm} 0.79^{*} \end{array}$
	Micro-F1	20% 40% 60% 80%	$ \begin{vmatrix} 78.25 \pm 0.27 \\ 78.69 \pm 0.20 \\ 78.81 \pm 0.27 \\ 78.96 \pm 0.41 \end{vmatrix} $	$\frac{\underline{82.04}\pm0.16}{\underline{82.24}\pm0.17}\\ \underline{\underline{82.31}}\pm0.26\\ \underline{\underline{82.41}}\pm0.59$	$75.78 {\pm} 0.22 \\ 76.22 {\pm} 0.32 \\ 76.32 {\pm} 0.41 \\ 76.48 {\pm} 0.82$	$78.15 \pm 0.14 \\ 78.37 \pm 0.20 \\ 78.48 \pm 0.26 \\ 78.63 \pm 0.48$	$78.90 {\pm} 0.16 \\79.60 {\pm} 0.21 \\79.80 {\pm} 0.31 \\80.12 {\pm} 0.63$	$\begin{array}{c} 80.34{\pm}0.17\\ 81.13{\pm}0.19\\ 81.31{\pm}0.24\\ 81.74{\pm}0.47\end{array}$	$\begin{array}{c} 81.71{\pm}0.12\\ 82.01{\pm}0.13\\ 82.07{\pm}0.29\\ 82.30{\pm}0.75\end{array}$	$\begin{array}{c} 81.35{\pm}0.18\\ 81.56{\pm}0.18\\ 81.64{\pm}0.35\\ 81.95{\pm}0.52\end{array}$	$79.90{\pm}0.14\\80.29{\pm}0.17\\80.47{\pm}0.36\\80.64{\pm}0.59$	$\begin{array}{c} \textbf{83.30}{\pm}0.23^{*}\\ \textbf{84.87}{\pm}0.29^{*}\\ \textbf{85.44}{\pm}0.31^{*}\\ \textbf{85.75}{\pm}0.47^{*}\\ \end{array}$

TABLE 4: Quantitative results ($\% \pm \sigma$) on node clustering. (bold: best; underline: runner-up)

Datasets	Metrics	DeepWalk	LINE-1st	LINE-2nd	ESim	HERec	Mp2vec	HIN2Vec	HetGNN	HeGAN	HEAD
Yelp	NMI ARI	$\begin{array}{c c} 34.73 {\pm} 1.32 \\ 40.39 {\pm} 1.24 \end{array}$	$\substack{38.66 \pm 0.00 \\ 42.47 \pm 0.00}$	$\substack{38.66 \pm 0.00 \\ 42.47 \pm 0.00}$	$35.12 {\pm} 0.06$ $38.42 {\pm} 0.05$	$\frac{39.01}{42.67} \pm 0.00$	$\begin{array}{c} 36.66{\pm}0.43\\ 41.56{\pm}0.84\end{array}$	$36.14{\pm}0.30$ $39.36{\pm}0.68$	$\begin{array}{c} 37.76 {\pm} 0.00 \\ 41.69 {\pm} 0.00 \end{array}$	$\begin{array}{c c} 37.31 \pm 0.32 \\ \underline{43.31} \pm 0.73 \end{array}$	$\substack{\textbf{40.06} \pm 0.00 \\ \textbf{44.02} \pm 0.00}$
ACM	NMI ARI	$\begin{array}{c} 56.28{\pm}0.18\\ 43.22{\pm}0.09\end{array}$	$^{21.58\pm0.00}_{15.45\pm0.00}$	$\substack{40.79 \pm 0.04 \\ 38.66 \pm 0.03}$	$\substack{48.21 \pm 0.00 \\ 35.45 \pm 0.00}$	$\begin{array}{c} 47.26{\pm}0.00\\ 36.39{\pm}0.00\end{array}$	$\substack{48.37 \pm 0.00 \\ 35.60 \pm 0.00}$	$\begin{array}{c} 45.46{\pm}0.00\\ 34.30{\pm}0.00\end{array}$	$56.28 {\pm} 0.09 \\ 43.45 {\pm} 0.10$	$\begin{array}{c} \underline{56.42} \pm 0.00 \\ 43.84 \pm 0.00 \end{array}$	$\frac{\textbf{56.64}{\pm}0.03}{\underline{43.51}{\pm}0.02}$
DBLP	NMI ARI	$\begin{array}{c} 73.22{\pm}0.25 \\ 78.45{\pm}0.26 \end{array}$	$\begin{array}{c} 78.08 {\pm} 0.11 \\ 82.01 {\pm} 0.08 \end{array}$	$67.35 {\pm} 0.24$ $72.87 {\pm} 0.34$	$\begin{array}{c} 75.67 {\pm} 0.13 \\ 81.09 {\pm} 0.17 \end{array}$	71.75±0.10 77.72±0.08	73.60±0.10 77.73±0.07	$\begin{array}{c} 77.44{\pm}0.06\\ 82.26{\pm}0.06\end{array}$	$\begin{array}{c} 76.95 {\pm} 0.00 \\ 82.52 {\pm} 0.00 \end{array}$	$\frac{78.98}{\underline{83.84}} \pm 0.03$	81.39 ±0.00 85.92 ±0.00
AMiner	NMI ARI	$\begin{array}{c} \underline{21.17} \pm 2.10 \\ \underline{14.61} \pm 4.46 \end{array}$	$\substack{11.93 \pm 2.33 \\ 4.81 \pm 4.69}$	$\substack{15.61 \pm 0.02 \\ 2.26 \pm 0.02}$	$\substack{21.03 \pm 0.17 \\ 14.20 \pm 0.23}$	$\substack{10.88 \pm 2.27 \\ 11.42 \pm 1.89}$	$\substack{18.56 \pm 1.72 \\ 2.70 \pm 1.90}$	$\substack{18.58 \pm 0.66 \\ -2.73 \pm 2.37}$	$\substack{19.66 \pm 0.37 \\ 3.41 \pm 2.32}$	20.75±0.17 12.17±0.14	21.42 ±1.39 14.63 ±2.41

where authors are divided into four classes based on the conferences they submitted to. The meta-path set we utilize is {*APA*, *APCPA*, *APTPA*}.

• AMiner [58]. This dataset comprises three types of nodes, *i.e.*, author (A), paper (P) and reference (R), where papers are divided into six classes *w.r.t.* the conferences they published on. The used meta-paths are {*PAP*, *PRP*}.

Please note that for the AMiner dataset, we do not add conference nodes to the heterogeneous information network. Because once the PCP meta-path is used, the classification performances of the meta-path based HIN embedding methods *w.r.t. F1*, all approach 1, losing the meaning of comparison.

5.1.2 Baselines

We consider three categories of methods: homogeneous network embedding (DeepWalk, LINE), meta-path based (ESim, HERec, Mp2vec, HIN2Vec) and meta-path free (Het-GNN, HeGAN) HIN embedding algorithms, summarized as follows.

• **DeepWalk** [23] performs truncated random walks, and employs the skip-gram model [20]. Here we ignore the heterogeneity of nodes and evaluate it on the whole HINs.

- LINE [25] exploits the first (LINE-1st) and second (LINE-2nd) order proximity in networks and preserves the local and global network structures meanwhile.
- ESim [7] captures semantic information from multiple meta-paths, and assigns weights to them. Here we conduct grid search to find the optimal weights for these meta-paths.
- **HERec** [16] designs a type-constrained strategy to filter node sequences of meta-path and utilizes skip-gram [20] to embed the heterogeneous information networks.
- **Mp2vec** [8] samples meta-path based random walks and employ skip-gram model [20] for semantic-preserving HIN embedding.
- **HIN2Vec** [41] converts HIN embedding into multi-task classification, then learns latent representations for both nodes and meta-paths jointly.
- **HetGNN** [44] utilizes random walk with restart to sample a fixed size of neighbors and group them according to their types, finally learns node embedding based on GNN framework.
- HeGAN [43] proposes a relation-aware discriminator and generator to embed HINs on the adversarial principle to

learn node distribution for better negative sampling.

Please note that we do not choose HAN [18] as baseline, because HAN is semi-supervised, while HEAD and other baselines are unsupervised methods, thus cannot be directly compared.

5.1.3 Implementation

We implement the proposed model HEAD with deep learning library PyTorch¹. All experiments are conducted on a Linux server with GPU (GeForce GTX 1080 Ti) and CPU (Intel Xeon E5-2620), and its operating system is Ubuntu 16.04.1. The Python and PyTorch versions are 3.6.9 and 1.3.1, respectively.

We perform Adaptive Moment Estimation (Adam) [59] to optimize our model and apply a grid search for hyper-parameters: the learning rate is tuned in $\{0.05, 0.01, 0.005, 0.001\}$ and the dropout ratio is tuned amongst $\{0.1, 0.2, \dots, 0.5\}$. Moreover, we utilize the two-layer MLP for all modules. To avoid gradients vanishing or exploding, we employ Batch Normalization [60] and set LeakyReLU [61] as the activation function. During each iteration of alternative training, we use a batch size of 32 and set the prior distribution to standard Gaussian. For the input of HEAD, we use node embeddings from Mp2vec for each meta-path.

The embedding dimension for all methods is set as 64. For random walk-based methods (*i.e.*, DeepWalk, HERec, HIN2Vec, Mp2vec and HetGNN), we set the number of walks per node to 10, the walk length to 100 and the window size to 5. For meta-path based methods, we select the same meta-paths as our model. Particularly, for ESim, we search the weights for all meta-paths; while for HERec and Mp2vec, we concatenate the embeddings of all meta-paths to ensure they utilize the same global information as ours. For LINE, we set the number of samples as 10000M. For fair comparison, HetGNN and HeGAN also use pre-trained embeddings from Mp2vec, which serve as node attributes for HetGNN meanwhile. In terms of other parameters, we follow the settings in their original papers.

5.2 Experimental Evaluation

5.2.1 Node Classification

In this task, we use different percentages of labeled nodes (20%, 40%, 60%, 80%) to train a logistic regression classifier whose input is the node embeddings learned by different models, and test the classifier on the remaining nodes. We evaluate the classification quality in terms of *Micro-F1* and *Macro-F1 w.r.t.* the node labels. Since the performance of logistic regression is affected by train/test split, we repeat 10 times, and report average and standard deviation results in Table 3. We use * to indicate that HEAD is significantly different from runner-up methods based on paired t-tests at the significance level of 0.01.

Based on results, we make the following observations.

• HEAD consistently and significantly outperforms all baselines and achieves performance gains over the best baseline by 2.2%, 2.9%, 1.4% and 14.4% on four datasets respectively, which demonstrates that disentangled intrinsic embeddings can characterize nodes in a robust manner. Note that the overwhelming performance superiority of HEAD over Mp2vec serving as the input of HEAD implies that HEAD is capable of enhancing existing meta-path based methods.

- Compared with the runner-ups, HEAD greatly improves the classification performance on the AMiner dataset. One feasible reason is that the used meta-paths, *i.e.*, PAP and PRP, are not informative. In this case, using the information of all meta-paths and purifying out intrinsic factors may reduce noise and learn robust representation.
- As the percentage of labeled nodes increases, our performance margins over the runner-ups become larger, since disentangled factors benefit more from the supervision, widening the gaps.
- Both HEAD and HeGAN employ the adversarial principle, thus their performance improvement illustrates the vital role of adversarial learning in enhancing representation robustness. Moreover, HEAD is consistently better than HeGAN, further indicating the effectiveness of disentangling the intrinsic factors.
- Among the baselines, HIN methods generally outperform homogeneous network methods (*i.e.*, DeepWalk and LINE), implying the importance of semantic-preserving embedding.

5.2.2 Node Clustering

In this task, we employ the K-Means algorithm on the embeddings produced by all algorithms to perform clustering, and evaluate the clustering quality in terms of *normalized mutual information* (*NMI*) and *adjusted rand index* (*ARI*) *w.r.t.* the node labels. Note that the larger *NMI* and *ARI* mean the better performance. Since the performance of K-Means is affected by initial centroids, we repeat the process for 10 times, and report average and standard deviation results in Table 4.

Similar conclusions to the node classification can be drawn, where HEAD generally shows better performance, which once again proves the effectiveness of intrinsic representations.

5.2.3 Relation Prediction

In this task, we predict the *co-author* relation of authors (*i.e.*, APA) on DBLP, and the *same-author* relation of papers (*i.e.*, PAP) on ACM and AMiner. Specifically, we randomly hide 20% of author-paper relations from the original network and label their corresponding relations as ground-truth positives, then randomly sample the same number of negative instances to form the test set together. Subsequently, we utilize the left network to learn node embeddings. After obtaining node embeddings, we train a logistic regression classifier using node embeddings of positive and negative instances from the left network. We adopt *F1*, *ACC* and *AUC* as evaluation metrics, shown in Table 5.

To investigate the impact of specific factors on the corresponding relation prediction, we evaluate the performances of the final intrinsic embedding and the concatenation of the final intrinsic embedding and corresponding meta-path specific embedding (*i.e.*, tested meta-paths: APA for DBLP, PAP for ACM and AMiner), denoted as HEAD and HEAD_s.

TABLE 5: Quantitative results (%) on relation prediction. (bold: best; underline: runner-up)

Datasets	Metrics	DeepWalk	LINE-1st	LINE-2nd	ESim	HERec	Mp2vec	HIN2Vec	HetGNN	HeGAN HEAD	$HEAD_s$
ACM (Same-Author)	F1 ACC AUC	58.36 59.32 63.39	64.88 64.64 66.34	61.37 62.72 66.08	59.35 61.45 64.91	64.96 63.73 65.74	63.92 63.77 67.24	60.84 62.68 67.01	59.40 61.43 65.57	$\begin{array}{c cccc} 59.80 & \underline{65.60} \\ 61.00 & \underline{65.04} \\ 65.26 & \underline{67.56} \end{array}$	65.92 65.74 69.15
DBLP (Co-Author)	F1 ACC AUC	56.59 56.57 59.66	56.36 57.12 61.00	54.41 58.27 61.36	55.98 57.67 60.66	57.31 57.55 60.54	<u>58.32</u> 57.60 59.69	50.52 57.91 <u>61.87</u>	56.29 56.97 60.50	$\begin{array}{c cccc} 55.32 & 57.39 \\ 56.90 & \underline{58.28} \\ 60.75 & 61.03 \end{array}$	59.19 59.43 62.00
AMiner (Same-Author)	F1 ACC AUC	62.19 62.80 65.57	59.11 59.49 63.31	56.57 58.04 61.13	58.09 59.62 64.19	62.39 61.73 62.03	60.54 62.73 65.92	62.06 61.91 <u>66.84</u>	62.38 61.91 65.78	62.7961.4861.0962.2563.7964.46	<u>62.47</u> 62.96 66.96



Fig. 3: Visualization of the embeddings learned by baselines and HEAD on DBLP. Color indicates the categories of authors.



Fig. 4: Disentanglement analysis of HEAD on DBLP.

We observe the following phenomena. (1) HEAD_s and HEAD generally outperform all baselines, which reveals that detaching intrinsic and specific factors innately learns a better structure- and semantics-preserving embedding space for relations. (2) There exists performance gain of HEAD_s over HEAD. Note that the task is to predict relations defined by given meta-paths, which is directly meta-path dependent. HEADs injects meta-path specific embeddings, which makes it perform better than HEAD. It implies that the learned specific embeddings are also useful for metapath related tasks.

5.2.4 Network Visualization

To examine the network representation intuitively, we visualize embeddings of author nodes in DBLP using the t-SNE algorithm [62] in Figure 3.

From the plots, we find that DeepWalk cannot effectively identify different author categories due to the ignorance of heterogeneity. On the other hand, Mp2vec and HeGAN can reasonably separate author categories, implying that they learn semantic-preserving embeddings to some degree. Compared to these baselines, HEAD correctly separate these authors with crisp boundaries, demonstrating that disentangling intrinsic factors from multiple meta-paths can exert a vital part in reducing noise.

5.3 Model Analysis

5.3.1 Disentanglement Analysis

First of all, we study whether HEAD can effectively disentangle intrinsic and specific factors of nodes. We use DBLP dataset as an example, and similar patterns have been observed on other datasets. In Figure 4, we visualize intrinsic and specific embedding under APCPA, subsequently present the performance changes of the final intrinsic embeddings and specific embeddings under all meta-paths on node clustering task.

In Figure 4 (a) and (b), the visualizations for intrinsic and specific embeddings are labeled by author categories which can be considered as intrinsic factors, and we find that the intrinsic embedding correctly separates nodes with clear boundaries, whereas the specific embedding is an extremely chaotic. We further label specific embeddings by meta-path types which can be treated as specific factors, and re-visualize them. As shown in Figure 4 (c), the visualization results have obvious clusters, conforming to our hypothesis that specific embedding reflects meta-path dependent semantics. From Figure 4 (d), we observe that the final intrinsic embedding converges quickly and performs well, while the performance of the specific embeddings fluctuates and slowly increases at a low level. The results are in





(c) Loss change of meta-path and



(a) Performances of HEAD and variants on node classification



(b) Performances of HEAD and variants on node clustering

Fig. 5: Analysis on adversarial learning of HEAD.

semantic classifier



Fig. 6: Performances of HEAD gradually adding paths.



Fig. 7: Improvements to existing meta-path based methods.

line with our assumption that the intrinsic factors involved in various meta-paths reveal the intrinsic characteristic of nodes.

5.3.2 Adversarial Learning

Since adversarial learners in HEAD guarantee representation disentanglement, we conduct an ablation study to validate the effectiveness of them. In particular, we compare HEAD with three variants named HEAD_{\mp} (*i.e.*, HEAD without meta-path discriminator), HEAD_{\sem} (*i.e.*, without semantic discriminator) and HEAD_{\mbox} (*i.e.*, without both discriminators). For HEAD and its variants, we separately evaluate intrinsic embeddings under every meta-path used in citation networks (*i.e.*, ACM, DBLP and AMiner datasets) on node clustering and classification tasks, and plot their performances of different meta-paths as boxes in Figure 5 (a) and (b).

With respect to the average performance, we observe that HEAD is consistently better than the remaining variants and HEAD_{\m&s} generally performs the worst. Such phenomenon is not surprising, and reveals that disregarding the adversarial learners leads to the intrinsic factors mottled, which harms the performance. In terms of the fluctuation range, HEAD is always minimal. We hypothesize that removing any module affects the mutual guidance of semantic purification among meta-paths, pushing the performances of various meta-paths away. Overall, these observations justify the necessity of adversarial learners.

sifier

Next, to understand the adversarial learning process of these two discriminators, we present the loss change of classifiers inside them on DBLP. As shown in Figure 5 (c), the loss of semantic classifier decreases, in the sense that the specific embedding becomes more meta-path dependent. Whereas the loss of meta-path classifier increases, the reason for which is the encoder tries to fool meta-path classifier by disentangling more intrinsic embedding. From Figure 5 (d), we observe the classification performances of real and fake semantics gradually become worse, indicating that the synthesized semantic is getting closer to the real one. After around 10k steps of adversarial training, the losses tend to converge. Note that the step here represents the mini-batch. In all, the experiment results reveal that these adversarial learners can stably and consistently improve performance through mutual adversarial learning.

5.3.3 Effects of Meta-paths

To investigate the effect of different meta-paths, we inject meta-paths to HEAD one by one from scratch, constantly concatenate intrinsic embedding under the new added meta-path and evaluate the resulted intrinsic embedding on node classification. We report their performances in Figure 6.

With the increase of injected meta-paths, the performance of HEAD consistently grows, which further confirms that HEAD can absorb intrinsic factors from diverse metapaths for better performance. Note that, although many HIN based methods [16], [63] can improve performance when adding more meta-paths, they need supervised information. But the experiments here imply our HEAD is able to robustly improve performance through adding meta-paths, even though we have no supervised information.



Fig. 9: Impact of standard deviation on Gaussian prior.

5.3.4 Improvements to existing methods

As mentioned before, our HEAD can purify the learned embeddings of different meta-paths, which indicates that HEAD can effectively absorbs intrinsic factors from various meta-paths and enhance the existing meta-path based methods. Therefore, we design comparative experiments for different input methods to prove that point.

Here we employ two typical models (Mp2vec and HERec) as input methods. Specifically, for each dataset, we train the proposed model HEAD for 100 epochs with Mp2vec and HERec as input methods, respectively. Then the input embedding and disentangled intrinsic embedding are respectively evaluated by node classification and clustering task *w.r.t. Micro-F1* and *NMI* metrics. Please note that the input embedding and intrinsic embedding evaluated here are the concatenation embedding of all meta-paths in the dataset. As shown in Figure 7, the blue bars correspond to the *Micro-F1* on the left Y axis and the green bars correspond to the *NMI* on the right Y axis.

It is obvious that HEAD can significantly improve the performance of two base models, Mp2vec and HERec, on node classification and clustering tasks on all datasets. This demonstrates that HEAD, as an unsupervised learner, can utilize the mutual guidance among meta-paths to extract the intrinsic factors of nodes. Moreover, the performance improvements on node classification seem more significant than clustering, which further confirms that the disentangled factors benefit more from the supervision.

5.4 Hyper-parameter Sensitivity

Finally, we investigate the sensitivity of hyper-parameters and report the results *w.r.t.* node classification on four datasets. In the experiments, we vary the studied hyperparameter and fix the others.

5.4.1 Dimension Selection

The embedding dimension is a key parameter to control the complexity and capacity. Therefore, we evaluate how it affects the classification performance. As shown in Figure 8, as we gradually increase the embedding dimension d, the performance grows since a larger dimension could enhance the representation capability. Nevertheless, when d is larger than the optimal value, increasing d will hurt the performance probably due to overfitting. Therefore, we employ the applicable embedding dimension 64 to balance the trade-off between performance and complexity.

5.4.2 Standard Deviation

The standard deviation of Gaussian prior reflects the distribution law that the intrinsic factors may follow. A small standard deviation σ will make the meta-path disentangler closer to vanilla auto-encoder, while a large σ gives the intrinsic factors more uncertainty. To investigate whether HEAD can benefit from modeling the uncertainty, we vary the standard deviation σ of Gaussian prior in the range of $\{0.001, 0.01, 0.1, 1, 10\}$, while keeping the other parameters the same. Figure 9 shows the experimental results on four datasets. HEAD achieves optimal performance at the standard deviation $\sigma = 1.0$ and is generally stable around it, and too small or large values would harm the model.

6 CONCLUSION

In this paper, we first investigate latent explanatory factors of nodes in HIN embedding and propose a novel solution HEAD with disentangled representation. Our HEAD purifies out intrinsic embedding from input embeddings of multiple meta-paths, which absorb informative characteristic from all meta-paths and reduce the noise interference. Specifically, given node embeddings of various metapaths, a meta-path disentangler in HEAD encodes them into intrinsic and specific spaces. Subsequently, two adversarial learners (*i.e.*, meta-path discriminator and semantic discriminator) utilize meta-path schemes as self-supervised information to further refine the disentanglement. Extensive experimental results verify the effectiveness of HEAD and its ability to attain robust node embedding.

7 ACKNOWLEDGMENTS

This research is supported in part by the National Natural Science Foundation of China (No. U20B2045, 61772082, 62002029), The Fundamental Research Funds for the Central Universities (2021RC28), and BUPT Excellent Ph.D. Students Foundation (CX2021105). In addition, Y. Ye's work is partially supported by the NSF under grants IIS-2107172, IIS-2027127, IIS-2040144, CNS-2034470, IIS-1951504, CNS-1940859, CNS-1814825 and OAC-1940855, and the DoJ/NIJ 2018-75-CX-0032.

REFERENCES

- Y. Sun and J. Han, "Mining heterogeneous information networks: a structural analysis approach," Acm Sigkdd Explorations Newsletter, pp. 20–28, 2013.
- [2] C. Shi, Y. Li, J. Zhang, Y. Sun, and S. Y. Philip, "A survey of heterogeneous information network analysis," *IEEE Transactions* on Knowledge and Data Engineering, pp. 17–37, 2016.
- [3] C. Shi and S. Y. Philip, *Heterogeneous information network analysis and applications*. Springer, 2017.
- [4] H. Cai, V. W. Zheng, and K. C.-C. Chang, "A comprehensive survey of graph embedding: Problems, techniques, and applications," *IEEE Transactions on Knowledge and Data Engineering*, pp. 1616–1637, 2018.
- [5] P. Cui, X. Wang, J. Pei, and W. Zhu, "A survey on network embedding," *IEEE Transactions on Knowledge and Data Engineering*, pp. 833–852, 2018.
- [6] J. Tang, M. Qu, and Q. Mei, "Pte: Predictive text embedding through large-scale heterogeneous text networks," in *SIGKDD*, 2015, pp. 1165–1174.
- [7] J. Shang, M. Qu, J. Liu, L. M. Kaplan, J. Han, and J. Peng, "Metapath guided embedding for similarity search in large-scale heterogeneous information networks," arXiv preprint arXiv:1610.09769, 2016.
- [8] Y. Dong, N. V. Chawla, and A. Swami, "metapath2vec: Scalable representation learning for heterogeneous networks," in *SIGKDD*, 2017, pp. 135–144.
- [9] L. Xu, X. Wei, J. Cao, and P. S. Yu, "Embedding of embedding (eoe): Joint embedding for coupled heterogeneous networks," in WSDM, 2017, pp. 741–749.
- [10] X. Wang, P. Cui, J. Wang, J. Pei, W. Zhu, and S. Yang, "Community preserving network embedding." in AAAI, 2017, pp. 203–209.
- [11] C. Yang, Y. Xiao, Y. Zhang, Y. Sun, and J. Han, "Heterogeneous network representation learning: Survey, benchmark, evaluation, and beyond," arXiv preprint arXiv:2004.00216, 2020.
- [12] Y. Zhang, Y. Xiong, X. Kong, S. Li, J. Mi, and Y. Zhu, "Deep collective classification in heterogeneous information networks," in WWW, 2018, pp. 399–408.
- [13] Y. Jacob, L. Denoyer, and P. Gallinari, "Learning latent representations of nodes for classifying in heterogeneous social networks," in WSDM, 2014, pp. 373–382.
- [14] Y. Sun, R. Barber, M. Gupta, C. C. Aggarwal, and J. Han, "Coauthor relationship prediction in heterogeneous bibliographic networks," in ASONAM, 2011, pp. 121–128.
- [15] H. Wang, F. Zhang, M. Hou, X. Xie, M. Guo, and Q. Liu, "Shine: Signed heterogeneous information network embedding for sentiment link prediction," in WSDM, 2018, pp. 592–600.
 [16] C. Shi, B. Hu, W. X. Zhao, and S. Y. Philip, "Heterogeneous infor-
- [16] C. Shi, B. Hu, W. X. Zhao, and S. Y. Philip, "Heterogeneous information network embedding for recommendation," *IEEE Transactions on Knowledge and Data Engineering*, pp. 357–370, 2018.
 [17] B. Hu, C. Shi, W. X. Zhao, and P. S. Yu, "Leveraging meta-
- [17] B. Hu, C. Shi, W. X. Zhao, and P. S. Yu, "Leveraging metapath based context for top-n recommendation with a neural coattention model," in *SIGKDD*, 2018, pp. 1531–1540.

- [18] X. Wang, H. Ji, C. Shi, B. Wang, Y. Ye, P. Cui, and P. S. Yu, "Heterogeneous graph attention network," in WWW, 2019, pp. 2022–2032.
- [19] Y. Sun, J. Han, X. Yan, P. S. Yu, and T. Wu, "Pathsim: Meta path-based top-k similarity search in heterogeneous information networks," *Proceedings of the VLDB Endowment*, pp. 992–1003, 2011.
- [20] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," arXiv preprint arXiv:1301.3781, 2013.
- [21] D. Zhang, J. Yin, X. Zhu, and C. Zhang, "Network representation learning: A survey," *IEEE transactions on Big Data*, 2018.
- [22] D. Zhou, J. He, H. Yang, and W. Fan, "Sparc: Self-paced network representation for few-shot rare category characterization," in *KDD*, 2018, pp. 2807–2816.
- [23] B. Perozzi, R. Al-Rfou, and S. Skiena, "Deepwalk: Online learning of social representations," in SIGKDD, 2014, pp. 701–710.
- [24] A. Grover and J. Leskovec, "node2vec: Scalable feature learning for networks," in SIGKDD, 2016, pp. 855–864.
- [25] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, and Q. Mei, "Line: Large-scale information network embedding," in WWW, 2015, pp. 1067–1077.
- [26] D. Wang, P. Cui, and W. Zhu, "Structural deep network embedding," in SIGKDD, 2016, pp. 1225–1234.
- [27] S. Cao, W. Lu, and Q. Xu, "Grarep: Learning graph representations with global structural information," in CIKM, 2015, pp. 891–900.
- [28] M. Ou, P. Cui, J. Pei, Z. Zhang, and W. Zhu, "Asymmetric transitivity preserving graph embedding," in SIGKDD, 2016, pp. 1105– 1114.
- [29] C. Yang, Z. Liu, D. Zhao, M. Sun, and E. Chang, "Network representation learning with rich text information," in *IJCAI*, 2015.
- [30] D. Zhang, J. Yin, X. Zhu, and C. Zhang, "Homophily, structure, and content augmented network representation learning," in *ICDM*, 2016, pp. 609–618.
- [31] D. Zhou, L. Zheng, J. Han, and J. He, "A data-driven graph generative model for temporal interaction networks," in *KDD*, 2020, pp. 401–411.
- [32] Z. Liu, D. Zhou, Y. Zhu, J. Gu, and J. He, "Towards finegrained temporal network representation via time-reinforced random walk," in AAAI, vol. 34, no. 04, 2020, pp. 4973–4980.
- [33] Z. Liu, D. Zhou, and J. He, "Towards explainable representation of time-evolving graphs via spatial-temporal graph attention networks," in *CIKM*, 2019, pp. 2137–2140.
- [34] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in NIPS, 2016, pp. 3844–3852.
- [35] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *ICLR*, 2017.
- [36] W. L. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in NIPS, 2017, pp. 1025–1035.
- [37] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *ICLR*, 2014.
- [38] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in NIPS, 2017, pp. 5998–6008.
- [39] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *ICLR*, 2018.
- [40] B. Hu, Z. Zhang, C. Shi, J. Zhou, X. Li, and Y. Qi, "Cash-out user detection based on attributed heterogeneous information network with a hierarchical attention mechanism," in AAAI, 2019, pp. 946– 953.
- [41] T.-y. Fu, W.-C. Lee, and Z. Lei, "Hin2vec: Explore meta-paths in heterogeneous information networks for representation learning," in CIKM, 2017, pp. 1797–1806.
- [42] H. Chen, H. Yin, W. Wang, H. Wang, Q. V. H. Nguyen, and X. Li, "Pme: projected metric embedding on heterogeneous networks for link prediction," in *SIGKDD*, 2018, pp. 1177–1186.
- [43] B. Hu, Y. Fang, and C. Shi, "Adversarial learning on heterogeneous information networks," in SIGKDD, 2019, pp. 120–129.
- [44] C. Zhang, D. Song, C. Huang, A. Swami, and N. V. Chawla, "Heterogeneous graph neural network," in *SIGKDD*, 2019, pp. 793–803.
- [45] X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, and P. Abbeel, "Infogan: Interpretable representation learning by information maximizing generative adversarial nets," in *NIPS*, 2016, pp. 2172–2180.

- [46] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, "beta-vae: Learning basic visual concepts with a constrained variational framework," *ICLR*, 2017.
- [47] A. Odena, C. Olah, and J. Shlens, "Conditional image synthesis with auxiliary classifier gans," in *ICML*, 2017, pp. 2642–2651.
- [48] H. Kim and Á. Mnih, "Disentangling by factorising," *arXiv preprint arXiv:1802.05983*, 2018.
- [49] T. D. Kulkarni, W. F. Whitney, P. Kohli, and J. Tenenbaum, "Deep convolutional inverse graphics network," in NIPS, 2015, pp. 2539– 2547.
- [50] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE transactions on pattern analysis and machine intelligence*, pp. 1798–1828, 2013.
- [51] A. A. Alemi, I. Fischer, J. V. Dillon, and K. Murphy, "Deep variational information bottleneck," arXiv preprint arXiv:1612.00410, 2016.
- [52] J. Ma, P. Cui, K. Kuang, X. Wang, and W. Zhu, "Disentangled graph convolutional networks," in *ICML*, 2019, pp. 4212–4221.
- [53] Y. Liu, X. Wang, S. Wu, and Z. Xiao, "Independence promoted graph disentangled networks," arXiv preprint arXiv:1911.11430, 2019.
- [54] D. J. Rezende, S. Mohamed, and D. Wierstra, "Stochastic backpropagation and approximate inference in deep generative models," arXiv preprint arXiv:1401.4082, 2014.
- [55] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," arXiv preprint arXiv:1312.6114, 2013.
- [56] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee, "Generative adversarial text to image synthesis," arXiv preprint arXiv:1605.05396, 2016.
- [57] D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling, "Semisupervised learning with deep generative models," in *NIPS*, 2014, pp. 3581–3589.
- [58] R. Bekkerman and A. McCallum, "Disambiguating web appearances of people in a social network," in WWW, 2005, pp. 463–470.
- [59] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," arXiv preprint arXiv:1412.6980, 2014.
- [60] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," arXiv preprint arXiv:1502.03167, 2015.
- [61] A. L. Maas, A. Y. Hannun, and A. Y. Ng, "Rectifier nonlinearities improve neural network acoustic models," in *ICML*, 2013.
- [62] L. v. d. Maaten and G. Hinton, "Visualizing data using t-sne," Journal of machine learning research, pp. 2579–2605, 2008.
 [63] C. Shi, Z. Zhang, P. Luo, P. S. Yu, Y. Yue, and B. Wu, "Semantic path
- [63] C. Shi, Z. Zhang, P. Luo, P. S. Yu, Y. Yue, and B. Wu, "Semantic path based personalized recommendation on weighted heterogeneous information networks," in *CIKM*, 2015, pp. 453–462.



Tianyu Zhao received the B.S. degree from Beijing University of Posts and Telecommunications in 2020. He is currently a master student in Beijing University of Posts and Communications. His current research interest is graph analysis, especially heterogeneous graph neural networks.



Xiao Wang is an Assistant Professor in the School of Computer Science, Beijing University of Posts and Telecommunications. He received his Ph.D. degree from the School of Computer Science and Technology, Tianjin University, in 2016. He was a postdoctoral researcher in Department of Computer Science and Technology, Tsinghua University. He got the China Scholarship Council Fellowship in 2014 and visited Washington University as a joint training student from 2014 to 2015. His current research inter-

ests include data mining, social network analysis and machine learning. Until now, he has published more than 50 papers in conferences, such as AAAI, IJCAI, KDD, and journals such as IEEE TKDE, IEEE Trans. on Cybernetics, etc.



Yanfang Ye is Theodore and Dana Schroeder Associate Professor at Case Western Reserve University. She received her Ph.D. in computer science from Xiamen University. She was an assistant professor and then associate professor in the department of computer science and electrical engineering (CSEE) at West Virginia University. Her research areas mainly include data mining, machine learning, cybersecurity, and health intelligence.



Ruijia Wang received the B.S. degree from Beijing University of Posts and Telecommunications in 2019. She is currently a Ph.D. student in Beijing University of Posts and Communications. Her current research interest is graph analysis, especially graph neural networks.



Chuan Shi received the B.S. degree from the Jilin University in 2001, the M.S. degree from the Wuhan University in 2004, and Ph.D. degree from the ICT of Chinese Academic of Sciences in 2007. He joined the Beijing University of Posts and Telecommunications as a lecturer in 2007, and is a professor and deputy director of Beijing Key Lab of Intelligent Telecommunications Software and Multimedia at present. His research interests are in data mining and machine learning. He has published more than 100 papers in

refereed journals and conferences.