

## Chapter 6

# Emerging Topics of Heterogeneous Graph Representation

**Abstract** Heterogeneous graph (HG) embedding, aiming to project HG into a low-dimensional space, has attracted considerable research attention. We have introduced some kinds of HG embedding methods, and there are also some other essential topics on HG embeddings. In this chapter, we will introduce three novel HG embedding methods. Specifically, to learn semantic-preserving and robust node representations, we study the problem of adversarial learning on HG. Also, we work with large-scale heterogeneous interaction graphs and focus on the problem of importance sampling on HG embedding. Moreover, we explore the intrinsic spaces of HG and propose a hyperbolic space based HG embedding method.

### 6.1 Introduction

Recently, some efforts begin to combine graph embedding methods with varied machine learning technologies to obtain better graph representations. These machine learning technologies include generative adversarial networks (GANs) [12, 26], importance sampling [5, 17], as well as hyperbolic representation [20, 22]. Specifically, GANs hinge on the idea of adversarial learning, in which discriminators and generators compete with each other to learn better underlying data distributions. Besides, importance sampling aims to reduce the computational and memory costs in machine learning. Moreover, hyperbolic spaces are non-Euclidean spaces, which are more suitable to model data with a hierarchical structure.

Some homogeneous graph embedding methods combined with these technologies have shown their effectiveness in modeling graphs. For example, GAN-based homogeneous graph embedding methods [6, 29] leverage generative and discriminative components to learn graph representation. The generative component aims to learn the underlying connectivity distribution in the graph, while the discriminative component aims to predict the probability of edge existence between nodes. Besides, to make large-scale graph representation learning possible, a natural idea is to sample a small set of important nodes from the whole graph. The importance

sampling strategies have achieved significant improvements for homogeneous graph representation learning [5, 14, 17, 39]. Moreover, some researchers begin to embed graphs into low-dimensional hyperbolic spaces. They find that embedding a graph in hyperbolic spaces would have a low distortion when the graph has a hierarchical structure [11, 22], and the hierarchical relation between nodes can be reflected by analyzing hyperbolic embedding [22].

In this chapter, we introduce three emerging HG embedding methods. First, **HG** embedding with **GAN**-based adversarial learning (named **HeGAN** [16]) is designed to learn semantic-preserving and robust node representations, which leverages the relation-aware discriminator and generator to fit the heterogeneous setting. Second, to reduce the cost of large-scale HG embedding, **Heterogeneous importance Sampling** (named **HeteSamp** [19]) is proposed, which introduces both type-dependent and type-fusion samplers with self-normalized and adaptive estimators to ensure the model efficiency. Third, assuming the underlying spaces of HG are hyperbolic spaces, **Hyperbolic Heterogeneous Network Embedding** (named **HHNE**) [33] is proposed to preserve the structure and semantic information of HG in hyperbolic spaces.

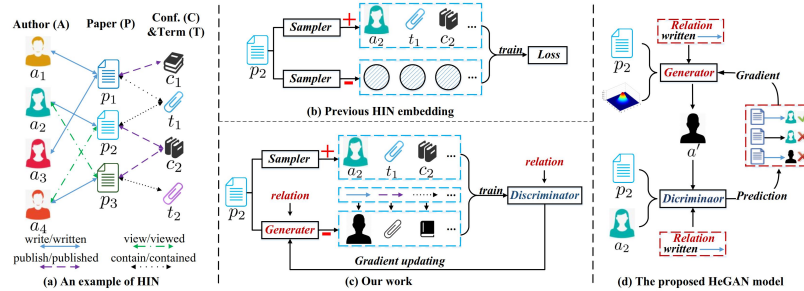
## 6.2 Adversarial Learning

### 6.2.1 Overview

GANs [12, 26] have been developed for learning robust latent representations in various applications [31, 36]. In recent researches, GANs only study homogeneous graphs [6, 24, 30, 37], such as the bibliographic data in Fig. 6.1a, so they do not consider the heterogeneity of nodes and edges, resulting in unsatisfactory performance on HG. In existing methods, positive (real) & negative (fake) nodes can only be differentiated by the network structure, which is shown in Fig. 6.1b. Therefore, it is urgent to design new formal discriminator and generator to distinguish and simulate the real and fake semantic-rich nodes involving various relationships. Generators are essentially picking an existing node from the original network according to the learned distribution, without the ability to generalize to “unseen” nodes. Not surprisingly, they do not generate the most representative fake nodes, as such nodes may not even appear in the network. Thus, it is important to design a generator that can efficiently produce latent fake samples.

We exploit the heterogeneity of HG in adversarial setting in order to learn semantic-preserving and robust node representations. Different from previous efforts [2, 6, 24, 30, 37], HeGAN raises two major novelties to address the challenges of adversarial learning on HG. It is not only relation-aware to capture rich semantics, but also equipped with a generalized generator that is effective and efficient. It proposes a new form of discriminator and generator, as illustrated in Fig. 6.1c. To a given relation, the discriminator can tell whether a node pair is real or fake, whereas the generator can produce fake node pairs that mimic real pairs. To further

improve the effectiveness and efficiency of sample generation, we propose a generalized generator, which is able to directly sample latent nodes from a continuous distribution.



**Fig. 6.1** Overview of HeGAN. (a) A toy example of HG for bibliographic data. (b) An example of previous HG embedding. (c) Relation-aware discriminator and generator in our work. (d) The framework of HeGAN for adversarial learning on HGs.

## 6.2.2 The HeGAN Model

### 6.2.2.1 Model Framework

HeGAN is based on GAN so that it consists of a discriminator and a generator. It is a novel framework for HG embedding with GAN-based adversarial learning. The discriminator and generator are relation-aware in HeGAN. The discriminator is able to tell apart the specified relationship between nodes, and the generator will try to generate nodes which has certain relation with other nodes. The node pair is considered true in the following cases: (i) It is a positive pair based on network topology; (ii) The pair is formed under the correct relation. HeGAN provides a generalized generator, which is able to directly sample latent nodes from a continuous distribution, such that (i) no softmax computation is necessary; and (ii) fake samples are not restricted to the existing nodes.

Existing studies typically model the distribution of nodes using some forms of softmax over all nodes in the original network. In terms of effectiveness, their fake samples are limited to the nodes in the network, whereas the most representative fake samples may fall “in between” the existing nodes in the embedding space. For the bibliographic data in Fig. 6.1a, given a paper  $p_2$ , it can only choose fake samples from  $V$ , such as  $o_1$  and  $o_3$ . However, both may not be adequately similar to real samples such as  $o_2$ . Towards a better sample generation, HeGAN introduces a generalized generator that can produce latent nodes not appearing in the embedding

space. The generator can sample fake nodes directly without using a softmax. The framework of HeGAN is in Fig. 6.1d.

### 6.2.2.2 Relation-aware Discriminator

In an HG, the relation-aware discriminator  $\sigma_{D-A; \theta}$  evaluates the relation  $A$  between the pair of nodes  $D$  and  $E$ . Specifically,  $D \in V$  is a given node and  $A \in \mathcal{R}$  is a given relation from a HG  $G$ ,  $e_E$  is the embedding of a sample node  $E$  (which can be real or fake), and  $\theta$  denotes the model parameters of  $\sigma$ . In essence,  $\sigma$  outputs a probability that the sample  $E$  is connected to  $D$  under the relation  $A$ . We quantify this probability as:

$$\sigma_{D-A; \theta}(e_E) = \frac{1}{1 + \exp(-e_D^T A e_E)} \quad (6.1)$$

where  $e_E \in \mathbb{R}^{3-1}$  is the embedding of the sample  $E$ ,  $e_D \in \mathbb{R}^{3-1}$  is the learnable embedding of node  $D$ , and  $A \in \mathbb{R}^{3-3}$  is a learnable relation matrix for relation  $A$ .  $\theta = \{e_D : D \in V, A : A \in \mathcal{R}\}$  forms the model parameters of  $\sigma$ , i.e., the union of all node embeddings and relation matrices learnt by  $\sigma$ . Naturally, the probability should be high when  $E$  is a positive sample related to  $D$  through  $A$ , or low when it is a negative sample. In general, a sample  $E$  forms a triple  $\langle D, E, A \rangle$  together with the given  $D$  and  $A$ , and each triple belongs to one of the three cases below with regard to its polarity. Each case also contributes to one part of the discriminator loss, inspired by the idea of conditional GAN [26].

Case 1: Connected under a given relation. That is, nodes  $D$  and  $E$  are indeed connected through the right relation  $A$  on the HG  $G$ , such as  $\langle D_2, E_2, A_2 \rangle$  shown in Fig. 6.1a. Such a triple is considered positive and can be modeled by the below loss.

$$L_1 = \mathbb{E}_{\langle D, E, A \rangle \sim \mathcal{G}} \log \sigma_{D-A; \theta}(e_E) \quad (6.2)$$

Here we draw the positive triple from the observed  $G$ , denoted as  $\langle D, E, A \rangle \sim \mathcal{G}$ .

Case 2: Connected under an incorrect relation. That is,  $D$  and  $E$  are connected on the HG under a wrong relation  $A^0 \neq A$ , such as  $\langle D_2, E_2, A^0 \rangle$ . The discriminator is expected to mark them as negative, as their connectivity does not match the desired semantics carried by the given relation  $A$ . We define this part of loss as follows:

$$L_2 = \mathbb{E}_{\langle D, E, A^0 \rangle \sim \mathcal{G} \times \mathcal{R}^0} \log (1 - \sigma_{D-A^0; \theta}(e_E)) \quad (6.3)$$

Here, we still draw the pair of nodes  $\langle D, E \rangle$  from  $G$ , but the negative relation  $A^0$  is drawn from  $\mathcal{R}^0 = \mathcal{R} \setminus A$  uniformly.

Case 3: Fake node from a relation-aware generator. That is, given a node  $D \in V$ , it can form a fake pair with the node  $E$  supplied by the generator  $\sigma_{D-A; \theta}$ , such as  $\langle D_2, E_2, A_2 \rangle$  in Fig. 6.1d. The generator is also relation-aware: It attempts to generate a fake node's embedding that mimics the real nodes connected to  $D$  under

the correct relation. Again, the discriminator aims to identify this triple as negative, which can be formulated as follows.

$$L_3 = E_{(D, A, E) \sim \mathcal{D}} \log(1 - \sigma(\mathbf{1}_{D-A}^0 \cdot \mathbf{1}_E^0)) \quad (6.4)$$

Note that the fake samples embedding  $\mathbf{1}_E^0$  is drawn from the generator's learnt distribution. On the other hand, the discriminator simply treats  $\mathbf{1}_E^0$  as non-learnable input, and only optimizes its own parameters.

The loss function of the discriminator is comprehensive by the above three parts:

$$L = L_1 + L_2 + L_3 + \frac{\lambda}{2} \|\mathbf{g}\|_2^2 \quad (6.5)$$

where  $\lambda \geq 0$  controls the regularization term to avoid over fitting. The parameters  $\mathbf{g}$  of the discriminator can be optimized by minimizing  $L$ .

### 6.2.2.3 Relation-aware Generalized Generator

Given a node  $D \in V$  and a relation  $A \in R$ , the generator  $\mathbf{1}_{D-A}^0$  aims to generate a fake node  $E$  likely to connect to  $D$  in the context of relation  $A$ . In other words,  $E$  should be as close as possible to a real node  $E$ , which has  $D-E \in \mathcal{G}$ . On the other hand, the generator is generalized, which means the fake  $E$  can be latent and not found in  $V$ .

To meet the two requirements, the generator should also employ relation-specific matrices (for relation-awareness), and generate samples from an underlying continuous distribution (for generalization). One of the suitable distributions is Gaussian distribution:

$$N(\mathbf{1}_D^> \cdot \mathbf{1}_A, \Sigma_A - \mathbf{f}^2 \cdot \mathbf{1}_A) \quad (6.6)$$

where  $\mathbf{1}_D \in \mathbb{R}^{|V|}$  and  $\mathbf{1}_A \in \mathbb{R}^{|R|}$  denote the node embedding  $D \in V$  and the relation matrix of  $A \in R$  for the generator. In other words, it is a Gaussian distribution with mean  $\mathbf{1}_D^> \cdot \mathbf{1}_A$  and covariance  $\Sigma_A - \mathbf{f}^2 \cdot \mathbf{1}_A$  for some choices of  $\mathbf{f}$ . Intuitively, the mean represents a fake node is likely to be connected by relation  $A$ , and the covariance represents potential deviations. As neural networks have shown strong ability in modeling complex structure [15], we integrate the multi-layer perceptron (MLP) into the generator for enhancing the expression of the fake samples. Hence, our generator is formulated as follows,

$$\mathbf{1}_{D-A}^0 = \sigma(\mathbf{W}_1 \cdot \mathbf{1}_D + \mathbf{b}_1, \mathbf{W}_2 \cdot \mathbf{1}_A + \mathbf{b}_2) \quad (6.7)$$

where we draw  $\mathbf{1}_E$  from the distribution  $N(\mathbf{1}_D^> \cdot \mathbf{1}_A, \Sigma_A - \mathbf{f}^2 \cdot \mathbf{1}_A)$ . Here,  $\mathbf{W}_1$  and  $\mathbf{W}_2$  respectively denote the weight matrix and the bias vector for each layer, and  $\sigma$  is an activation function. The parameter set of the generator is thus  $\mathcal{G}_{D-A} = \{D \in V, A \in R, \mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \mathbf{f}\}$ , i.e., the union of all node embeddings and relation matrices, as well as the parameters of MLP.

As motivated earlier, the generator aims to fool the discriminator by generating close-to-real fake samples, such that the discriminator gives high score to them.

$$L = E_{\mathcal{D}-\mathcal{A}} \left[ \log \frac{1 + \exp(\mathcal{D}(\mathcal{A}))}{1 + \exp(\mathcal{D}(\mathcal{G}))} \right] + \lambda \|\theta\|_2^2 \quad (6.8)$$

where  $\lambda \geq 0$  controls the regularization term. The parameters of the generator can be optimized by minimizing  $L$ .

---

#### Algorithm 6.1 Model training for HeGAN.

---

Input: HG  $G$ , number of generator and discriminator trainings per epoch  $\tau$ , number of samples  $\mathcal{S}_B$

- 1: Initialize  $\mathcal{D}$  and  $\mathcal{G}$  for  $\mathcal{D}$  and  $\mathcal{G}$ , respectively
- 2: while not converged do
- 3:   for  $i = 0; i < \tau$  do • Discriminator training
- 4:     Sample a batch of triples  $\mathcal{E}_B$ ,  $\mathcal{D} = \mathcal{E}(\mathcal{A})$
- 5:     Generate  $\mathcal{S}_B$  fake nodes  $\mathcal{G}(\mathcal{A})$  for each  $\mathcal{D} = \mathcal{A}$
- 6:     Sample  $\mathcal{S}_B$  relations  $\mathcal{R}$  for each  $\mathcal{D} = \mathcal{E}$
- 7:     Update  $\mathcal{D}$  according to Eq. 6.5
- 8:   end for
- 9:   for  $i = 0; i < \tau$  do • Generator training
- 10:     Sample a batch of triples  $\mathcal{E}_B$ ,  $\mathcal{D} = \mathcal{E}(\mathcal{A})$
- 11:     Generate  $\mathcal{S}_B$  fake nodes  $\mathcal{G}(\mathcal{A})$  for each  $\mathcal{D} = \mathcal{A}$
- 12:     Update  $\mathcal{G}$  according to Eq. 6.8
- 13:   end for
- 14: end while
- 15: return  $\mathcal{D}$  and  $\mathcal{G}$

---

#### 6.2.2.4 Model Training

We adopt the iterative optimization strategy to train HeGAN. The iterative optimization methods follow the steps below: At each iteration, first, and then the generator generates fake samples to optimize and thus improve the discriminator. Next,  $\mathcal{G}$  is updated, followed by optimizing  $\mathcal{D}$  in order to produce increasingly better fake samples as evaluated by the discriminator. The above steps repeat until the model converges. The model training for HeGAN is outlined in Algorithm 6.1.

### 6.2.3 Experiments

#### 6.2.3.1 Experimental Settings

We conduct extensive experiments on three benchmark datasets [10, 15], namely DBLP (with four types of nodes: author, paper, conference and term), Yelp (user, business, service, star and reservation), AMiner (author, paper, conference and

reference. We organize them into HGs, as summarized in Table 6.1. We perform node classification and link prediction on three datasets.

Table 6.1 Description of datasets.

| Datasets | #Nodes  | #Edges  | #Node types | #Labels |
|----------|---------|---------|-------------|---------|
| DBLP     | 37,791  | 170,794 | 4           | 4       |
| Yelp     | 3,913   | 38,680  | 5           | 3       |
| Aminer   | 312,776 | 599,951 | 4           | 6       |

For the baselines, we consider three categories of network embedding methods: Traditional (Deepwalk [25], LINE [28]), GAN-based (GraphGAN [30], ANE [6]) and HG (HERec-HNE [27], HIN2vec [10], Metapath2vec [8]) embedding algorithms.

### 6.2.3.2 Node Classification

We use 80% of the labeled nodes to train a logistic regression classifier, and test the classifier on the remaining 20% of the nodes. We report accuracy on the test set in Table 6.2.

Table 6.2 Performance comparison on node classification. (bold: best; underline: runner-up).

| Methods      | DBLP     |          |          | Yelp     |          |          | Aminer   |          |          |
|--------------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
|              | Micro-F1 | Macro-F1 | Accuracy | Micro-F1 | Macro-F1 | Accuracy | Micro-F1 | Macro-F1 | Accuracy |
| Deepwalk     | 0.9201   | 0.9242   | 0.9298   | 0.8262   | 0.7551   | 0.8145   | 0.9519   | 0.9460   | 0.9529   |
| LINE-1st     | 0.9239   | 0.9213   | 0.9285   | 0.8229   | 0.7440   | 0.8126   | 0.9776   | 0.9713   | 0.9788   |
| LINE-2nd     | 0.9144   | 0.9172   | 0.9236   | 0.7591   | 0.5518   | 0.7571   | 0.9469   | 0.9341   | 0.9471   |
| GraphGAN     | 0.9198   | 0.9210   | 0.9286   | 0.8098   | 0.7268   | 0.7820   | -        | -        | -        |
| ANE          | 0.9143   | 0.9153   | 0.9189   | 0.8232   | 0.7623   | 0.7932   | 0.9256   | 0.9203   | 0.9221   |
| HERec-HNE    | 0.9214   | 0.9228   | 0.9299   | 0.7962   | 0.7713   | 0.7912   | 0.9801   | 0.9726   | 0.9784   |
| HIN2vec      | 0.9141   | 0.9115   | 0.9224   | 0.8352   | 0.7610   | 0.8200   | 0.9799   | 0.9775   | 0.9801   |
| Metapath2vec | 0.9288   | 0.9296   | 0.9360   | 0.7953   | 0.7884   | 0.7839   | 0.9853   | 0.9860   | 0.9857   |
| HeGAN        | 0.9381   | 0.9375   | 0.9421   | 0.8524   | 0.8031   | 0.8432   | 0.9864   | 0.9873   | 0.9883   |

In node classification, HeGAN consistently outperforms the best baseline with statistical significance. It is also worth noting that our performance margins over the best baseline become smaller compared to the node clustering task, since in classification all methods are helped by the supervision, narrowing their gaps.

### 6.2.3.3 Link Prediction

In this task, we predict user-business links on Yelp, and author-paper links on DBLP and AMiner. We randomly hide 20% of such links from the original network as the ground truth positives, and randomly sample disconnected node pairs of the given form as negative instances. The ground truth serves as our test set. We adopt two

Table 6.3 Performance comparison on link prediction. (bold: best; underline: runner-up).

|            | Methods      | DBLP   |        |        | Yelp          |               |        | AMiner |        |        |
|------------|--------------|--------|--------|--------|---------------|---------------|--------|--------|--------|--------|
|            |              | Acc    | AUC    | F1     | Acc           | AUC           | F1     | Acc    | AUC    | F1     |
| Logistic   | Deepwalk     | 0.5441 | 0.5630 | 0.5208 | 0.7161        | 0.7825        | 0.7182 | 0.4856 | 0.5182 | 0.4618 |
|            | LINE-1st     | 0.6546 | 0.7121 | 0.6685 | <u>0.7226</u> | <u>0.7971</u> | 0.7099 | 0.5983 | 0.6413 | 0.6080 |
|            | LINE-2nd     | 0.6711 | 0.6500 | 0.6208 | 0.6335        | 0.6745        | 0.6499 | 0.5604 | 0.5114 | 0.4925 |
|            | GraphGAN     | 0.5241 | 0.5330 | 0.5108 | 0.7123        | 0.7625        | 0.7132 | -      | -      | -      |
| Regression | ANE          | 0.5123 | 0.5430 | 0.5280 | 0.6983        | 0.7325        | 0.6838 | 0.5023 | 0.5280 | 0.4938 |
|            | HERec-HNE    | 0.7123 | 0.7823 | 0.6934 | 0.7087        | 0.7623        | 0.6923 | 0.7089 | 0.7776 | 0.7156 |
|            | HIN2vec      | 0.7180 | 0.7948 | 0.7006 | 0.7219        | 0.7959        | 0.7240 | 0.7142 | 0.7874 | 0.7264 |
|            | Metapath2vec | 0.5969 | 0.5920 | 0.5698 | 0.7124        | 0.7798        | 0.7106 | 0.7069 | 0.7623 | 0.7156 |
|            | HeGAN        | 0.7290 | 0.8034 | 0.7119 | 0.7240        | 0.8075        | 0.7325 | 0.7198 | 0.7957 | 0.7389 |
| Inner      | Deepwalk     | 0.5474 | 0.7231 | 0.6874 | 0.5654        | 0.8164        | 0.6953 | 0.5309 | 0.6064 | 0.6799 |
|            | LINE-1st     | 0.6647 | 0.7753 | 0.7363 | 0.6769        | 0.7832        | 0.7199 | 0.6113 | 0.6899 | 0.7123 |
|            | LINE-2nd     | 0.4728 | 0.4797 | 0.6325 | 0.4193        | 0.7347        | 0.5909 | 0.5000 | 0.4785 | 0.6666 |
|            | GraphGAN     | 0.5532 | 0.6825 | 0.6214 | 0.5702        | 0.7725        | 0.6894 | -      | -      | -      |
| Product    | ANE          | 0.5218 | 0.6543 | 0.6023 | 0.5432        | 0.7425        | 0.6324 | 0.5421 | 0.6123 | 0.6623 |
|            | HERec-HNE    | 0.5123 | 0.7473 | 0.6878 | 0.5323        | 0.6756        | 0.7066 | 0.6063 | 0.6912 | 0.6798 |
|            | HIN2vec      | 0.5775 | 0.8295 | 0.6714 | 0.6273        | 0.8340        | 0.4194 | 0.5348 | 0.6934 | 0.6824 |
|            | Metapath2vec | 0.4775 | 0.6926 | 0.6287 | 0.5124        | 0.6324        | 0.6702 | 0.6243 | 0.7123 | 0.6953 |
|            | HeGAN        | 0.7649 | 0.8712 | 0.7837 | 0.7391        | 0.8298        | 0.7705 | 0.6505 | 0.7431 | 0.7752 |

ways to perform link prediction, namely logistic regression and inner product. The evaluation metrics consist of accuracy, AUC and F1, as shown in Table 6.3.

We observe that, the performance of HeGAN over the best baseline is much larger with inner product than with logistic regression. It is hypothesized that HeGAN inately learns a much better structure and semantics preserving embedding space than the baseline methods, since the inner product only relies on the learnt representations without resorting to any external supervision.

The complete method and more experiments can be found in [16].

## 6.3 Importance Sampling

### 6.3.1 Overview

To make large-scale graph representation learning possible, researchers have proposed several sampling strategies on Graph Neural Networks (GNNs), including node-wise neighborhood sampling [14, 35] and layer-wise neighborhood sampling [5, 17, 39] to accelerate the training process of GNNs by reducing the number of edges in computation. The former is to sample neighbors for each node while the latter is to sample neighbors from the whole graph. Unfortunately, both the node- and layer-wise sampling only deal with homogeneous graphs, which are inadequate in many real-world scenarios such as E-commerce graphs, as they (1) deal with homogeneous graphs only; (2) take all the nodes of a graph as initial candidates; (3) still incur rapidly increasing cost with more layers. There are two challenges we have to face with as follows.



First, how to design an effective sampler that works with the heterogeneous neighborhoods? While a sampler can be easily defined on a per-node, per-type basis, it is not clear how we can sample neighborhoods in a batch, given that the candidates consist of different types of nodes. Two alternatives could work, namely, type-dependent and type-fusion sampling. In the former, a sampler is deployed for each neighbor type, and these type-based samplers sample from their respective sub-neighborhood on their own. In the latter, a single sampler is deployed, which treats the entire common neighborhood of the batch as its candidates. Intuitively, the type-fusion sampler would work better, as it takes the influences of total types of interactions into account and models the sampling distribution over all types jointly.

Second, how to design the corresponding effective estimators with the sampled heterogeneous neighborhoods? With the neighbors being sampled based on the global importance in a batch, traditional importance sampling could introduce unwanted variance because of the imbalance between the local importance to the given target node and the global importance to the given batch. To address this challenge, we propose self-normalized estimators and adaptive estimators respectively. The former is to adjust the estimators by self-normalizing importance of sampled neighborhoods while the latter is to automatically learn the global importance (i.e., the importance distribution of candidate neighborhoods) by modeling both structural and attributed information.

Finally, we test the proposed sampling strategies on two public datasets, and the experimental results demonstrate that the proposed framework can achieve statistically significant improvements. Compared with the full model without any sampling, we achieve a memory cost reduction by up to 92.48% and a time cost reduction by up to 85.95%, while maintaining the same level of accuracy on two real-world datasets.

### 6.3.2 The HeteSamp Model

This section introduces the architecture of the proposed HeteSamp model, which consists of the general heterogeneous interaction graph embedding (named HIGE) and the heterogeneous sampling framework. The heterogeneous sampling framework leverages type-dependent and type-fusion strategies to overcome the expensive time cost and high computational complexity. Notice that, we name the large-scale heterogeneous graph as heterogeneous interaction graph to emphasize the rapidly increasing scale. Moreover, as HeteSamp consists of multiple samplers and estimators, we name the specific strategies according to their adopted sampler and estimator.

#### 6.3.2.1 The General HIGE

To tackle with HIGs, a general idea is to reconstruct node embeddings by propagating information from its heterogeneous neighbors, and backwardly propagate

Fig. 6.2 The sampling processes of type-dependent and type-fusion sampling.

gradients [32, 38, 4], as shown in Fig. 6.2a. Specifically, given a node  $E_8$  with type  $q^1 E_8^0$  and its edge  $E_8 E_9$  of type  $A$ , the aggregated information of node  $E_8$  from type- $A$  neighborhoods is

$$\mathbf{G}_{E_8-A}^0 = \frac{1}{N_{E_8-A}} \sum_{E_9 \in N_{E_8-A}} F^1_{E_8 E_9 A} \mathbf{E}_9 \quad (6.9)$$

where  $\mathbf{G}_{E_8-A}^0$  is the aggregated information,  $N_{E_8-A}$  is the type- $A$  neighborhoods of node  $E_8$ ,  $E_9 \in N_{E_8-A}$  is a specific type- $A$  neighbor of node  $E_8$ ,  $F^1_{E_8 E_9 A}$  denotes the weight of edge  $E_8 E_9$ ,  $\mathbf{E}_9 \in \mathbb{R}^3$  denotes the embedding of node  $E_9$ , 3 is the dimension of  $\mathbf{E}_9$ .  $F^1_{E_8 E_9 A}$  between node  $E_8$  and  $E_9$  is calculated by

$$F^1_{E_8 E_9 A} = \sigma(\mathbf{G}_{E_8-A} \cdot \mathbf{1}_A) \quad (6.10)$$

where  $\sigma$  is an activation function,  $\mathbf{G}_{E_8-A} \in \mathbb{R}^{3 \times 3_A}$  is the edge features between node  $E_8$  and  $E_9$ ,  $3_A$  is the length of type- $A$  edge features,  $\mathbf{1}_A \in \mathbb{R}^{3_A \times 1}$  and  $\mathbf{1}_A \in \mathbb{R}$  are the weight and bias parameters shared by type- $A$  edges.

We then take all semantics into account by aggregating information from different-typed neighborhoods and constructing the global embedding  $\mathbf{G}_{E_8}^0 \in \mathbb{R}^3$  of node  $E_8$  as

$$\mathbf{G}_{E_8}^0 = \sigma(\mathbf{2} \gg \bigoplus_{E_9-A} \mathbf{G}_{E_8-A}^0 - \mathbf{G}_{E_8-A}^0, \mathbf{1}_{q^1 E_8^0}, \mathbf{1}_{q^1 E_8^0}) \quad (6.11)$$

where  $\sigma$  is an activation function,  $\mathbf{2} \gg$  is the concatenation option,  $\mathbf{1}_{q^1 E_8^0} \in \mathbb{R}^{j \times 3}$  denotes the projection matrix of type- $E_8^0$  nodes,  $\mathbf{1}_{q^1 E_8^0}$  denotes the bias and  $j$  denotes the total number of edge types in the heterogeneous graph. The computational complexity of the general model is linear with the scale of edges and nodes on heterogeneous graphs.

### 6.3.2.2 Batch-Wise Heterogeneous Sampling

To reduce computational overhead and memory cost of the training process, a naive idea is to sample several neighbors rather than aggregating information from all neighborhoods. Previous works [14, 35] on homogeneous works usually adopt node-wise sampling to sample several neighbors per node for learning. Based on the global and local structural information, current works [5, 17, 39] propose layer-wise sam-

pling strategies which consider the whole neighborhood as the candidates. However, these strategies are to deal with homogeneous graphs but ignore the abundant semantics of heterogeneous nodes and edges. Moreover, these strategies can only deal with small graphs because the overhead for node-wise and layer-wise sampling during optimization could be quite expensive, or even unaffordable.

In this section, Eq. (6.9) can be re-written as

$$\mathbf{I}_{\mathcal{E}_B-A}^0 = \mathbb{E}_{\mathcal{E}_B} \left[ \frac{\sum_{j \in \mathcal{E}_B} \mathbf{I}_{\mathcal{E}_B-A}^j}{|\mathcal{E}_B|} \right] \quad (6.12)$$

where  $\mathcal{E}_B$  denotes the target nodes in the mini-batch,  $\mathbf{I}_{\mathcal{E}_B-A}^j$  denotes the corresponding sampling probability in the mini-batch, or in other words, the importance in this batch, and  $\sum_{j \in \mathcal{E}_B} \mathbf{I}_{\mathcal{E}_B-A}^j$  equals to  $\frac{1}{|\mathcal{N}_{\mathcal{E}_B-A}|}$ . Thereby, we attempt to make heterogeneous sampling for each batch. The batch-wise heterogeneous sampling probability is defined as

$$\mathbf{I}_{\mathcal{E}_B-A}^j = \frac{|\mathcal{N}_{\mathcal{E}_B-A}^j|}{|\mathcal{N}_{\mathcal{E}_B-A}|} \quad (6.13)$$

where  $\mathcal{N}_{\mathcal{E}_B-A}$  is the sampled neighbor. As shown in Fig. 6.2a, we sample instances from the union neighborhood of all the target nodes in a batch, where the union neighborhood denotes the union of the individual neighbors of each target node. Since sampling is now done for each batch instead of each target node, batch-wise sampling is often a good choice to reduce computational overhead and memory cost of the training process. Compared with layer-wise sampling methods [5, 17, 39] which require the total nodes as candidates and have to load the whole graph structure before training, our batch-wise heterogeneous sampling focuses on the neighborhood union in the current batch. Compared with node-wise sampling which sample neighborhoods for each target node, our batch-wise heterogeneous sampling contains the advantages of reducing the overhead of sampling.

### 6.3.2.3 Type-Dependent Sampling Strategy

By now, we have defined the general batch-wise heterogeneous sampling. Different from traditional importance sampling, the neighborhoods of each batch connect with target nodes based on different-typed interactions. A straightforward idea is to design a sampler for each type, respectively. For example, in Fig. 6.2b, there are four types of candidates, and we sample neighbors from each type of candidates respectively by using one sampler. By adopting type-dependent sampling, we consider  $\mathcal{E}_B$  as a set of  $\mathcal{E}_B^t$ , and sample type- $t$  neighborhoods with the corresponding sampler  $\mathbf{I}_{\mathcal{E}_B-A}^t$ . The remaining problem for type-dependent sampling is how to design the exact form of this sampler so as to keep low variance for efficient training. Here we define the average information of  $\mathbf{I}_{\mathcal{E}_B-A}^0$  as

12

$$\bar{y}_A = \frac{1}{j} \sum_{j \in E_A} \tilde{y}_j = \frac{1}{j} \sum_{j \in E_A} \tilde{y}_j \quad \text{where } \tilde{y}_j = \frac{y_j - \mu}{\sigma} \quad (6.14)$$

where  $j \in E_A$  is the number of type  $A$  samples. Then, the variance of  $\bar{y}_A$  is calculated by

$$\text{Var}(\bar{y}_A) = \frac{1}{j} \sum_{j \in E_A} \frac{\sum_{j \in E_A} (y_j - \mu)^2}{j} \quad (6.15)$$

To ensure minimizing the variance, a better sampler is shown as

$$\mathbb{E}_A^0 = \frac{\sum_{j \in E_A} (y_j - \mu)^2}{\sum_{j \in E_A} (y_j - \mu)^2} \quad (6.16)$$

### 6.3.2.4 Type-Fusion Sampling Strategy

Essentially, the type-dependent strategy pays attention to the influence of same-type neighborhoods, without jointly considering the effect of heterogeneous types. Thus, they only reduce the individual variance of each type, lacking a global picture on the overall variance of the batch. To address this weakness, we propose a type-fusion sampling strategy which considers the entire neighborhoods as the candidates.

The aggregation in Eq. (6.11) are to gather information from different-typed neighborhoods, and can be rewritten as

$$\mathbb{E}_B^0 = \frac{\sum_{j \in E_B} (y_j - \mu)^2}{\sum_{j \in E_B} (y_j - \mu)^2} \quad (6.17)$$

where,  $A \in \mathbb{R}^{3 \times 3}$  denotes the relation-wise projection matrix. Under the type-fusion strategy, the average information-rich batch over all types is given by

$$\bar{y}_B = \frac{1}{j} \sum_{j \in E_B} \tilde{y}_j = \frac{1}{j} \sum_{j \in E_B} \tilde{y}_j \quad \text{where } \tilde{y}_j = \frac{y_j - \mu}{\sigma} \quad (6.18)$$

As shown in Fig. 6.2c, the different-typed neighborhoods are sampled according to the type-fusion sampling distribution. Similar to that in Section 6.3.2.3, the corresponding sampler is defined as follows.

$$\mathbb{E}_B^0 = \frac{\sum_{j \in E_B} (y_j - \mu)^2}{\sum_{j \in E_B} (y_j - \mu)^2} \quad (6.19)$$

where  $\mathbb{E}_B^0$  is the structure-based type-fusion sampler.

While the above sampling strategies only consider link-based weight as the importance of nodes but ignoring the attributed-based importance on edges, we further pay attention to the edge features within HIGs and propose the corresponding sampler as follows

$$\mathbb{Q}^1 E_g^0 = \frac{\sum_{E_g} \frac{1}{A_{2R}} \sum_{E_{g2}} \frac{1}{E_g} \frac{?^1 E_{gj} E_g - A_{5^1} G E_g - E_g - A^0}{?^1 E_{gj} E_g - A_{5^1} G E_g - E_g - A^0}}{A_{2R} E_g 2^N \dots} \quad (6.20)$$

where  $\mathbb{Q}^1 E_g^0$  denotes the adaptive sampler,  $G E_g - E_g - A^0$  denotes the importance of edge features which is calculated by  $G E_g - E_g - A$ ,  $G - \lambda$  and,  $G - A R^{3A - 1}$  is the type-A learnable parameter.

### 6.3.2.5 Heterogeneous Self-Normalized and Adaptive Estimators

In type-dependent or type-fusion strategies, the batch-wise estimator is

$$\mathbf{6}_{E_g} = \frac{1}{A_{2R} E_g 2^N \dots} \frac{\sum_{E_g} \frac{?^1 E_{gj}^0}{@ E_{gj}^0} F^1 E_g - E_g - A_{E_g, A^-}}{?^1 E_{gj}^0} \quad (6.21)$$

where  $\mathbf{6}_{E_g}$  is the approximated information,  $n$  is the number of samples,  $\dots$  is the sampled type-A neighborhoods in this batch,  $?^1 E_{gj}^0$  is equal to  $\frac{1}{j N_{E_g - A_j}}$  and  $@ E_{gj}^0$  denotes the heterogeneous samplers, such as  $\mathbb{Q}_1$  and  $\mathbb{Q}_2$ . However, this could increase the variance resulted from the imbalance  $?^1 E_{gj}^0$  and  $@ E_{gj}^0$ .

To address such problem, for structure-based  $\mathbb{Q}^1 E_g^0$ , a promising way is to balance the weights based on self-normalized importance, similar to that in [23]. The corresponding estimator can be computed as

$$\mathbf{6}_{B=-E} = \frac{\sum_{E_g} \frac{1}{A_{2R} E_g 2^N \dots} \frac{c^1 E_g^0}{E_g 2^N \dots} F^1 E_g - E_g - A_{E_g, A^-}}{c^1 E_g^0} \quad (6.22)$$

where  $\mathbf{6}_{B=-E}$  denotes the self-normalized information,  $c^1 E_g^0 = \frac{?^1 E_{gj}^0}{@ E_{gj}^0}$ ,  $N_A$  is the sampled neighborhoods of type-A. Besides, for the adaptive sampling, we add the variance to the loss function and explicitly minimize the variance during training to achieve a better  $\mathbb{Q}^1 E_g^0$ .

### 6.3.2.6 Optimization Framework

The overall loss function consists of four parts

$$\mathcal{L} = \mathcal{L}_{C0B:-, U} + \mathcal{L}_{4?-:, V} + \mathcal{L}_{1^0, b+0A_{E_g:-, 1^0} @} \quad (6.23)$$

where  $\mathcal{L}_i$  is the loss value in  $i$ -th batch,  $\mathcal{L}_{\text{sup}}$  is the loss from supervised learning in the same batch,  $\mathcal{L}_{\text{prop}}$  is the embedding propagation loss with sampling in the  $i$ -th batch,  $\sigma^2$  is the variance of sampled information,  $\lambda$  is the regularization of all latent parameters,  $\alpha$ ,  $\beta$  and  $\gamma$  are three hyper-parameters. Notice that, for type-dependent sampling and structure-based type-fusion sampling,

### 6.3.3 Experiments

#### 6.3.3.1 Experimental Settings

**Datasets.** We evaluate the empirical performance of our method on two real-world heterogeneous graphs, including the small Aminer graph consisting of 41,523 nodes and 199,429 edges, and the large Alibaba graph consisting of 4,527,222 nodes and 49,785,900 edges.

**Baseline.** We first compare our various importance sampling strategies with the general HIGE model and HAN [32] without any sampling. We also compare our model with state-of-the-art sampling algorithms on GCNs including Fast-GCN [17] and AS-GCN [5] to showcase our advantage of effectiveness and scalability, whereas the sampling-based models not only achieve superior accuracy, but also scale to large graphs. Finally, to study the performance of our proposed variance reduction, we also substitute our variance reduction sampler with a uniform sampler.

**Parameter Settings.** For all baselines and our methods, we set  $\alpha = 0.1$ ,  $\beta = 0.1$ ,  $\gamma = 0.1$ ,  $U = 0.4$ . The maximum iteration is 100 for Aminer and 5 for Alibaba. We adopt Micro-F1 and Macro-F1 as the evaluation metrics for node classification on Aminer while adopting F1 and AUC as the evaluation metrics for link prediction on Alibaba. All the metrics are positively related to the performance of methods.

#### 6.3.3.2 Empirical Validation

We perform empirical validation on the small Aminer graph for node classification and the much larger Alibaba interaction graph for link prediction. Notice that the non-HIGE-based baselines, GraphsAGE, Fast-GCN and AS-GCN are to deal with node classification, and these baselines can only work on Aminer; they cannot perform on Alibaba graph when the task is link prediction and this graph is too large. Furthermore, they are likely to suffer from insufficient memory on the large-scale Alibaba graph even if we implement the modified models for the task of link prediction.

(a) Aminer Memory (b) Alibaba Memory (c) Aminer Time (d) Alibaba Time

Fig. 6.3 Average memory cost and running time per iteration of training.

### 6.3.3.3 Effectiveness

We report the results in Table 6.4 and Table 6.5, respectively. We progressively sample more neighbors per batch and the sampling rate ranges from about 3% to 24% on smaller Aminer and 1.25% to 10% on the larger Alibaba graph. We make the following observations.

First, VarR-TF-AS generally achieves the best performance, whereas VarR-TF-SN comes as a close competitor. In particular, VarR-TF-AS performs as well as or even better than the original HIGE. This phenomenon is reasonable. On the one hand, HIGE aggregates information from all neighbors which may introduce noisy information while our sampling strategies are to sample valuable neighbors for training. On the other hand, we optimize the sampler by reducing the variance loss which enhances the similarity limitation between the sampled neighbors and its target nodes.

Second, compared to the corresponding uniform samplers, the variance reduction samplers guarantee more stable estimators and thus produce better results. In the VarR-\* methods, type-fusion strategies outperform type-dependent methods, as the former consider all types jointly rather than independently, and reduce the variance of the whole batch rather than a single type.

Third, our proposed general HIGE can perform competitively with or even better than HAN. However, HAN has higher time complexity. For Alibaba dataset, the time cost of HAN (about 1 day) is quite larger than HIGE (about 5 hours).

Forth, all sampling strategies for HIGE perform significantly better than GCN-based models and GraphsAGE. On the one hand, sampling for HIGE takes the graph heterogeneity into consideration, whereas the two GCN-based models do not make use of such information. On the other hand, the sampling size or number of layers of Fast-GCN and AS-GCN may be not enough. However, even under current settings, Fast-GCN and AS-GCN are already several times slower than our method, as we shall see in the efficiency study.

Table 6.4 Micro/Macro-F1 scores for node classification on Aminer. Excluding HIGE and HAN, the best method is bolded and the second best is underlined.

| Sampling size<br>per batch | Micro-F1  |           |            |             | Macro-F1  |           |            |             |
|----------------------------|-----------|-----------|------------|-------------|-----------|-----------|------------|-------------|
|                            | 128<br>3% | 256<br>6% | 512<br>12% | 1024<br>24% | 128<br>3% | 256<br>6% | 512<br>12% | 1024<br>24% |
| HIGE-Nil                   | 0.1990    |           |            |             | 0.1961    |           |            |             |
| HIGE                       | 0.9646    |           |            |             | 0.9593    |           |            |             |
| HAN                        | 0.9512    |           |            |             | 0.9508    |           |            |             |
| GraphSAGE                  | 0.2022    | 0.2039    | 0.2125     | 0.2119      | 0.1989    | 0.2023    | 0.2036     | 0.2073      |
| Fast-GCN                   | 0.2117    | 0.2244    | 0.2318     | 0.2361      | 0.1850    | 0.1898    | 0.2116     | 0.2119      |
| AS-GCN                     | 0.2361    | 0.2307    | 0.2390     | 0.2390      | 0.2005    | 0.2028    | 0.2004     | 0.2068      |
| Unif-TD                    | 0.3043    | 0.4123    | 0.4256     | 0.4221      | 0.2638    | 0.3907    | 0.3553     | 0.3962      |
| Unif-TF                    | 0.1780    | 0.2229    | 0.3785     | 0.6296      | 0.1266    | 0.1952    | 0.3081     | 0.5727      |
| VarR-TD                    | 0.8086    | 0.7990    | 0.8971     | 0.9123      | 0.7913    | 0.7894    | 0.8945     | 0.9133      |
| VarR-TF                    | 0.9461    | 0.9671    | 0.9712     | 0.9675      | 0.9424    | 0.9651    | 0.9696     | 0.9664      |
| VarR-TF-SN                 | 0.9659    | 0.9643    | 0.9612     | 0.9684      | 0.9637    | 0.9629    | 0.9631     | 0.9603      |
| VarR-TF-AS                 | 0.9650    | 0.9676    | 0.9705     | 0.9687      | 0.9649    | 0.9667    | 0.9602     | 0.9671      |

Table 6.5 F1 and AUC scores for purchase prediction on the Alibaba graph. Excluding HIGE and HAN, the best method is bolded and the second best is underlined.

| Sampling size<br>per batch | F1           |              |            |             | ROC-AUC      |              |            |             |
|----------------------------|--------------|--------------|------------|-------------|--------------|--------------|------------|-------------|
|                            | 512<br>1.25% | 1024<br>2.5% | 2048<br>5% | 4096<br>10% | 512<br>1.25% | 1024<br>2.5% | 2048<br>5% | 4096<br>10% |
| HIGE-Nil                   | 0.3994       |              |            |             | 0.5134       |              |            |             |
| HIGE                       | 0.5663       |              |            |             | 0.7715       |              |            |             |
| HAN                        | 0.5618       |              |            |             | 0.7704       |              |            |             |
| Unif-TD                    | 0.4017       | 0.4226       | 0.4352     | 0.4371      | 0.5768       | 0.5826       | 0.5908     | 0.5924      |
| Unif-TF                    | 0.4008       | 0.4122       | 0.4125     | 0.4451      | 0.5731       | 0.5790       | 0.5862     | 0.5977      |
| VarR-TD                    | 0.4841       | 0.5003       | 0.5274     | 0.5682      | 0.6207       | 0.6504       | 0.6925     | 0.7475      |
| VarR-TF                    | 0.5769       | 0.5908       | 0.5709     | 0.5833      | 0.7648       | 0.7671       | 0.7653     | 0.7796      |
| VarR-TF-SN                 | 0.5780       | 0.5883       | 0.5802     | 0.5798      | 0.7669       | 0.7625       | 0.7660     | 0.7742      |
| VarR-TF-AS                 | 0.5729       | 0.5913       | 0.5806     | 0.5844      | 0.7674       | 0.7641       | 0.7676     | 0.7799      |

#### 6.3.3.4 Efficiency

We first investigate the efficiency of our sampling strategies in memory cost. As shown in Fig. 6.3a and Fig. 6.3b, our models incur by up to 92.48% less memory. Further note that the differences in both the number of edges and memory cost are more prominent on the larger Alibaba graph, indicating excellent scalability of our sampling strategies. Second, in terms of the running time, as shown in Fig. 6.3c and Fig. 6.3d, VarR-TF-AS and VarR-TF-SN require less time than HIGE to attain close performance on the two datasets. Compared to HIGE, our models incur by up to 84.39% less time cost. Notice that when dealing with large-scale Alibaba dataset, the time cost is quite larger (about 1 day) than HIGE (about 5 hours).

The specific designs of HeteSamp can be found in [19].



## 6.4 Hyperbolic Representation

### 6.4.1 Overview

HG embedding has attracted considerable research attention recently. Most HG embedding methods choose Euclidean spaces to represent HG graphs, which is because Euclidean spaces are the natural generalization of our intuition-friendly, and visual three-dimensional space. However, a fundamental problem is that what are the appropriate or intrinsic underlying spaces of HGs? Therefore, we wonder whether Euclidean spaces are the intrinsic spaces of HGs?

Recently, hyperbolic spaces have gained momentum in the context of network science [20]. Hyperbolic spaces are spaces of constant negative curvature [3]. A superiority of hyperbolic spaces is that they expand faster than Euclidean spaces [22]. Therefore, it is easy to model complex data with low-dimensional embedding in hyperbolic spaces. Due to the characteristic of hyperbolic spaces, [20] assumes hyperbolic spaces underlie complex network and finds that data with power-law structure is suitable to be modeled in hyperbolic spaces. Also, some researchers begin to embed different data in hyperbolic spaces. For instance, [7] embeds text in hyperbolic spaces. [22] and [11] learn the hyperbolic embeddings of homogeneous networks. However, it is unknown whether HGs are suitable to be embedded in hyperbolic spaces.

In this section, we analyze the relation distribution in HGs and propose HHNE, which is able to preserve the structure and semantic information in hyperbolic spaces. HHNE leverages the meta-path guided random walk to generate heterogeneous neighborhoods to capture the structure and semantic relations in HGs. Then the proximity between nodes is measured by the distance in hyperbolic spaces. Also, HHNE is able to maximize the proximity between the neighborhood nodes while minimize the proximity between the negative sampled nodes. The optimization strategy of HHNE is derived to optimize hyperbolic embeddings iteratively.

### 6.4.2 The HHNE model

#### 6.4.2.1 Model Framework

HHNE leverages the meta-path guided random walk to obtain neighbors for each node to capture the structure and semantic relations in HGs. Also, HHNE learns the embeddings by maximizing the proximity between the neighborhood nodes and minimizing the proximity between the negative sampled nodes. Moreover, we derive the optimization strategy of HHNE to upgrade the hyperbolic embeddings.

6.4.2.2 Hyperbolic HG Embedding

To design HG embedding methods in hyperbolic spaces, we make use of the Poincaré ball model to describe hyperbolic spaces. Let  $\mathbb{D}^3 = \{g \in \mathbb{R}^3 : \|g\| < 1\}$  be the open 3-dimensional unit ball. The Poincaré ball model is defined by the manifold  $\mathbb{D}^3$  equipped with the following Riemannian metric tensor

$$g_{\mathbb{D}^3} = \frac{2}{1 - \|g\|^2} g^E \quad \text{where } g^E := \frac{2}{1 - \|g\|^2} \quad (6.24)$$

where  $g^E$  denotes the Euclidean metric tensor.

HHNE aims to learn the representation of nodes to preserve the structure and semantic correlations in hyperbolic spaces. Given an HG  $(\mathcal{V}, \mathcal{E})$  with  $\mathcal{V} = \{v_i\}_{i=1}^n$ , HHNE is interested in learning the embeddings  $\{v_i\}_{i=1}^n \in \mathbb{D}^3$ . HHNE preserves the structure by facilitating the proximity between a node and its neighborhoods. HHNE uses meta-path guided random walks [8] to obtain heterogeneous neighborhoods of a node. In meta-path guided random walks, the node sequences are restrained by the node types which are defined by meta-paths. Specifically, let  $\mathcal{C}_v$  and  $\mathcal{C}_e$  as the types of nodes and edges, respectively, given a meta-path

$P = v_1 \xrightarrow{e_1} v_2 \xrightarrow{e_2} \dots \xrightarrow{e_{l-1}} v_l \xrightarrow{e_l} v_{l+1}$  the transition probability at step  $i$  is defined as follows:

$$P(v_{i+1} | v_i, e_i) = \begin{cases} \frac{1}{\#\{v_{i+1} \in \mathcal{C}_{v_{i+1}}\}} & \text{if } v_{i+1} \in \mathcal{C}_{v_{i+1}} \\ 0 & \text{otherwise} \end{cases} \quad (6.25)$$

where  $v_{i+1}$  is node  $v_{i+1}$  with type  $\mathcal{C}_{v_{i+1}}$ , and  $\#\{v_{i+1} \in \mathcal{C}_{v_{i+1}}\}$  denotes the number of nodes in the neighborhood of node  $v_i$ . The meta-path guided random walk strategy ensures that the semantic relationships between different types of nodes can be properly incorporated into HHNE.

In order to preserve the proximity between nodes and its neighborhoods in hyperbolic spaces, HHNE uses distances in Poincaré ball model to measure their proximity. Given nodes embeddings  $\{v_i\}_{i=1}^n \in \mathbb{D}^3$ , the distance in Poincaré ball is given by:

$$d_{\mathbb{D}^3}(v_i, v_j) = \cosh^{-1} \left( 1 + 2 \frac{\|v_i - v_j\|^2}{(1 - \|v_i\|^2)(1 - \|v_j\|^2)} \right) \quad (6.26)$$

It is worth noting that as the Poincaré ball model is defined in metric spaces, the distance in Poincaré ball meets the triangle inequality and can well preserve the transitivity in HG. Then, HHNE uses a probability to measure the neighborhood of node  $v_i$  as follows:

$$P(v_j | v_i) = \frac{1}{\sum_{v_k \in \mathcal{N}(v_i)} \exp(-d_{\mathbb{D}^3}(v_i, v_k))}$$

where  $\mathcal{N}(v_i) = \{v_j \in \mathcal{V} : (v_i, v_j) \in \mathcal{E}\}$ . Then the object of HHNE is to maximize the probability as follows:

$$\arg \max_{\{v_i\}_{i=1}^n} \log P(v_j | v_i) \quad (6.27)$$

To achieve efficient optimization, HHNE leverages the negative sampling proposed in [21], which basically samples a small number of negative objects to enhance the influence of positive objects. For a given node  $E$ , HHNE aims to maximize the proximity between  $E$  and its neighborhood  $\mathcal{C}$  while minimizes the proximity between  $E$  and its negative sampled nodes. Therefore, the objective function Eq. (6.27) can be rewritten as follows:

$$L^1 = \log f \approx \sum_{c \in \mathcal{C}} \log f_{\theta}(E, c) - \sum_{c \in \mathcal{C}^+} \log f_{\theta}(E, c) \quad (6.28)$$

where  $\mathcal{C}^+$  is a pre-defined distribution from which a negative node  $c$  is drawn from for  $n$  times. HHNE builds the node frequency distribution by drawing nodes regardless of their types.

### 6.4.2.3 Optimization

As the parameters of the model live in a Poincaré ball which has a Riemannian manifold structure, the back-propagated gradient is a Riemannian gradient. It means that the Euclidean gradient based optimization, such as  $\nabla_{\theta} L^1$ , makes no sense as an operation in the Poincaré ball, because the addition operation is not defined in this manifold. Instead, HHNE can optimize Eq. (6.28) via a Riemannian stochastic gradient descent (RSGD) optimization method [1]. In particular,  $T_{\theta} \mathcal{B}^3$  denote the tangent space of a node embedding  $\theta$  in  $\mathcal{B}^3$ , and HHNE can compute the Riemannian gradient  $\nabla_{\theta} L^1$  of  $L^1$ . Using RSGD, HHNE can be optimized by maximizing Eq. (6.28), and a node embedding can be updated in the form of:

$$\theta \leftarrow \exp_{\theta}^{-1}[\nabla_{\theta} L^1] \quad (6.29)$$

where  $\exp_{\theta}^{-1}$  is exponential map in the Poincaré ball. The exponential map is given by [11]:

$$\exp_{\theta}^{-1} \mathbf{v} = \frac{\frac{1}{\|\mathbf{v}\|} \cosh^{-1} \left( \frac{\|\mathbf{v}\|}{\|\mathbf{v}\|_{\theta}} \right) \mathbf{v}}{\frac{1}{\|\mathbf{v}\|} \cosh^{-1} \left( \frac{\|\mathbf{v}\|}{\|\mathbf{v}\|_{\theta}} \right) + \frac{\|\mathbf{v}\|_{\theta}}{\|\mathbf{v}\|} \sinh^{-1} \left( \frac{\|\mathbf{v}\|}{\|\mathbf{v}\|_{\theta}} \right)} \quad (6.30)$$

As the Poincaré ball model is a conformal model of hyperbolic spaces  $\mathbb{H}^3$ ,  $\mathbb{H}^2$ , the Riemannian gradient is obtained by rescaling the Euclidean gradient  $\nabla_{\theta} L^1$  by the inverse of the metric tensor, i.e.:

$$\nabla_{\theta} L^1 = \frac{1}{\|\mathbf{v}\|} \nabla_{\theta} L^1 \quad (6.31)$$

Furthermore, the gradients of Eq. (6.28) can be derived as follows:

$$\frac{mL}{m \setminus_{D^<}} = \frac{4}{U \frac{W^2}{1}} \mathbb{1}_{E \setminus D^<} \mathcal{F}^1 \mathcal{Z}_{D^<} \setminus_{D^<} \quad (6.32)$$

$$\frac{\setminus_{2_c}}{V_{<}} = \frac{k \setminus_{2_c} k^2}{U^2} \frac{2h \setminus_{2_c} \setminus_{D^<} \setminus_{D^<} \setminus_{D^<}}{\setminus_{D^<}} \setminus_{D^<}$$

$$\frac{mL}{m \setminus_{2_c}} = \frac{\tilde{\Theta}}{V_{<} \frac{W^2}{1}} \mathbb{1}_{E \setminus D^<} \mathcal{F}^1 \mathcal{Z}_{D^<} \setminus_{D^<} \quad (6.33)$$

$$\frac{\setminus_{D^<}}{U} = \frac{k \setminus_{D^<} k^2}{U^2} \frac{2h \setminus_{2_c} \setminus_{D^<} \setminus_{D^<} \setminus_{D^<}}{\setminus_{D^<}} \setminus_{2_c}$$

where  $U = 1 - k \setminus_{2_c} k^2$ ,  $V_{<} = 1 - k \setminus_{D^<} k^2$ ,  $W = 1 - \frac{2}{UV} k \setminus_{2_c} \setminus_{D^<} k^2$  and when  $< = 0$ ,  $D^0 = E$ .  $\mathbb{1}_{E \setminus D^<}$  is an indicator function to indicate whether  $D$  is  $E$ . Then, HHNE can be updated by using Eq. (6.32)-(6.33) iteratively.

### 6.4.3 Experiments

#### 6.4.3.1 Experimental Setup

**Datasets.** The basic statistics of the two HGs used in our experiments are shown in Table 6.6.

**Table 6.6** Statistics of datasets.

|           |       |       |      |       |       |
|-----------|-------|-------|------|-------|-------|
| DBLP      | # A   | # P   | # V  | # P-A | # P-V |
|           | 14475 | 14376 | 20   | 41794 | 14376 |
| MovieLens | # A   | # M   | # D  | # M-A | # M-D |
|           | 11718 | 9160  | 3510 | 64051 | 9160  |

**Baselines.** HHNE is compared with the following state-of-the-art methods: (1) the homogeneous graph embedding methods, i.e., DeepWalk [25], LINE [28] and node2vec [13]; (2) the heterogeneous graph embedding methods, i.e., metapath2vec [8]; (3) the hyperbolic homogeneous graph embedding methods, i.e., PoincaréEmb [22].

**Parameter Settings.** For random walk based methods DeepWalk, node2vec, metapath2vec and HHNE, we set neighborhood size as 5, walk length as 80, walks per node as 40. For LINE, metapath2vec, PoincaréEmb and HHNE, we set the number of negative samples as 10. For methods based on meta-path guided random walks, we use “APA” for relation “P-A” in network reconstruction and link prediction experiments in DBLP; “APVPA” for relation “P-V” in above experiments in DBLP; “AMDMA”

for all relation in above experiments in MovieLens. In visualization experiment, in order to focus on analyzing the relation of “A” and “P”, we use “APA”.

### 6.4.3.2 Network Reconstruction

**Table 6.7** AUC scores for network reconstruction.

| Dataset   | Edge   | Dimension | Deepwalk | LINE(1st) | LINE(2nd) | node2vec | metapath2vec  | PoincaréEmb   | HHNE          |
|-----------|--------|-----------|----------|-----------|-----------|----------|---------------|---------------|---------------|
| DBLP      | P-A    | 2         | 0.6933   | 0.5286    | 0.6740    | 0.7107   | 0.6686        | 0.8251        | <b>0.9835</b> |
|           |        | 5         | 0.8034   | 0.5397    | 0.7379    | 0.8162   | 0.8261        | 0.8769        | <b>0.9838</b> |
|           |        | 10        | 0.9324   | 0.6740    | 0.7541    | 0.9418   | 0.9202        | 0.8921        | <b>0.9887</b> |
|           |        | 15        | 0.9666   | 0.7220    | 0.7868    | 0.9719   | 0.9500        | 0.8989        | <b>0.9898</b> |
|           |        | 20        | 0.9722   | 0.7457    | 0.7600    | 0.9809   | 0.9623        | 0.9024        | <b>0.9913</b> |
|           | 25     | 0.9794    | 0.7668   | 0.7621    | 0.9881    | 0.9690   | 0.9034        | <b>0.9930</b> |               |
|           | P-V    | 2         | 0.7324   | 0.5182    | 0.6242    | 0.7595   | 0.7286        | 0.5718        | <b>0.8449</b> |
|           |        | 5         | 0.7906   | 0.5500    | 0.6349    | 0.8019   | 0.9072        | 0.5529        | <b>0.9984</b> |
|           |        | 10        | 0.8813   | 0.7070    | 0.6333    | 0.8922   | 0.9691        | 0.6271        | <b>0.9985</b> |
|           |        | 15        | 0.9353   | 0.7295    | 0.6343    | 0.9382   | 0.9840        | 0.6446        | <b>0.9985</b> |
|           |        | 20        | 0.9505   | 0.7369    | 0.6444    | 0.9524   | 0.9879        | 0.6600        | <b>0.9985</b> |
|           | 25     | 0.9558    | 0.7436   | 0.6440    | 0.9596    | 0.9899   | 0.6760        | <b>0.9985</b> |               |
|           | M-A    | 2         | 0.6320   | 0.5424    | 0.6378    | 0.6402   | 0.6404        | 0.5231        | <b>0.8832</b> |
|           |        | 5         | 0.6763   | 0.5675    | 0.7047    | 0.6774   | 0.6578        | 0.5317        | <b>0.9168</b> |
|           |        | 10        | 0.7610   | 0.6202    | 0.7739    | 0.7653   | 0.7231        | 0.5404        | <b>0.9211</b> |
| 15        |        | 0.8244    | 0.6593   | 0.7955    | 0.8304    | 0.7793   | 0.5479        | <b>0.9221</b> |               |
| 20        |        | 0.8666    | 0.6925   | 0.8065    | 0.8742    | 0.8189   | 0.5522        | <b>0.9239</b> |               |
| 25        | 0.8963 | 0.7251    | 0.8123   | 0.9035    | 0.8483    | 0.5545   | <b>0.9233</b> |               |               |
| MoiveLens | M-A    | 2         | 0.6626   | 0.5386    | 0.6016    | 0.6707   | 0.6589        | 0.6213        | <b>0.9952</b> |
|           |        | 5         | 0.7263   | 0.5839    | 0.6521    | 0.7283   | 0.7230        | 0.7266        | <b>0.9968</b> |
|           |        | 10        | 0.8246   | 0.6114    | 0.6969    | 0.8308   | 0.8063        | 0.7397        | <b>0.9975</b> |
|           | M-D    | 15        | 0.8784   | 0.6421    | 0.7112    | 0.8867   | 0.8455        | 0.7378        | <b>0.9972</b> |
|           |        | 20        | 0.9117   | 0.6748    | 0.7503    | 0.9186   | 0.8656        | 0.7423        | <b>0.9982</b> |
|           |        | 25        | 0.9345   | 0.7012    | 0.7642    | 0.9402   | 0.8800        | 0.7437        | <b>0.9992</b> |

A good HG embedding method should ensure that the learned embeddings can preserve the original HG structure. The reconstruction error in relation to the embedding dimension is then a measure for the capacity of the model. More specifically, we use network embedding methods to learn feature representations. Then for each type of links in the HG, we enumerate all pairs of objects that can be connected by such a link and calculate their proximity [18], i.e., the distance in Poincaré ball model for HHNE and PoincaréEmb. Finally, we use the AUC [9] to evaluate the performance of each embedding method. For example, for link type “write”, we calculate all pairs of authors and papers in DBLP and compute the proximity for each pair. Then using the links between authors and papers in real DBLP network as ground-truth, we compute the AUC value for each embedding method.

The results are shown in Table 6.7. As we can see, HHNE consistently performs the best in all the tested HGs. The results demonstrate that HHNE can effectively preserve the original network structure and reconstruct the network, especially on the reconstruction of P-V and M-D edges. Also, please note that HHNE achieves very promising results when the embedding dimension is very small. This suggests that regarding hyperbolic spaces underlying HG is reasonable and hyperbolic spaces have strong ability of modeling network when the dimension of spaces is small.

**Table 6.8** AUC scores for link prediction.

| Dataset   | Edge Dimension | Deepwalk | LINE(1st) | LINE(2nd) | node2vec | metapath2vec | PoincaréEmb   | HHNE          |               |
|-----------|----------------|----------|-----------|-----------|----------|--------------|---------------|---------------|---------------|
| DBLP      | P-A            | 2        | 0.5813    | 0.5090    | 0.5909   | 0.6709       | 0.6536        | 0.6742        | <b>0.8777</b> |
|           |                | 5        | 0.7370    | 0.5168    | 0.6351   | 0.7527       | 0.7294        | 0.7381        | <b>0.9041</b> |
|           |                | 10       | 0.8250    | 0.5427    | 0.6510   | 0.8469       | 0.8279        | 0.7699        | <b>0.9111</b> |
|           |                | 15       | 0.8664    | 0.5631    | 0.6582   | 0.8881       | 0.8606        | 0.7743        | <b>0.9111</b> |
|           |                | 20       | 0.8807    | 0.5742    | 0.6644   | 0.9037       | 0.8740        | 0.7806        | <b>0.9106</b> |
|           | 25             | 0.8878   | 0.5857    | 0.6782    | 0.9102   | 0.8803       | 0.7830        | <b>0.9117</b> |               |
|           | P-V            | 2        | 0.7075    | 0.5160    | 0.5121   | 0.7369       | 0.7059        | 0.8257        | <b>0.9331</b> |
|           |                | 5        | 0.7197    | 0.5663    | 0.5216   | 0.7286       | 0.8516        | 0.8878        | <b>0.9409</b> |
|           |                | 10       | 0.7292    | 0.5873    | 0.5332   | 0.7481       | 0.9248        | 0.9113        | <b>0.9619</b> |
|           |                | 15       | 0.7325    | 0.5896    | 0.5425   | 0.7583       | 0.9414        | 0.9142        | <b>0.9625</b> |
|           |                | 20       | 0.7522    | 0.5891    | 0.5492   | 0.7674       | 0.9504        | 0.9185        | <b>0.9620</b> |
|           | 25             | 0.7640   | 0.5846    | 0.5512    | 0.7758   | 0.9536       | 0.9192        | <b>0.9612</b> |               |
|           | M-A            | 2        | 0.6278    | 0.5053    | 0.5712   | 0.6349       | 0.6168        | 0.5535        | <b>0.7715</b> |
|           |                | 5        | 0.6353    | 0.5636    | 0.5874   | 0.6402       | 0.6212        | 0.5779        | <b>0.8255</b> |
|           |                | 10       | 0.6680    | 0.5914    | 0.6361   | 0.6700       | 0.6332        | 0.5984        | <b>0.8312</b> |
| 15        |                | 0.6791   | 0.6184    | 0.6442    | 0.6814   | 0.6382       | 0.5916        | <b>0.8319</b> |               |
| 20        |                | 0.6868   | 0.6202    | 0.6596    | 0.6910   | 0.6453       | 0.5988        | <b>0.8318</b> |               |
| 25        | 0.6890         | 0.6256   | 0.6700    | 0.6977    | 0.6508   | 0.5995       | <b>0.8309</b> |               |               |
| MoiveLens | M-D            | 2        | 0.6258    | 0.5139    | 0.6501   | 0.6299       | 0.6191        | 0.5856        | <b>0.8520</b> |
|           |                | 5        | 0.6482    | 0.5496    | 0.6607   | 0.6589       | 0.6332        | 0.6290        | <b>0.8967</b> |
|           |                | 10       | 0.6976    | 0.5885    | 0.7499   | 0.7034       | 0.6687        | 0.6518        | <b>0.8984</b> |
|           |                | 15       | 0.7163    | 0.6647    | 0.7756   | 0.7241       | 0.6702        | 0.6715        | <b>0.9007</b> |
|           |                | 20       | 0.7324    | 0.6742    | 0.7982   | 0.7412       | 0.6746        | 0.6821        | <b>0.9000</b> |
| 25        | 0.7446         | 0.6957   | 0.8051    | 0.7523    | 0.6712   | 0.6864       | <b>0.9018</b> |               |               |

### 6.4.3.3 Link Prediction

Link prediction aims to infer the unknown links in an HG given the observed HG structure, which can be used to test the generalization performance of a network embedding method. The experimental setting is similar to [34]. For each type of edge, 20% of edges are removed randomly from the network while ensuring that the rest network structure is still connected. The proximity of all pair of nodes are calculated in the test. AUC is used as the evaluation metric.

From the results in Table 6.8, HHNE outperforms the baselines upon all the dimensionality, especially in the low dimensionality. The results can demonstrate the generalization ability of HHNE. In DBLP dataset, the results of HHNE in 10 dimensionality exceed all the baselines in higher dimensionality results. In MovieLens dataset, HHNE with only 2 dimensionality surpasses baselines in all dimensionality. Besides, both of LINE(1st) and PoincaréEmb preserve proximities of node pairs linked by an edge, while LINE(1st) embed network into Euclidean spaces and Poincaré embed network into hyperbolic spaces. PoincaréEmb performs better than LINE(1st) in most cases, especially in dimensionality lower than 10, suggesting the superiority of embedding network into hyperbolic spaces. Because HHNE can preserve high-order network structure and handle different types of nodes in HG, HHNE is more effective than PoincaréEmb.

More detailed introduction of HHNE can be found in [33].

## 6.5 Conclusion

In this chapter, we have introduced three emerging HG embedding topics, as well as the related HG embedding methods. For adversarial learning, HeGAN is designed for learning semantic-preserving and robust HG embedding based on the adversarial principle. Besides, for importance sampling, HeteSamp studies the problem of accelerating large-scale HG embedding with importance sampling. Moreover, for hyperbolic representation, HHNE makes the efforts to embed HGs in hyperbolic spaces, and the optimization strategies are derived to optimize the hyperbolic embedding. We hope more HG embedding methods with deeper insights can be proposed to discover the abundant semantic of HG in the future.

## References

1. Bonnabel, S., et al.: Stochastic gradient descent on riemannian manifolds. *IEEE Trans. Automat. Contr.* **58**(9), 2217–2229 (2013)
2. Cai, X., Han, J., Yang, L.: Generative adversarial network based heterogeneous bibliographic network representation for personalized citation recommendation. In: *AAAI*, pp. 5747–5754 (2018)
3. Cannon, J.W., Floyd, W.J., Kenyon, R., Parry, W.R., et al.: Hyperbolic geometry. *Flavors of geometry* **31**, 59–115 (1997)
4. Cen, Y., Zou, X., Zhang, J., Yang, H., Zhou, J., Tang, J.: Representation learning for attributed multiplex heterogeneous network. In: *KDD*, pp. 1358–1368 (2019)
5. Chen, J., Ma, T., Xiao, C.: Fastgcn: Fast learning with graph convolutional networks via importance sampling. In: *ICLR* (2018)
6. Dai, Q., Li, Q., Tang, J., Wang, D.: Adversarial network embedding. In: *AAAI*, pp. 2167–2174 (2018)
7. Dhingra, B., Shallue, C., Norouzi, M., Dai, A., Dahl, G.: Embedding text in hyperbolic spaces. In: *TextGraphs*, pp. 59–69 (2018)
8. Dong, Y., Chawla, N.V., Swami, A.: metapath2vec: Scalable representation learning for heterogeneous networks. In: *KDD*, pp. 135–144 (2017)
9. Fawcett, T.: An introduction to roc analysis. *Pattern recognition letters* **27**(8), 861–874 (2006)
10. Fu, T.y., Lee, W.C., Lei, Z.: Hin2vec: Explore meta-paths in heterogeneous information networks for representation learning. In: *CIKM*, pp. 1797–1806 (2017)
11. Ganea, O., Becigneul, G., Hofmann, T.: Hyperbolic entailment cones for learning hierarchical embeddings. In: *ICML*, pp. 1646–1655 (2018)
12. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: *NeurIPS*, pp. 2672–2680 (2014)
13. Grover, A., Leskovec, J.: node2vec: Scalable feature learning for networks. In: *KDD*, pp. 855–864 (2016)
14. Hamilton, W.L., Ying, Z., Leskovec, J.: Inductive representation learning on large graphs. In: *NeurIPS*, pp. 1025–1035 (2017)
15. He, X., Liao, L., Zhang, H., Nie, L., Hu, X., Chua, T.S.: Neural collaborative filtering. In: *WWW*, pp. 173–182 (2017)
16. Hu, B., Fang, Y., Shi, C.: Adversarial learning on heterogeneous information networks. In: A. Teredesai, V. Kumar, Y. Li, R. Rosales, E. Terzi, G. Karypis (eds.) *KDD*, pp. 120–129 (2019)
17. Huang, W., Zhang, T., Rong, Y., Huang, J.: Adaptive sampling towards fast graph representation learning. In: *NeurIPS*, pp. 4563–4572 (2018)

