# **Prohibited Item Detection on Heterogeneous Risk Graphs**

Yugang Ji Beijing University of Posts and Telecommunications Beijing, China jiyugang@bupt.edu.cn Chuan Shi\* Beijing University of Posts and Telecommunications Beijing, China shichuan@bupt.edu.cn

# ABSTRACT

Prohibited item detection, which aims to detect illegal items hidden on e-commerce platforms, plays a significant role in evading risks and preventing crimes for online shopping. While traditional solutions usually focus on mining evidence from independent items, they cannot effectively utilize the rich structural relevance among different items. A naïve idea is to directly deploy existing supervised graph neural networks to learn node representations for item classification. However, the very few manually labeled items with various risk patterns introduce two essential challenges: (1) How to enhance the representations of enormous unlabeled items? (2) How to enrich the supervised information in this few-labeled but multiple-pattern business scenario? In this paper, we construct item logs as a Heterogeneous Risk Graph (HRG), and propose the novel Heterogeneous Self-supervised Prohibited item Detection model (HSPD) to overcome these challenges. HSPD first designs the heterogeneous self-supervised learning model, which treats multiple semantics as the supervision to enhance item representations. Then, it presents the directed pairwise labeling to learn the distance from candidates to their most relevant prohibited seeds, which tackles the binary-labeled multi-patterned risks. Finally, HSPD integrates with self-training mechanisms to iteratively expand confident pseudo labels for enriching supervision. HSPD has been deployed on Taobao platform, and the extensive offline and online experimental results on three real-world HRGs demonstrate that HSPD consistently outperforms the state-of-the-art alternatives.

# CCS CONCEPTS

 Computing methodologies → Machine learning; Artificial intelligence;
 Information systems → World Wide Web.

# **KEYWORDS**

prohibited item detection; graph neural network; heterogeneous self-supervised learning; self-training

#### ACM Reference Format:

Yugang Ji, Chuan Shi, and Xiao Wang. 2021. Prohibited Item Detection on Heterogeneous Risk Graphs. In Proceedings of the 30th ACM Int'l Conf. on

CIKM '21, November 1-5, 2021, Virtual Event, Australia.

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8446-9/21/11...\$15.00

https://doi.org/10.1145/3459637.3481945



Xiao Wang

Beijing University of Posts and

Telecommunications

Beijing, China

xiaowang@bupt.edu.cn

Figure 1: The typical workflow of prohibited item detection on the e-commerce.

Information and Knowledge Management (CIKM '21), November 1–5, 2021, Virtual Event, Australia. ACM, New York, NY, USA, 11 pages. https://doi. org/10.1145/3459637.3481945

# **1 INTRODUCTION**

In the era of prosperous e-commerce, online shopping has become popular for its advantages in convenience, detailed information as well as competitive prices. However, besides numerous interesting products, as shown in Figure 1, there are various items against laws hidden on the platforms, including protected wildlife, pornographic materials, illegal medicines, and many others. For example, there have been more than 1 million products claiming to cure coronavirus trying to sell on Amazon since COVID-19<sup>1</sup>. Meanwhile, 1.35 million listings of wildlife products attacked Taobao<sup>2</sup> in 2019. The selling of these illegal items would bring huge risks to platforms like expansive fines and create salient personal and social issues such as increasing crime rate and rampant poaching.

Prohibited item detection, which aims at searching and deleting illegal items hidden on e-commerce platforms, has played an essential and vital role in evading risks and preventing crimes for online shopping [7, 28]. Conventional industrial solutions prefer to formulate this problem as a typical classification task and directly deploy traditional machine learning or deep learning algorithms [4] to extract confident evidences from independent instances. Obviously, these solutions require laboring feature engineering and adequate supervised information of item logs. However, prohibited item detection work becomes harder and harder in recent years, because of the following two reasons.

First, the attributes of instances are weak due to the adversarial actions. Since unstructured features (e.g., texts and images) of prohibited items can be easily transformed very similar to those of normal ones, traditional feature engineering suffers from heavy

<sup>\*</sup>The corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

<sup>&</sup>lt;sup>1</sup>https://www.bbc.com/news/technology-51675183

<sup>&</sup>lt;sup>2</sup>https://www.endwildlifetraffickingonline.org/our-progress

adversarial efforts by illegal sellers. Meanwhile, current solutions rarely utilize rich structural relations between different items which indeed reflect the high risks in someway. For instance, items belonging to the same seller are likely to share consistent risks. Modeling such structure information can learn more robust representations of items [1, 12?] and make the adversarial efforts unaffordable or unavailable [8]. Second, the supervised information for training is weak in both the size and the expression of labels, limited by manpower costs. There are very few labeled items. Besides, the labels can only describe the riskiness while neglecting the diverse patterns of prohibited items. For example, the risk patterns in pornographic products include porn sounds and videos, adult novels, online private services and website accounts, etc, while illegal medicines can be divided into prescription drugs, fake medicines, illegal aphrodisiacs and so on. Obviously, nodewise labels cannot make full use of the diversity of same-labeled but different-pattern items, leading to incomplete training and overfitting.

To tackle the first problem (i.e., weak attributes), a promising direction is to model the rich structures as a graph and learn robust representations of items via advanced graph neural networks [6, 17, 31]. Aware of both the risk and heterogeneity, in this paper, we design a Heterogeneous Risk Graph (HRG) to connect different items via multiple risk-relevant relations. Taking Figure 2 as an example, there are three kinds of high-level semantics including "same visitor", "same seller", and "relevant seller" between items on Taobao platform. Focusing on representation learning on HRGs, a naïve idea is to introduce heterogeneous graph neural networks[10, 13, 33] to make use of both heterogeneous structures and attributed information. Unfortunately, these semi-supervised models often require plenty of supervised information, however, there are very few labeled but too many unlabeled instances on HRGs. Recently, self-supervised learning has been introduced into graphs [14, 18, 29, 30] to learn robust representation of unlabeled objects for training. However, these methods mainly deal with homogeneous graphs but fail to preserve abundant semantics of HRGs. This remains the first challenge, namely, how to enhance item representation by fully exploiting abundant semantics on HRGs?

Focusing on the second problem (i.e., weak supervision), it has become popular to integrate with self-training mechanism [18, 41] to expand the scale of labels. Given labeled and unlabeled instances, a typical self-training pipeline make up of three steps, (1) pretrain a model over labeled instances, (2) assigning "pseudo-labels" to highly confident unlabeled instances, (3) include these pseudo labels into the labeled set for next round of training. Existing works [18, 30] are to expand nodewise supervised information. However, for prohibited item detection, as the same-class items may be very different, the performances of both pretraining and labeling processes are limited. In other words, current nodewise self-training mechanisms just learn label-level similarity while neglecting the pattern-level relevance in real-world scenarios. The easy negative and hard positive supervision would lead to overfitting [36]. Thereby, the second challenge is how to enrich supervised information in this few labeled but various patterned business scenario?

In this paper, we are the first to introduce heterogeneous graph to model risk-relevant structures of item logs, and propose a novel Heterogeneous Self-supervised framework for Prohibited item Detection (HSPD). In this model, we treat the semantics between items as self supervision and design an effective heterogeneous self-supervised learning on HRGs which factorizes and disentangles semantics within relations to enhance robust representations of items. And then, we transform the detection process as a metric learning task between items to be predicted and existing prohibited ones, making full use of various pattern-level relevance. We further design the directed pairwise self-training mechanism to iteratively generate more supervision to improve the generalization performance. Obviously, HSPD can be widely used in many realworld scenarios where objects are rarely labeled or objects belong to different patterns or classes.

In a nutshell, the contributions of this paper are:

- The problem of prohibited item detection is very significant to prevent crimes and protect healthy online shopping. To our best knowledge, we are the first to introduce heterogeneous graph modeling to address this problem.
- We design the effective HSPD consisting of heterogeneous self-supervised learning and directed pairwise self-training, which can simultaneously enhance item representations and enrich supervised information to overcome the challenges of weak attributes and weak supervision during detecting.
- We evaluate our model in three industrial scenarios, including protected wildlife, illegal medicines and pornographic products. All experimental results consistently demonstrate the effectiveness of our designs and the improvements to the second best baseline in the AP and Max-F1 metrics are respectively up to 9.90% and 9.40%.

## 2 RELATED WORK

The related work includes the heterogeneous graph neural networks, the self-supervised learning on graphs and so on.

Heterogeneous graph neural networks. Recent years have witnessed the success of GNNs which have the ability to model graphstructured data, naturally capturing both graph structures and attributes on graphs [6, 17, 35]. GNNs usually generate contextual node representations via neighborhood aggregation. Under this framework, various GNN architectures have been proposed [5, 6, 17, 31, 34]. However, as item logs in real-word scenarios are often connected by multiple relations, these homogeneous GNNs fail to model the heterogeneity within such HRGs. Recently, some studies have attempted to deploy GNNs on heterogeneous graphs [10, 13, 26, 33]. RGCN [26] utilizes multiple linear projection weights for each edge type. HAN [33] and HGT [10] incorporate attention mechanisms into heterogeneous graphs and hierarchically aggregate information from different-typed neighborhoods. On Taobao platform, heterogeneous GNNs have been introduced to address various real-world tasks including recommendation [3, 20], user alignment [39] and so on. However, when detecting prohibited items, existing methods cannot thoroughly take advantage of the abundant information, because of too few labeled data.

**Self-supervised learning on graphs.** Self-supervised learning [19], which is a general learning framework that relies on pretext tasks that can be formulated using unlabeled data, has shown its advantages in graph modeling for the fantastic data efficiency and generalization ability. DGI [32] proposes to maximize the mutual



Figure 2: An toy example of HRG construction on Taobao platform.

information between the local node representation and the global graph context. GMI [25] further proposes to maximize the mutual information of both features and edges between inputs and outputs of the encoder. Focusing on the inherent structures and attributes, GraphCL [37] designs the contrastive learning with four types of graph augmentations (i.e., node dropping, edge perturbation, attribute masking and subgraphs), to enhance robust node representations. GPT-GNN [9] proposes to consider both the edge and node attributes as self-supervision and models the generation for pre-training GNNs. Recently, some works [18, 30] propose to adopt self-supervised learning to help downstream tasks, indicating that the self-supervision should be consistent with the supervised learning. While the above methods are designed to handle homogeneous networks, they fail to fully exploit the semantics within HRGs to enhance item representation. Some researches attempt to model the self-supervised learning on heterogeneous graphs. SE-LAR [11] treats meta-paths as the self-supervised information while still constructing node representation in a homogeneous manner. DMGI [23] introduces DGI [32] into attributed multiplex networks and treat types of edges as supervision, however, it focuses on modeling the global properties, which is unaffordable and likely useless in prohibited item detection because of the web-scale data.

**Other related works.** Self-training [41], which is to iteratively assign confident predictions as the the supervised information, has been proved effective in GNNs. M3S [30] evaluates the confidence by matching the labels and clusters, while Pedronette *et al.* [24] focus on the ranking information. However, these methods cannot directly introduced in HRGs because of neglecting the semantics. On the other line, due to the ability to learn the relevance (i.e., metric learning) of candidates to the existing classes, pairwise labeling [15] has been introduced in few-shot learning containing multiple classes with few labels.

#### **3 PRELIMINARIES**

This section introduces the general structure of item logs, the construction of heterogeneous risk graphs, and then formalizes the problem of prohibited item detection on HRGs. The key notations are shown in Table 1.

As shown in Figure 2(a), item log generally consists of not only the unstructured attributes (e.g., the title and images) but also several related objects like its visitor, its seller, its seller's IP address and Mobile number (MID) and etc, indicating the structural relevance between items. Obviously, we can directly connect items according

Table	1:	Not	ations
-------	----	-----	--------

Notation	Description
G	the input HRG
$\mathcal{V}, \mathcal{E}$	the item/relation set of $G$
$\mathcal R$	the relation type set of ${\cal E}$
$r,\psi\in\mathcal{R}$	the relation type / semantic of ${\cal R}$
$\mathcal{N}_{i}^{r}$	the type- $r$ neighbors of $v_i$
m	the sample size of neighborhoods
Ν	the number of labels
d	the dimension of attributes
Т	the number of self-training epochs
$X \in \mathbb{R}^d$	the attributes of items
$ ilde{m{h}}_i \in \mathbb{R}^{ \mathcal{R} d}$	the base representation of item $v_i$
$ ilde{m{h}}^r_i \in \mathbb{R}^d$	the type- <i>r</i> representation from $\mathcal{N}_i^r$
$H_{i,\psi} \in \mathbb{R}^d$	the disentangled semantic- $\psi$ factor of $v_i$
$H_i \in \mathbb{R}^{ \mathcal{R} d}$	the self-supervised representation of $v_i$
$I_{i,j,\psi} \in \{0,1\}$	the type- $\psi$ connection between $v_i$ and $v_j$
$Y_i, Y_{i,j} \in \{0, 1\}$	nodewise/ pairwise label of $v_i / v_i$ and $v_j$
$Y^t, Z^t \in \{0,1\}$	training / pseudo label set at $t^{th}$ self-training

to their same factors. However, some relations like "same-category" could hide the risk trace. In other words, prohibited items are often related to normal ones in these relations, leading to very noisy structures. Focusing on modeling rich structure information related to prohibited item detection, taking Figure 2(b) as an example, we empirically connect different items via three kinds of relations, namely, (1) **Same seller** describes that both connected items belong to the same seller, capturing the risk from same sellers. (2) **Same visitor** describes that both source and target nodes (i.e., items) have been visited by some same visitors, indicating the relevant risk to consumers. (3) **Relevant seller** is to connect the items of relevant sellers to overcome the multiple fake identifications of adversarial sellers. In addition, there are some other risk-relevant relations and we choose the three representative relations for discussion.

DEFINITION 1. Heterogeneous Risk Graph (HRG). An HRG is denoted as  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}, X\}$  where  $\mathcal{V}$  is the item set,  $\mathcal{E}$  is the relation set among items and X denotes the attributes of nodes. There is a relation-type mapping function on HRGs, namely  $\psi : \mathcal{E} \to \mathcal{R}$  where  $\mathcal{R}$  denotes the relation types including "same-seller", "same-visitor", "relevant-seller" and some others where prohibited items are likely to connect with each other. Notice that, the relations on HRGs are in the form of meta-paths [27] and meta-graphs [38] to describe risk-relevant semantics, rather than general edges on heterogeneous graphs.

As shown in Figure 2(c), HRG is able to preserve the risk-relevant semantics within item logs for detecting prohibited items rather than keeping all connections.

DEFINITION 2. Prohibited item detection on HRGs. Given an HRG  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}, X\}$ , label set  $\mathcal{Y} = \{0, 1\}$  as well as the supervised information  $Y \in \mathcal{Y}^N$  where N denotes the size of manually labeled items and  $N \ll |\mathcal{V}|$ , the goal of prohibited item detection is to learn the judgement function  $\mathcal{H} : \mathcal{H}(\mathcal{V}) \to \mathcal{Y}^{|\mathcal{V}|}$ . Notice that the values 0 and 1 in  $\mathcal{Y}$  are respectively to mark the normal and prohibited items.

This problem contains two major characteristics, making it more challenging than traditional node classification. On the one hand, there are too few labeled items (e.g., the ratio  $\frac{N}{|V|}$  is 0.05% in our online scenario) to ensure the robust representation of items. On the other hand, same-labeled items may belong to different patterns (e.g., coats and bags in Figure 2(c) are all prohibited but they are quite different to each other.), implying that nodewise labeling and training cannot make full use of supervised information.

## 4 METHODOLOGY

In this section, we present the HSPD to tackle prohibited item detection on HRGs, making full use of both structures and attributes. We begin with an overview, before zooming into the details.

#### 4.1 Overview

The system architecture of HSPD in Figure 3 consists of three main modules, including HRG construction, heterogeneous selfsupervised learning for enhancing item representations, and directed pairwise self-training for enriching supervised information to handle prohibited item detection and overcome the two challenges. Specifically, (1) we first construct the numerous item logs as an HRG which connects items via multiple risk-relevant relations. (2) To enhance robust representations of unlabeled items, HSPD treats the semantics of relations as the self-supervised information, and propose the heterogeneous self-supervised learning which disentangles semantic-level factors of relations as the robust representations of items. (3) To enrich supervised information, HSPD generates the directed pairwise labels to model the pattern-level relevance of prohibited items via metric learning, and deploy selftraining framework to expand confident pseudo labels to improve the generalization. Notice that, as the details of HRG modeling have been described in Section 3, we mainly introduce the representation enhancing and supervision enhancing in this section.

# 4.2 Representation Enhancing via Heterogeneous Self-Supervised Learning

We begin with the heterogeneous self-supervised learning for enhancing robust representations of enormous unlabeled items. Although self-supervised learning has been introduced to deal with few-labeled graphs, however, existing works are on homogeneous graphs which fail to incorporate the abundant semantics when generating node representations in HRGs. In order to fully preserve the rich semantics on HRGs, we treat the types of relations, i.e.,



Figure 3: System architecture of HSPD.

risk-relevant semantics between items as the self-supervision and present the self-supervised training on heterogeneous graphs, as shown in Figure 4.

At first, given the HRG G as well as item  $v_i$ , we respectively gather information from type-aware neighborhoods, namely,

$$\hat{h}_{i}^{r} = AGGREGATE\left(\left\{\sigma(X_{j}W^{r} + b^{r}) \middle| j \in \mathcal{N}_{i}^{r}\right\}\right), \tag{1}$$

where  $\mathbf{\tilde{h}}_{i}^{r}$  denotes the gathered neighborhood information of  $\mathcal{N}_{i}^{r}$ ,  $\mathcal{N}_{i}^{r}$  denotes the *r*-typed neighborhoods of item  $v_{i}$ ,  $X_{j}$  denotes the features of item *j*,  $\mathbf{W}^{r}$  and  $b^{r}$  are the type-wise parameter and bias to model the semantics within relations,  $\sigma(\cdot)$  denotes the activation function and we adopt RELU in this paper.  $AGGREGATE(\cdot)$  is the pooling operation for neighborhoods where we deploy the mean-pooling to keep the full information.

Thereby, the base representation  $\tilde{h}_i$  of item  $v_i$  is generated by gathering information from multiple neighborhoods via heterogeneous aggregator  $HeteAGG(\cdot)$ , namely,

$$\tilde{\boldsymbol{h}}_i = HeteAGG(\{\boldsymbol{X}_j | j \in \mathcal{N}_i^r, r \in \mathcal{R}\}) = [\tilde{\boldsymbol{h}}_i^{r_1} \| \tilde{\boldsymbol{h}}_i^{r_2} \| \cdots \| \tilde{\boldsymbol{h}}_i^{r_{|\mathcal{R}|}}].$$
(2)

With the assumption that all heterogeneous neighborhoods contain latent factors to result in current connections at different levels, inspired by [21], we factorize and disentangle these factors by designing the semantic-aware self-attention mechanism, namely,

$$H_{i,\psi} = \alpha_{i,\psi} h_i W_{\psi} + X_i, \tag{3}$$

where  $H_{i,\psi} \in \mathbb{R}^d$  with dimension *d* denotes the latent factor of type- $\psi$  semantics,  $\tilde{h}_i \in \mathbb{R}^{|\mathcal{R}| \times d}$  is the embedding generated by Eq. (2), and  $W_{\psi} \in \mathbb{R}^{d \times d}$  is the semantic-aware projection parameter,  $\alpha_{i,\psi} \in \mathbb{R}^{1 \times |\mathcal{R}|}$  denotes the corresponding importance of the multiple neighbors, defined as

$$\boldsymbol{\alpha}_{i,\psi} = \operatorname{softmax} \left( \tanh(\tilde{\boldsymbol{h}}_i \boldsymbol{W}_{\psi,\alpha}) \boldsymbol{w}_{\psi,\alpha} \right), \tag{4}$$

where  $W_{\psi,\alpha} \in \mathbb{R}^{d \times d_{\psi}}$  and  $w_{\psi,\alpha} \in \mathbb{R}^{d_{\psi} \times 1}$  are two projection parameters in self-attention, and we adopt softmax(·) to normalize the importance of these multiple information.



Figure 4: Enhancing representation of items via heterogeneous self-supervised learning.

Furthermore, we treat the heterogeneous relations between items as self-supervised information and focus on prediction the semantics between items. We randomly mask several relations between neighbors as the positive links and adopt negative sampling to generate the corresponding unconnected pairs. The sizes of unconnected and connected pairs are the same. Given the semantic-aware pair < i, j,  $\psi$  >, the heterogeneous self-supervised cross-entropy loss is defined as

$$\mathcal{L}_{SS} = -\sum_{\langle i,j,\psi \rangle} I_{i,j,\psi} \log(p_{i,j,\psi}) + (1 - I_{i,j,\psi}) \log(1 - p_{i,j,\psi}), \quad (5)$$

where  $I_{i,j,\psi} \in \{0,1\}$  denotes the connection of  $v_i$  and  $v_j$  under type- $\psi$  semantic, and  $p_{i,j,\psi}$  denotes the probability, namely,

$$p_{i,j,\psi} = \boldsymbol{H}_{i,\psi}^T \boldsymbol{H}_{j,\psi}.$$
 (6)

Notice that, in this paper, we concatenate all disentangled factors as the robust representation  $H_i^T = \|_{\psi} H_{i,\psi}^T$  for supervised learning.

# 4.3 Supervision Enriching via Directed Pairwise Self-Training

Besides enhancing the inputs, another alternative is to enrich supervised information for better generalization. As shown in Figure 5, taking the pattern-level relevance into consideration, we present the directed pairwise labeling to evaluate the relevance from candidate items to their related prohibited items, transforming node classification into metric learning. These connections are indeed easy positive or hard negative training pairs, making the metric learning effective. Furthermore, we introduce the corresponding pairwise self-training mechanism to expand confident pseudo labels.

4.3.1 Directed Pairwise Labeling. With the robust inputs, a natural idea is to consider the detection as a binary classification task and train an effective supervised model to judge candidate items. However, as mentioned in Section 1, despite the binary *Y*, there are various patterns of prohibited items and these patterns are unique in some way, e.g., porn sounds and videos, adult novels, online private services and website accounts in pornographic products. Besides, the normal items could be very different from each other. Therefore, traditional approaches [18, 30] in node classification have to suffer from poor generalization.

Inspired by metric learning [15, 36], we propose to transform the nodewise classification as the proximity between items. Different from traditional edge labeling which treats the source and target nodes equally, the relevance between normal items contributes



Figure 5: Enriching supervised information via directed pairwise self-training.

little to risk detection. Thereby, we design the directed pairwise labeling in the form of  $\{ \langle v_i, v_j \rangle | Y_i = 1 \}$  as follows.

- Collect all the prohibited items  $\{v_i | v_i \in Y, Y_i = 1\}$  as seeds;
- Collect all the labeled items  $\{v_j | v_j \in Y\}$  as candidates;
- Select the pair < v<sub>i</sub>, v<sub>j</sub> > where v<sub>j</sub> is connected to v<sub>i</sub> or v<sub>i</sub> is the several most similar ones to v<sub>j</sub>;
- Label the pairs  $\langle v_i, v_j \rangle$  with  $Y_{i,j} = Y_j$ .

Thereby, the loss with directed pairwise labels is defined as

$$\mathcal{L}_{PW} = -\sum_{\langle i,j \rangle} Y_{i,j} \cdot \log(\hat{Y}_{i,j}) + (1 - Y_{i,j})\log(1 - \hat{Y}_{i,j}), \quad (7)$$

where  $\hat{Y}_{i,j}$  denotes the risk probability of  $v_i$ , calculated by

$$\hat{Y}_{i,j} = MLP([\boldsymbol{g}_i \| \boldsymbol{g}_j \| \boldsymbol{g}_i \cdot \boldsymbol{g}_j]), \tag{8}$$

where MLP(·) denotes the Multiple Layer Perceptron which outputs the link probability,  $g_i \in \mathbb{R}^d$  denotes the embedding of  $v_i$ , namely,

$$g_i = HeteAGG(H_o | o \in \mathcal{N}_i^{\psi}), \tag{9}$$

where  $H_o \in \mathbb{R}^{d|\mathcal{R}|}$  denotes the concatenation of all latent factors generated by Eq. (3). Notice that, since each item can be assigned with a certain number of pairs, the directed pairwise labeling contributes to expanding supervision, as well.

4.3.2 Pairwise Self-Training Strategy. Different from metric learning which generates extra information from labels themselves, selftraining is another strong alternative to deal with weak supervision. However, the current methods mainly focus on node-level self-training, which conflicts with our pairwise setting. Here we introduce self-training into pairwise supervised learning to generate labels via multiple pairwise instances rather than single nodes.

At first, we respectively generate n directed pairs of training set Y and prediction set Z, and then rewrite the loss in Eq. (7) as

$$\mathcal{L}_{ST}^{t} = -\sum_{\langle i,j \rangle} Y_{i,j}^{t} \cdot \log(\hat{Y}_{i,j}^{t}) + (1 - Y_{i,j}^{t})\log(1 - \hat{Y}_{i,j}^{t}), \quad (10)$$

where *t* denotes the epoch of self-training, and there are *T* total epochs,  $Y_{i,j}^t$  and  $\hat{Y}_{i,j}^t$  respectively denote the ground truth and the probability, and  $Y^0 = Y$ .

And then, for each epoch of self-training, we predict the probability for the candidate prediction set  $Z^t$  as follows

$$\hat{Z}_{i}^{t+1} = \frac{1}{k} \sum_{\langle i,j \rangle \in \mathbb{Z}^{t}} \hat{Z}_{i,j}^{t},$$
(11)

where  $\hat{Z}_{i,j}^{t+1}$  denotes the probability, namely the confidence to judge  $v_i$  as the prohibited item, and vice versa, and k denotes the size of

**Algorithm 1:** The proposed HSPD model.

**Input:** HRG  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}, X\}$ , labels *Y*, unlabeled items *Z*, directed pairs *k*, epoch *T*, threshold  $\kappa_0$  and  $\kappa_1$ ; **Output:** Optimized parameters  $\Theta$ ;

- 1 Initialize parameters  $\Theta$ ,  $t \leftarrow 0$ ;
- <sup>2</sup> Randomly sample several type-aware connected and unconnected quad set  $I = \{\langle i, j, \psi, I_{i,j,\psi} \rangle\};$
- 3 for each <  $i, j, \psi, I_{i,j,\psi}$  > do
- 4 Calculate base embedding  $\tilde{h}_i$  and  $\tilde{h}_j$  by Eq. (2);
- 5 Calculate disentangled  $\tilde{H}_{i,\psi}$  and  $\tilde{H}_{j,\psi}$  by Eq. (3);
- 6 Calculate the probability of connections  $p_{i,j,\psi}$  by Eq. (6); 7 end
- 8 Minimize  $\mathcal{L}_{SS}$  in Eq. (5) and obtain H;
- 9 Generate directed pairwise labels  $\{Y_{i,i}^0 | Y_i = 1\}$  from Y;
- 10 Generate candidate pairs  $\{Z_{i,i}^0 | Y_j = 1\}$  from Z and Y;
- 11 while t < T do
- <sup>12</sup> Calculate item robust representation  $g_i$  and  $g_j$  by Eq. (2) for all training pairs in *Y*;
- 13 Calculate supervised loss  $\mathcal{L}_{ST}^t$  by Eq. (10);
- 14 Minimize  $\mathcal{L}^t$  in Eq. (14) with fixed  $\mathcal{L}_{SS}$  by Adam;
- 15 Predict  $\hat{Z}^{t+1}$  by Eq. (11);
- 16 Update  $Y^{t+1}$  and  $Z^{t+1}$  by Eq. (12) and Eq. (13);
- 17  $t \leftarrow t + 1;$
- 18 end

pairs. Notice that,  $Z^0 = Z$ . We then rank the candidates and expand the confident predictions as the pseudo labels for training in the next epoch and remove them from the candidate set, namely,

$$Y^{t+1} = Y^t \cup \{ \tilde{Z}_{i,i}^{t+1} | \gamma_{i-}^{t+1} \le \kappa_1 \text{ or } \gamma_{i+}^{t+1} \le \kappa_0 \},$$
(12)

$$Z^{t+1} = \{ Z^t_{i,j} | \kappa_0 < \gamma^{t+1}_{i,+} < \kappa_1 \},$$
(13)

where  $\gamma_{i,:}^{t+1}$  denotes the ranking order of item  $v_i$  at  $t + 1^{th}$  epoch, the signs + and – under  $\gamma$  respectively indicate ascending and descending, and  $\kappa_0$  and  $\kappa_1$  are the borderline of normal items and risk items. Notice that, labels in risk scenarios are often imbalanced, here we set  $\kappa_0$  and  $\kappa_1$  differently, where  $\kappa_0 = 1000$  and  $\kappa_1 = 100$ . For the same reason, we do not require  $Z^T = \{\phi\}$  as too large pseudo labels are likely to accumulate very many errors.

## 4.4 The Unified Framework

By now, we have introduced both the heterogeneous self-supervised learning and the directed pairwise self-training mechanism to respectively enhance item representation and enrich supervised information. The overall optimized objective minimize both the self-supervised loss  $\mathcal{L}_{SS}$  and self-training loss  $\mathcal{L}_{ST}^{t}$ , defined as

$$\mathcal{L}^{t} = \mathcal{L}_{ST}^{t} + \beta \mathcal{L}_{SS} + \xi \Omega(\Theta), \qquad (14)$$

where  $\beta$  is the weight of self-supervised learning tasks,  $\xi$  denotes the regularization of all learnt parameters  $\Theta$  and  $\Omega$  denotes the L2 regularization. Notice that  $\mathcal{L}_{ST}^t = \mathcal{L}_{PW}$  if t = 0. Since the HRG  $\mathcal{G}$ is web-scale, we minimize the loss in two-step optimization and the details are in Algorithm 1.

#### 4.5 Complexity Analysis

The computational complexity of HSPD consists of two major parts, including heterogeneous self-supervised learning and directed pairwise self-training. For the former, the complexity is  $O(m|I||\mathcal{V}||\mathcal{R}|d_{att}d + |\mathcal{R}|^2d)$  where |I| and *m* respectively denote the the size of self-supervised labels and the size of each-typed neighborhood of each item,  $d_{att}$  and *d* respectively denote the dimension of attributes and outputs. For the latter, the complexity of  $t^{th}$  epoch is  $O(mk|Y^t||\mathcal{R}|^2d^3)$  where *k* denote the size of pairs of each labeled item,  $|Y^t|$  denotes the size of labeled items at this epoch, *d* denotes the size of output embedding of items. Obviously, both the two parts are linear with the scale of an HRG, demonstrating the scalability of our HSPD.

# **5 EXPERIMENTS**

In this section, we conduct experiments on three real-world risk scenarios to evaluate the empirical performance of our method, against seven state-of-the-art alternatives. And then, we perform both the variant and the parameter analysis to showcase the effectiveness of our design choices and key factors.

#### 5.1 Datasets

We collect the one-month real-world web-scale datasets in three risk scenarios including the protected wildlife (i.e., "*Wildlife*"), the illegal medicines (i.e., "*Medicine*") and the pornographic materials (i.e., "*Pornography*"), from the Taobao platform<sup>3</sup>. For each risk dataset, we empirically construct an HRG to preserve abundant semantics within billions of item logs. We adopt word2vec [22] to embed the titles of items as 64-dimensional numerical features.

Next, we introduce how to construct training, validation and test instances. For offline experiments, the instances are randomly divided into training, validation and test with rate 8:1:1. To get more robust results, we vary the size of each training sets from 20% to 80%. The detailed statistics of these datasets are described in Table 2. Besides offline experiments, we also evaluate the performance of our method by designing online testing of 9-day online dataset.

Obviously, there are several unique characteristics in our risk scenarios, compared to traditional binary classification. First, datasets are large enough and contain various relations between nodes. Second, the rate of manually labeled items to the whole instances is less than 0.2% while the balance rate between illegal (label = 1) and legal items (label = 0) is even very small to 2% in the risk of pornographic materials. These characteristics of data bring great challenges to our model designs.

#### 5.2 Experimental Settings

*5.2.1 Baselines.* We compare with seven representative baseline methods including *two* conventional classification algorithms (i.e., LR and GBDT [4]) which currently deployed for prohibited item detection on Taobao platform, *five* outstanding GNNs, as well as our nodewise HSPD (i.e., HSPD<sub>NW</sub>).

 Logistic Regression (LR). This is a fundamental classification algorithm used in industry for its good interpretability.

<sup>&</sup>lt;sup>3</sup>https://www.taobao.com/

Tal	ole	2:	Desc	ript	ion	of	datasets.
-----	-----	----	------	------	-----	----	-----------

Wildlife	Medicine	Pornography			
141,205,673	153,246,207	195,699,994			
251,584,066	518,175,838	96,180,405			
14,970,328	20,726,675	9,788,736			
7,793,425,568	9,720,654,470	11,699,840,416			
57.07	66.95	60.32			
262,281	394,306	408,978			
250,925	375,648	403,298			
11,356	18,658	5,680			
0.19%	0.26%	0.21%			
4.5%	5.0%	1.4%			
	Wildlife 141,205,673 251,584,066 14,970,328 7,793,425,568 57.07 262,281 250,925 11,356 0.19% 4.5%	WildlifeMedicine141,205,673153,246,207251,584,066518,175,83814,970,32820,726,6757,793,425,5689,720,654,47057.0766.95262,281394,306250,925375,64811,35618,6580.19%0.26%4.5%5.0%			

- **GBDT** [4]. This is also a classic machine learning algorithm which can detect the latent relevance of numerical and discrete features for classification.
- **GraphSAGE** [6]. This is a representative GNN model which exploits both structures and attributes via neighborhood aggregation to construct node representations.
- MTL [6]. This is a unified GNN framework which integrates multiple self-supervised learning tasks (e.g., clustering, graph partition and graph completion), and the semisupervised node classification together to utilize structural information as much as possible.
- GATNE [2]. This is an inductive heterogeneous GNN model used in industry which learns node representation considering the semantics within both of nodes and edges.
- HAN [33]. This is a heterogeneous GNN which models both node-level and semantic-level importance and designs a hierarchical message passing from heterogeneous neighbors.
- HGT [10]. This is a heterogeneous graph transformer which designs the heterogeneous mutual attention mechanism to aggregate information considering both edge and node types.
- **HSPD**<sub>*NW*</sub>. This is the modified version of HSPD which construct loss function with nodewise labels. We compare with this variant to showcase the effectiveness of heterogeneous self-supervision.
- HSPD. This is our proposed model consisting of heterogeneous self-supervised learning and directed pairwise selftraining to handle the few labeling problem in prohibited item detection scenarios

5.2.2 Implementation Details. All baselines and our HSPD are implemented with Tensorflow 1.12 on PAI<sup>4</sup> with Tesla GeForce GTX 1080 Ti Cluster. As the scales of datasets are quite large, we utilize AliGraph [40] API to load graphs and do sampling over HRGs in a distributed system. For a fair comparison, we randomly initialize model parameters with Gaussian distribution and optimize the model with Adam [16]. We set the batch size to 1024 for each worker, the number of workers to 8, the learning rate to 0.005, the feature embedding *d* to 64, the regularization weight  $\xi$  to 0.01 and the dropout rate to 0.4, the weight  $\beta$  to 1. In homogeneous GNNs (e.g., GraphSAGE and MTL), we randomly sample 5 neighbors for

eacn item. In heterogeneous GNNs (e.g., HAN, GATNE, HGT and our HSPD), we set the sample size as 5 for each relation. The maximum iteration of all the nodewise baselines is set to 300. For our proposed HSPD, we generate 6 pairs from each candidate to the risk seeds according to there structural and attributed similarity. We set the maximum iteration of HSPD to 100 and the epochs *T* of self-training strategy to 5. For each self-training process,  $\kappa_0$  is set to 1000 and  $\kappa_1$  is set to 100. We further discuss the hyper-parameter sensitivity in Section 5.7.

5.2.3 Evaluation Metrics. In our offline experiments, we calculate the **Max-F1** (the max F1 value by varying the threshold of recalled items) and **A**verage **P**recision (AP) to evaluate the global performance of identifying all test instances, which are the general metrics in the current system. In our online experiments, limited by manpower cost, we choose ACC@10000 (the accuracy of the top 10,000 recalled items which is manually reviewed) to measure the effectiveness of our HSPD. The larger values of Max-F1, AP or ACC@10000 indicate the better performance.

## 5.3 Performance Evaluation

We start by evaluating the detection performance of all the baselines and our HSPD on the three real-world risk scenarios. The overall Max-F1 and AP results of different methods under different scales of labels are presented in Table 3, from which the following observations can be made:

First, HSPD performs the best in all three risk scenarios with all different sizes of training sets. Compared with the baselines except HSPD<sub>NW</sub>, the improvement is prominent from 3.08% up to 9.40% in the Max-F1 metric and is from 3.76% up to 9.90% in the AP metric. Besides, our variant HSPD<sub>NW</sub> is better than the best baselines as well. The reason is twofold: (1) HSPD and HSPD<sub>NW</sub> fully combine both the features and semantic information to enhance node representation by designing the heterogeneous self-supervised learning. (2) HSPD enriches supervised information via the asymmetric pairwise labeling to discover patterns and self-training framework to expand confident pseudo labels, leading to the advantages over HSPD<sub>NW</sub>.

Second, HSPD has the ability to handle the imbalance and small scale of labels in prohibit item detection. On the one hand, Compared to *Medicine*, although *Wildlife* has few labels and *Pornography* has more imbalance labels, the improvements of our HSPD to baselines on the two datasets are both more obvious. On the other hand, the pairwise self-training mechanism performs more significant in the fewer labeling *Wildlife*, by comparing the improvement from HGT, HSPD<sub>NW</sub> to our HSPD on all the three datasets.

Third, Modeling the structural and the semantic information within the complex datasets contributes much to address the problem of prohibited item detection. By comparing with the classic LR and GBDT, almost all other methods achieve better performance on the three datasets. Furthermore, The supervised heterogeneous GNNs (i.e., HAN, HGT, HSPD<sub>NW</sub> and our HSPD) outperform homogeneous GNNs (GraphSAGE and MTL) because of the semantic modeling.

<sup>&</sup>lt;sup>4</sup>https://www.aliyun.com/product/bigdata/product/learn

Table 3: Performance of baselines and HSPD for risk detection on the three datasets. The best performance is in bold and the second best except  $HSPD_{NW}$  is underlined. Relative improvements of HSPD w.r.t. the second best are reported as well.

Dataset	Metric	Rate	LR	GBDT	GraphSAGE	MTL	GATNE	HAN	HGT	HSPDNW	HSPD	Improv.
Wildlife		20%	0.2958	0.5886	0.6829	0.7071	0.5610	0.6956	0.7354	0.7719	0.7883	7.20%
		40%	0.3044	0.6183	0.7122	0.7284	0.5732	0.7397	0.7721	0.8052	0.8318	7.74%
	AP	60%	0.3060	0.6414	0.7166	0.7447	0.5787	0.7518	0.7804	0.8085	0.8445	8.22%
		80%	0.3063	0.6414	0.7348	0.7525	0.5837	0.7686	0.7834	0.8216	0.8610	9.90%
		20%	0.3741	0.5638	0.6475	0.6683	0.5440	0.6632	0.6705	0.6940	0.7145	6.56%
		40%	0.3758	0.5770	0.6695	0.6775	0.5537	0.6862	0.7078	0.7252	0.7548	6.63%
	Max-F1	60%	0.3723	0.5932	0.6777	0.6888	0.5646	0.6936	0.7135	0.7348	0.7650	7.21%
		80%	0.3726	0.5932	0.6842	0.7019	0.5609	0.7020	0.7245	0.7502	0.7926	9.40%
		20%	0.5617	0.7594	0.7616	0.7747	0.6570	0.7942	0.8064	0.8200	0.8378	3.89%
	AD	40%	0.5630	0.7686	0.7727	0.7929	0.6675	0.8115	0.8270	0.8364	0.8575	3.69%
	AP	60%	0.5667	0.7752	0.7724	0.8009	0.6719	0.8208	0.8289	0.8455	0.8698	4.93%
Madiaina		80%	0.5646	0.7791	0.7788	0.7960	0.6732	0.8283	0.8430	0.8606	0.8747	3.76%
medicine		20%	0.5274	0.6977	0.6855	0.7043	0.5987	0.7242	0.7393	0.7341	0.7656	3.56%
	Mor E1	40%	0.5281	0.7097	0.6984	0.7201	0.6148	0.7341	0.7547	0.7591	0.7779	3.08%
	wiax-F1	60%	0.5273	0.7123	0.6989	0.7277	0.6210	0.7514	0.7559	0.7665	0.7819	3.44%
		80%	0.5264	0.7184	0.7005	0.7143	0.6201	0.7523	0.7656	0.7794	0.7958	3.94%
		20%	0.3181	0.5047	0.5099	0.5856	0.6161	0.6701	0.7082	0.7345	0.7699	8.72%
	AP	40%	0.3198	0.5442	0.5435	0.6167	0.6264	0.6801	0.7405	0.7852	0.8028	8.41%
		60%	0.3214	0.5485	0.5802	0.6263	0.6293	0.7005	0.7666	0.8026	0.8252	7.64%
Dornography		80%	0.3213	0.5590	0.5867	0.6659	0.6255	0.7433	0.7839	0.8032	0.8314	6.06%
Pornograpny		20%	0.3938	0.5183	0.5150	0.5707	0.5990	0.6372	0.6980	0.6967	0.7290	4.44%
	Max-F1	40%	0.3882	0.5304	0.5443	0.6043	0.6062	0.6732	0.7083	0.7291	0.7543	6.49%
		60%	0.3933	0.5340	0.5593	0.5961	0.6096	0.6804	0.7261	0.7429	0.7674	5.70%
		80%	0.3905	0.5316	0.5801	0.6387	0.6026	0.7133	0.7332	0.7470	0.7682	4.78%
0.77								0.77				
0.77		H	ISPD-w\o-S'		1 💻 🔳		HSPD-w/o-S				HSPD-w	\o-ST
0.76		H H	ISPD-w\0-S	W 0.82			HSPD-w/o-P	W 0.76		_	HSPD-w	\o-PW
0.73	0.75			0.80	0- HSPD			0.75 H			HSPD	



Figure 6: Performance comparison of HSPD and its variants on the three risk scenarios with 20% supervised information.

# 5.4 Variant Analysis

HSPD is to fully utilize the heterogeneous self-supervision and pairwise self-training to enhance node representation and enrich supervised information. Here we analyze three HSPD variants to evaluate the effectiveness of our design choices. (1) **HSPD-w\o-ST** removes the self-training process but utilizes both pairwise labels and self-supervision. (2) **HSPD-w\o-SS** removes the heterogeneous self-supervision. (3) **HSPD-w\o-PW** replaces pairwise labels with nodewise labels.

In Figure 6, we showcase the Max-F1 and AP performance of HSPD and its variants on all three datasets with 20% supervised information. There are two main observations as follows. (1) First, our

proposed HSPD outperforms all variants with an obvious improvement. Compared with  $HSPD-w \ o-SS$ , the improvements mainly result from the robust node representations by designing heterogeneous self-supervised learning to fully exploit both structural and semantic information within the complex datasets. Compared with  $HSPD-w \ o-PW$ , the advantage of pairwise labeling is proved as well, due to the ability of learning distance between risk items and candidates. Compared with  $HSPD-w \ o-ST$ , our HSPD introduces the self-training mechanism which can help us to avoid the over-fitting and learn a robust model. (2) Second, the  $HSPD-w \ o-SS$  often performs worse against  $HSPD-w \ o-ST$  and HSPD. This phenomenon is reasonable and explicable. Due to the noise and weak attributes but quite a few labels of items, the base representations of items



Figure 7: Performance comparison of HSPD by varying the number of pairwise labels and the epochs of self-training.



Figure 8: Case study on Wildlife.



Figure 9: The results of online testing.

would be limited, leading to the inaccurate prediction as well as the worse generation of pseudo labels during self-training.

# 5.5 Parameter Analysis

In this section, we investigate the effect of both the number of pairwise labels and the epochs of self-training, which are two key factors in enriching supervised information. We respectively vary k, the number of pairwise labels of each candidate from 3 to 12, and vary the epochs of self-training T from 0 to 4, and report the corresponding AP results with 20% self-supervised information on all three datasets in Figure 7. Notice that, epoch=0 indicates the pairwise learning without self-training.

There are two main observations. First, the performance of HSPD is related to the size of pairwise labels. By setting *T* as 0, We can easily find that there is an obvious improvement when we increase *k* from 3 to 6, indicating that too few pairwise cannot help training. Besides, too large may lead to expensive computational cost but obtain little increase. Here we set k = 6 to achieve a balance. Second, self-training mechanisms help HSPD to keep both outstanding and robust performance. With the self-training epoch *T* increases, the performance overall increase to be stable. In this paper, we set *T* as 4 to achieve a robust performance but avoid too heavy costs.

## 5.6 Case Study

We showcase some representative cases of detecting prohibited items on the *Wildlife* dataset. As shown in Figure 8, there are five patterns, including bamboo partridges, selaginella, red coral, hawksbill and illegal traps. Due to the advantages in evaluating the relevance between predictions and seeds via directed pairwise labelling, our HSPD obviously outperforms HAN and can discover various patterns of prohibited items to help manually reviewing.

#### 5.7 Online Experiments

We deploy HSPD on Taobao platform for online prohibited pornographic product detection and compare HSPD with GBDT via online testing. In fact, online service is very more challenging where the label rate is about 0.05%. The online results are shown in Figure 9. For daily results, we rank the candidate items with  $\hat{Y}$  and select top-10000 items for manually checking. The long-term observations show that HSPD outperforms GBDT in all 9 days. This phenomenon demonstrates the high industrial practicability of HSPD.

#### 6 CONCLUSION

In this paper, we study the problem of prohibited item detection which plays an important and essential role in ensuring the health of online shopping. In order to solve the challenges of too few manual labels, we are the first to model the large-scale item logs as a HRG and introduce the self-supervised learning and self-training in HRGs to address this problem, and then propose the novel HSPD. This approach considers the semantics of relations as the selfsupervision and generates the disentangled factors of items as the robust representation. Moreover, the directed pairwise self-training is designed in HSPD to enrich the supervised information. Extensive results on both offline and online experiments demonstrate the effectiveness of our proposed model.

#### ACKNOWLEDGMENTS

This research is supported in part by the National Natural Science Foundation of China (No. U20B2045, 61772082, 62002029, 62172052, U1936104, 61972442). This work is also supported by the Fundamental Research Funds for the Central Universities (2021RC28).

#### REFERENCES

- Hongyun Cai, Vincent W. Zheng, and Kevin Chen-Chuan Chang. 2018. A Comprehensive Survey of Graph Embedding: Problems, Techniques, and Applications. *IEEE Transactions on Knowledge and Data Engineering* 30, 9 (2018), 1616–1637.
- [2] Yukuo Cen, Xu Zou, Jianwei Zhang, Hongxia Yang, Jingren Zhou, and Jie Tang. 2019. Representation Learning for Attributed Multiplex Heterogeneous Network. In Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2019, Anchorage, AK, USA, August 4-8, 2019. ACM, 1358–1368.
- [3] Shaohua Fan, Junxiong Zhu, Xiaotian Han, Chuan Shi, Linmei Hu, Biyu Ma, and Yongliang Li. 2019. Metapath-guided Heterogeneous Graph Neural Network for Intent Recommendation. In Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2019, Anchorage, AK, USA, August 4-8, 2019. ACM, 2478–2486.
- [4] Jerome H Friedman. 2001. Greedy function approximation: a gradient boosting machine. Annals of statistics (2001), 1189–1232.
- [5] Alberto García-Durán and Mathias Niepert. 2017. Learning Graph Representations with Embedding Propagation. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA. 5119–5130.
- [6] William L. Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA. 1024–1034.
- [7] Guoxiu He, Yangyang Kang, Zhe Gao, Zhuoren Jiang, Changlong Sun, Xiaozhong Liu, Wei Lu, Qiong Zhang, and Luo Si. 2019. Finding Camouflaged Needle in a Haystack?: Pornographic Products Detection via Berrypicking Tree Model. In Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2019, Paris, France, July 21-25, 2019. ACM, 365–374.
- [8] Shifu Hou, Yanfang Ye, Yangqiu Song, and Melih Abdulhayoglu. 2017. HinDroid: An Intelligent Android Malware Detection System Based on Structured Heterogeneous Information Network. In Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Halifax, NS, Canada, August 13 - 17, 2017. ACM, 1507–1515.
- [9] Ziniu Hu, Yuxiao Dong, Kuansan Wang, Kai-Wei Chang, and Yizhou Sun. 2020. GPT-GNN: Generative Pre-Training of Graph Neural Networks. In KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020. ACM, 1857–1867.
- [10] Ziniu Hu, Yuxiao Dong, Kuansan Wang, and Yizhou Sun. 2020. Heterogeneous Graph Transformer. In WWW '20: The Web Conference 2020, Taipei, Taiwan, April 20-24, 2020. ACM / IW3C2, 2704–2710.
- [11] Dasol Hwang, Jinyoung Park, Sunyoung Kwon, Kyung-Min Kim, Jung-Woo Ha, and Hyunwoo J. Kim. 2020. Self-supervised Auxiliary Learning with Meta-paths for Heterogeneous Graphs. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.
- [12] Yugang Ji, Chuan Shi, Fuzhen Zhuang, and Philip S. Yu. 2019. Integrating Topic Model and Heterogeneous Information Network for Aspect Mining with Rating Bias. In Advances in Knowledge Discovery and Data Mining - 23rd Pacific-Asia Conference, PAKDD 2019, Macau, China, April 14-17, 2019, Proceedings, Part I (Lecture Notes in Computer Science, Vol. 11439). Springer, 160–171.
- [13] Yugang Ji, Mingyang Yin, Hongxia Yang, Jingren Zhou, Vincent W. Zheng, Chuan Shi, and Yuan Fang. 2021. Accelerating Large-Scale Heterogeneous Interaction Graph Embedding Learning via Importance Sampling. ACM Transactions on Knowledge Discovery from Data 15, 1 (2021), 10:1–10:23.
- [14] Wei Jin, Tyler Derr, Haochen Liu, Yiqi Wang, Suhang Wang, Zitao Liu, and Jiliang Tang. 2020. Self-supervised Learning on Graphs: Deep Insights and New Direction. CoRR abs/2006.10141 (2020).
- [15] Jongmin Kim, Taesup Kim, Sungwoong Kim, and Chang D. Yoo. 2019. Edge-Labeling Graph Neural Network for Few-Shot Learning. In IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019. Computer Vision Foundation / IEEE, 11–20.
- [16] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, Yoshua Bengio and Yann LeCun (Eds.).
- [17] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net.
- [18] Qimai Li, Zhichao Han, and Xiao-Ming Wu. 2018. Deeper Insights Into Graph Convolutional Networks for Semi-Supervised Learning. In Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018. AAAI Press, 3538–3545.
- [19] Xiao Liu, Fanjin Zhang, Zhenyu Hou, Zhaoyu Wang, Li Mian, Jing Zhang, and Jie Tang. 2020. Self-supervised Learning: Generative or Contrastive. CoRR

abs/2006.08218 (2020).

- [20] Yuanfu Lu, Yuan Fang, and Chuan Shi. 2020. Meta-learning on Heterogeneous Information Networks for Cold-start Recommendation. In KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020. ACM, 1563–1573.
- [21] Jianxin Ma, Chang Zhou, Hongxia Yang, Peng Cui, Xin Wang, and Wenwu Zhu. 2020. Disentangled Self-Supervision in Sequential Recommenders. In KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020. ACM, 483–491.
- [22] Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient Estimation of Word Representations in Vector Space. In 1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings.
- [23] Chanyoung Park, Donghyun Kim, Jiawei Han, and Hwanjo Yu. 2020. Unsupervised Attributed Multiplex Network Embedding. In *The Thirty-Fourth AAAI* Conference on Artificial Intelligence, AAAI 2020, New York, NY, USA, February 7-12, 2020. AAAI Press, 5371–5378.
- [24] Daniel Carlos Guimarães Pedronette and Longin Jan Latecki. 2021. Rank-based self-training for graph convolutional networks. *Information Processing and Man*agement 58, 2 (2021), 102443.
- [25] Zhen Peng, Wenbing Huang, Minnan Luo, Qinghua Zheng, Yu Rong, Tingyang Xu, and Junzhou Huang. 2020. Graph Representation Learning via Graphical Mutual Information Maximization. In WWW '20: The Web Conference 2020, Taipei, Taiwan, April 20-24, 2020. ACM / IW3C2, 259–270.
- [26] Michael Sejr Schlichtkrull, Thomas N. Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and Max Welling. 2018. Modeling Relational Data with Graph Convolutional Networks. In The Semantic Web - 15th International Conference, ESWC 2018, Heraklion, Crete, Greece, June 3-7, 2018, Proceedings (Lecture Notes in Computer Science, Vol. 10843). Springer, 593–607.
- [27] Chuan Shi, Yitong Li, Jiawei Zhang, Yizhou Sun, and S Yu Philip. 2016. A survey of heterogeneous information network analysis. *IEEE Transactions on Knowledge* and Data Engineering 29, 1 (2016), 17–37.
- [28] Kaisong Song, Yangyang Kang, Wei Gao, Zhe Gao, Changlong Sun, and Xiaozhong Liu. 2021. Evidence Aware Neural Pornographic Text Identification for Child Protection. In AAAI.
- [29] Fan-Yun Sun, Jordan Hoffmann, Vikas Verma, and Jian Tang. 2020. InfoGraph: Unsupervised and Semi-supervised Graph-Level Representation Learning via Mutual Information Maximization. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.
- [30] Ke Sun, Zhouchen Lin, and Zhanxing Zhu. 2020. Multi-Stage Self-Supervised Learning for Graph Convolutional Networks on Graphs with Few Labeled Nodes. In The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, New York, NY, USA, February 7-12, 2020. AAAI Press, 5892–5899.
- [31] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2017. Graph Attention Networks. *CoRR* abs/1710.10903 (2017).
- [32] Petar Velickovic, William Fedus, William L. Hamilton, Pietro Liò, Yoshua Bengio, and R. Devon Hjelm. 2019. Deep Graph Infomax. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net.
- [33] Xiao Wang, Houye Ji, Chuan Shi, Bai Wang, Yanfang Ye, Peng Cui, and Philip S. Yu. 2019. Heterogeneous Graph Attention Network. In *The World Wide Web Conference, WWW 2019, San Francisco, CA, USA, May 13-17, 2019.* ACM, 2022–2032.
- [34] Xiao Wang, Ruijia Wang, Chuan Shi, Guojie Song, and Qingyong Li. 2020. Multi-Component Graph Convolutional Collaborative Filtering. In The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, New York, NY, USA, February 7-12, 2020. AAAI Press, 6267–6274.
- [35] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. 2021. A Comprehensive Survey on Graph Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems*. 32, 1 (2021), 4–24.
- [36] Hong Xuan, Abby Stylianou, and Robert Pless. 2020. Improved Embeddings with Easy Positive Triplet Mining. In IEEE Winter Conference on Applications of Computer Vision, WACV 2020, Snowmass Village, CO, USA, March 1-5, 2020. IEEE, 2463–2471.
- [37] Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. 2020. Graph Contrastive Learning with Augmentations. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.
- [38] Wentao Zhang, Yuan Fang, Zemin Liu, Min Wu, and Xinming Zhang. 2020. mg2vec: Learning Relationship-Preserving Heterogeneous Graph Representations via Metagraph Embedding. *IEEE Transactions on Knowledge and Data Engineering* (2020), 1–1. https://doi.org/10.1109/TKDE.2020.2992500
- [39] Vincent W. Zheng, Mo Sha, Yuchen Li, Hongxia Yang, Yuan Fang, Zhenjie Zhang, Kian-Lee Tan, and Kevin Chen-Chuan Chang. 2018. Heterogeneous Embedding Propagation for Large-Scale E-Commerce User Alignment. In IEEE International Conference on Data Mining, ICDM 2018, Singapore, November 17-20, 2018. IEEE

- Computer Society, 1434–1439.
  [40] Rong Zhu, Kun Zhao, Hongxia Yang, Wei Lin, Chang Zhou, Baole Ai, Yong Li, and Jingren Zhou. 2019. AliGraph: A Comprehensive Graph Neural Network Platform. *Proceedings of VLDB Endowment.* 12, 12 (2019), 2094–2105.
- [41] Barret Zoph, Golnaz Ghiasi, Tsung-Yi Lin, Yin Cui, Hanxiao Liu, Ekin Dogus Cubuk, and Quoc Le. 2020. Rethinking Pre-training and Self-training. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.