

Negative Sampling with Adaptive Denoising Mixup for Knowledge Graph Embedding

Xiangnan Chen, Wen Zhang, Zhen Yao, Mingyang Chen, and Siliang Tang*

Zhejiang University, China

{xnchen2020, zhang.wen, 22151303, mingyangchen, siliang}@zju.edu.cn

Abstract. Knowledge graph embedding (KGE) aims to map entities and relations of a knowledge graph (KG) into a low-dimensional and dense vector space via contrasting the positive and negative triples. In the training process of KGEs, negative sampling is essential to find high-quality negative triples since KGs only contain positive triples. Most existing negative sampling methods assume that non-existent triples with high scores are high-quality negative triples. However, negative triples sampled by these methods are likely to contain noise. Specifically, they ignore that non-existent triples with high scores might also be true facts due to the incompleteness of KGs, which are usually called false-negative triples. To alleviate the above issue, we propose an easily pluggable denoising mixup method called **DeMix**, which generates high-quality triples by refining sampled negative triples in a self-supervised manner. Given a sampled unlabeled triple, DeMix firstly classifies it into a marginal pseudo-negative triple or a negative triple based on the judgment of the KGE model itself. Secondly, it selects an appropriate mixup partner for the current triple to synthesize a partially positive or a harder negative triple. Experimental results on the knowledge graph completion task show that the proposed DeMix is superior to other negative sampling techniques, ensuring corresponding KGEs a faster convergence and better link prediction results.

Keywords: Knowledge graph embeddings · Negative sampling · Mixup.

1 Introduction

Recently, knowledge graphs (KGs) have been successfully profitable in many practical applications, including question answering [12], information retrieval [26] and dialogue systems [28]. However, the KGs constructed manually or automatically still suffer from incompleteness. Thus, completing KGs through efficient representation learning has been a hot topic. For flexible and efficient KG representation learning, knowledge graph embedding (KGE)[4,30] aims to represent entities and relations in KGs with real-valued vectors, also called embeddings. KGEs have shown promising performance in KG related tasks, such as triple classification [14] and link prediction [18]. They generally follow the same

* Corresponding author.

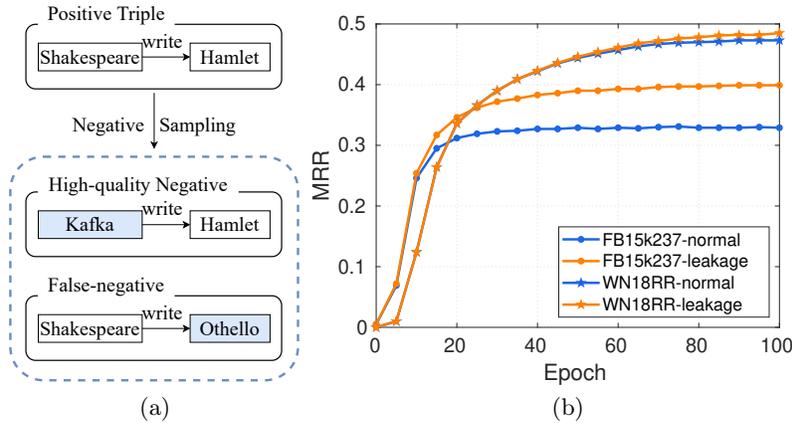


Fig. 1. (a) The example of sampling false-negative triples. (b) Testing MRR performance v.s. Epoch based on RotatE. Normal means self-adversarial negative sampling. Leakage means ensuring sampled negative triples are not contained in the validation or test set. We regard triples in the validation and test sets as false-negative triples during training.

training paradigm. Specifically, they define a score function to measure the plausibility of triples through calculation with entity and relation embeddings, then learn the embeddings with the training objective of enlarging the gap between the scores of the positive and negative triples.

Since KGs only contain positive triples, negative sampling methods for KGEs usually regard non-existent triples as negative triples. Formally, given a fact (h, r, t) of a KG \mathcal{G} , negative sampling methods sample an entity e among all candidate entities and replace the head entity h or tail entity t with e to form a corrupted triple $(e, r, t) \notin \mathcal{G}$ or $(h, r, e) \notin \mathcal{G}$ as a negative triple. Previous works [23,31] have proved the quality of negative triples significantly affects the performance of KGEs, such as low-quality negative triples can cause gradient vanishing problems[23]. Therefore, searching for high-quality negative triples is not only necessary but also vital for learning KGEs.

Many negative sampling methods [31,6] assume that non-existent triples with high scores are high-quality negative triples. These methods optimize the mechanism of negative sampling from different perspectives to search for non-existent triples with high scores. For example, KBGAN[6] and IGAN[23] introduce a generative adversarial network (GAN)[11] to generate negative triples with high scores, and NSCaching[31] introduces a caching mechanism to store corrupted triples with large scores, etc. However, these negative sampling methods for KGEs neglect the issue of sampling noisy triples, especially when they regard non-existing triples with high scores as high-quality negative triples. Because those corrupted triples with high scores might also be true facts with high probability due to the complementary capability of KGEs, which are usually called false-negative triples. As shown in Figure 1(a), given a positive triple (*Shake-*

speare, write, Hamlet), the false-negative triple such as (*Shakespeare, write, Othello*) may be sampled, which will give imprecise supervision signals to KGE models’ training. Furthermore, we set up a toy experiment ¹ to quantify the effect of false-negative triples on KGEs’ training. As shown in Figure 1(b), compared with normal negative sampling, negative sampling with data leakage can improve MRR by 21.3% and 1.2% on FB15K237 and WN18RR respectively. This phenomenon indicates that although the probability of sampling false-negative triples is very low, the false-negative triples do mislead the learning of KGE and degrade the inference ability of KGE models.

Therefore, it is important to consider the challenging denoising problem when sampling high-quality negative triples. In this paper, to address the above issue, we propose a novel and easily pluggable framework **DeMix** which could generate high-quality triples that are beneficial for models’ training by refining negative triples. Specifically, DeMix contains two modules, namely Marginal Pseudo-Negative triple Estimator (MPNE) and Adaptive Mixup (AdaMix). Given sampled corrupted triples (h, r, e) or (e, r, t) , the MPNE module firstly leverages the current predictive results of the KGE model to estimate whether these corrupted triples contain noisy triples, then divides them into marginal pseudo-negative triples which more likely contain noisy triples and true-negative ones which are more likely negative triples. Then, in order to refine corrupted triples, the AdaMix module selects an appropriate mixup partner e' for e , then mixes them in the entity embedding space. Overall, DeMix generates partially positive triples for marginal pseudo-negative triples and harder negative triples for true negative ones as high-quality triples to help the training of KGEs.

In summary, the contributions of our work are as follows:

- We propose a simple and efficient denoising framework, named **DeMix**. It generates high-quality triples by refining sampled negative triples without additional information.
- We design two pluggable modules MPNE and AdaMix, which are general to be combined with other negative sampling methods and are efficient for not wasting training time on additional parameters.
- We conduct extensive experiments on two benchmark datasets to illustrate the effectiveness and the superiority of our whole framework and each module.

2 Preliminaries

Notations We use lower-case letters h , r , and t to represent the head entity, relation, and tail entity in triples respectively, and the corresponding boldface lower-case letters \mathbf{h} , \mathbf{r} and \mathbf{t} indicate their embeddings. Let \mathcal{E} and \mathcal{R} represent the set of entities and relations respectively, we denote a knowledge graph as $\mathcal{G} = \{\mathcal{E}, \mathcal{R}, \mathcal{T}\}$ where $\mathcal{T} = \{(h, r, t)\} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$ is the set of facts in KG. For training

¹ We use the official open source code of RotatE model with the best config parameters: <https://github.com/DeepGraphLearning/KnowledgeGraphEmbedding>

of KGEs, a dataset \mathcal{D} of triples with labels is used $\mathcal{D} = \{(h_i, r_i, t_i), y_i\}_{i=1}^{n_p+n_u}$, which includes n_p positive triples with label $y = 1$ and n_u negative triples with label $y = 0$. n_p usually is equivalent to the size of fact set that $n_p = |\mathcal{T}|$. Since negative triples are not included in \mathcal{G} , thus they are usually created by negative sampling methods.

Negative Sampling for KGE Unlike the negative sampling at the instance level in the computer vision domain [9], the negative sampling for KGE is sampling and replacing entities in facts. More specifically, based on the fact set \mathcal{T} , the negative sampling for KGE sample entities $e \in \mathcal{E}$ and replace either h or t in facts with e to construct a set of corrupted triples \mathcal{T}_u which is computed as follows:

$$\mathcal{T}_u = \bigcup_{(h,r,t) \in \mathcal{T}} \{(e, r, t) \notin \mathcal{T}\} \cup \{(h, r, e) \notin \mathcal{T}\}. \quad (1)$$

Following the closed-world assumption [14], the majority of negative sampling methods for KGE treat all non-existent triples as negative triples. Thus,

$$\mathcal{D} = \{((h, r, t), y = 1) | (h, r, t) \in \mathcal{T}\} \cup \{((h', r, t'), y = 0) | (h', r, t') \in \mathcal{T}_u\}. \quad (2)$$

Therefore, the optimization objective for KGE models learning based on closed-world assumption(CWA) can be formulated as follows [18]:

$$\mathcal{L} = -\log \sigma(f(h, r, t)) - \sum_{(h', r, t') \in \mathcal{T}_u} \log \sigma(-f(h', r, t')), \quad (3)$$

where $f(h, r, t)$ denotes the score function of KGEs to assess the credibility of (h, r, t) .

KGE Score Function The two most typical score functions for a triple (h, r, t) are:

(1) The **distance**-based score function which evaluates the score of triples based on the Euclidean distance between vectors, such as the score function of RotatE [18]:

$$f(h, r, t) = \gamma - \|\mathbf{h} \circ \mathbf{r} - \mathbf{t}\|, \quad (4)$$

where γ is the margin, \circ indicates the hardmard product.

(2) The **similarity**-based score function which evaluates the score of triples based on dot product similarity between vectors, such as the score function of DistMult [27]:

$$f(h, r, t) = \mathbf{h}^\top \text{diag}(\mathbf{M}_r) \mathbf{t}, \quad (5)$$

where $\text{diag}(\mathbf{M}_r)$ represents the diagonal matrix of the relation r .

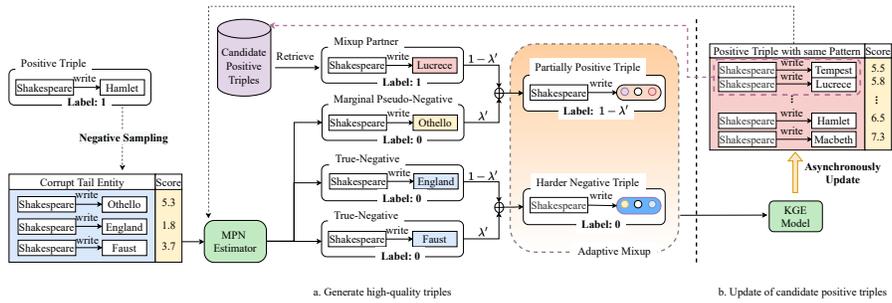


Fig. 2. An overview of the DeMix framework based on corrupting tail entity.

3 Methodology

In this section, we introduce the proposed method **DeMix**, which is a novel and easily pluggable framework for generating high-quality triples. Recalling the denoising challenge in sampling high-quality negative triples, we design two modules to address the above challenge, namely Marginal Pseudo-Negative Triple Estimator (MPNE) and Adaptive Mixup (AdaMix) modules. The MPNE module leverages the current predictive results of the KGE models to divide unlabeled corrupted triples into pseudo-negative triples and true-negative triples. Then, the AdaMix module selects a suitable mixup partner for each corrupted triple and mixes them in the entity embedding space to generate partially correct triples or harder negative triples to help train the KGE model. The overview of DeMix is shown in Figure 2.

3.1 Marginal Pseudo-Negative Triple Estimator (MPNE)

Motivated by sampled corrupted triples with high scores more likely to contain noise due to the incompleteness of KGs, we aim to leverage the predictive results of the KGE model itself to recognize noisy triples. In this work, we divide sampled corrupted triples into two subsets, such as marginal pseudo-negative triples which are likely to contain noise according to the current predictive results of the KGE model, and true-negative triples which are more likely negative triples. Because the KGE model has different learning levels for each relation pattern due to the long-tail distribution in KGs[25] and the phenomenon that the discriminative ability of the KGE model is continuously improving during the training process. We need to design a module to dynamically estimate noisy triples. Considering negative sampling for KGE is to replace the head entity or tail entity in a positive triple, the unreplaced binary terms in a positive triple are invariant. We call the invariant binary terms a pattern. So \mathcal{T}_u can be reformulated as follows:

$$\begin{aligned} \mathcal{T}_u = & \{(e, pattern) | e \in \mathcal{E}, pattern = (r, t)\} \\ & \cup \{(pattern, e) | e \in \mathcal{E}, pattern = (h, r)\}. \end{aligned} \quad (6)$$

The set of positive triples with the same pattern $((h, r)$ or (r, t)) as follows:

$$\begin{aligned}\mathcal{T} &= \mathcal{T}_{pattern}^{rt} \cup \mathcal{T}_{pattern}^{hr}, \\ \mathcal{T}_{pattern}^{rt} &= \{(e_i, r, t) \in \mathcal{T} | e_i \in \mathcal{E}\}, \\ \mathcal{T}_{pattern}^{hr} &= \{(h, r, e_j) \in \mathcal{T} | e_j \in \mathcal{E}\}.\end{aligned}\quad (7)$$

Then we treat corrupted triples whose scores are close to the scores of the positive samples with the same pattern as marginal pseudo-negative triples \mathcal{T}_{mpn} . So \mathcal{T}_{mpn} can be formulated as follows:

$$\begin{aligned}\mathcal{T}_{mpn} &= \{((h', r, t'), y = 0) | (h', r, t') \in \mathcal{T}_u, \\ &\quad -\delta_T + [f^{\mathcal{T}_{pattern}}]_{min} \leq f(h', r, t') \leq [f^{\mathcal{T}_{pattern}}]_{mean}\},\end{aligned}\quad (8)$$

where $f^{\mathcal{T}_{pattern}}$ is the collection of scores of positive triples with the same pattern, i.e., $f^{\mathcal{T}_{pattern}} = \{f(h, r, t) | (h, r, t) \in \mathcal{T}_{pattern}\}$, and $[X]_{min}$ and $[X]_{mean}$ is the minimum and mean value of X . δ_T is a hyper-parameter controlling the estimation range at the T -th training epoch. Specifically, $\delta_T = \delta \cdot \min(\beta, T/T_0)$, where T_0 denotes the threshold of stopping increase, δ and β are hyper-parameters. Notably, inspired by the characteristics of complex relations in KG, namely 1-N, N-1, 1-1, and N-N defined in TransH [24], we record the num of positive triples with the same pattern. When $|\mathcal{T}_{pattern}|$ is lower, the probability that the corrupted triples based on this pattern are noisy triples is lower. So after sampling negative triples \mathcal{T}_u , we set a threshold μ to decide whether to estimate the set \mathcal{T}_{mpn} from \mathcal{T}_u as follows:

$$\mathcal{T}_u = \begin{cases} \mathcal{T}_{mpn} \cup \mathcal{T}_{\bar{u}} & \text{if } |\mathcal{T}_{pattern}| \geq \mu, \\ \mathcal{T}_{\bar{u}} & \text{if } |\mathcal{T}_{pattern}| < \mu, \end{cases}\quad (9)$$

where $\mathcal{T}_{\bar{u}}$ contains triples which are regarding as true-negative triples.

3.2 Adaptive Mixup(AdaMix)

To address the denoising challenge in sampling high-quality negative triples, an intuitive approach is to directly label the triples in \mathcal{T}_{mpn} as 1. However, when the KGE model does not have the strong distinguishable ability, especially in the early training stage. this approach can give the wrong supervisory signal to the KGE model. Inspired by the recent progress of mixup [13], we aim to adapt the mixup technique to alleviate the wrong supervisory signal issue. Since the triples in \mathcal{T}_{mpn} are more likely to be positive triples than $\mathcal{T}_{\bar{u}}$, we develop an adaptive mixup mechanism to guide the selection of a suitable mixup partner for each corrupted triple, to generate high-quality triples containing rich information. In specific, we first build a pool of candidate positive triples \mathcal{T}_{cap} for each pattern from $\mathcal{T}_{pattern}$. \mathcal{T}_{cap} contains the positive triples around the current learned boundary of the KGE model. \mathcal{T}_{cap} can be formulated as follows:

$$\begin{aligned}\mathcal{T}_{cap} &= \{((h, r, t), y = 1) | (h, r, t) \in \mathcal{T}_{pattern}, \\ &\quad f(h, r, t) \leq [f^{\mathcal{T}_{pattern}}]_{mean}\}.\end{aligned}\quad (10)$$

Since the embeddings of the KGE model are continually updated, we thus propose to update \mathcal{T}_{cap} every epoch. Then, we select a mixup partner with the same pattern for each corrupted triple in \mathcal{T}_u . For marginal pseudo-negative triples in \mathcal{T}_{mpn} , which are more likely to be positive but annotated by negative, we uniformly choose a mixup partner from \mathcal{T}_{cap} to generate a partially positive triple to take more precise supervision to the KGE model. Besides, for true-negative triples in $\mathcal{T}_{\bar{u}}$, we select another true-negative triple as a mixup partner to construct harder negative triples using a non-existent mixing entity in KG. It is worth noting that since the pattern of the mixup partner is the same as the pattern of the corrupted triple, the actual mixing object is mixing between corrupted candidate entity and the entity in the corresponding position of the mixup partner. Let $\mathbf{e}_i, \mathbf{e}_j$ denote the entity to be mixed in the corrupted triple and the corresponding mixup partner respectively, the overall mixup partner selection is formulated as follows:

$$(\mathbf{e}_j, y_j) \sim \begin{cases} \text{Uniform}(\mathcal{T}_{cap}) & \text{if } (\mathbf{e}_i, y_i) \in \mathcal{T}_{mpn}, \\ \text{Uniform}(\mathcal{T}_{\bar{u}}) & \text{if } (\mathbf{e}_i, y_i) \in \mathcal{T}_{\bar{u}}. \end{cases} \quad (11)$$

Finally we mix each corrupted triple $(h'_i, r_i, t'_i) \in \mathcal{T}_{mpn} \cup \mathcal{T}_{\bar{u}}$ with its corresponding mixup partner (h_j, r_j, t_j) in the entity embedding space of the KGE model to generate an augmented triple set $\widehat{\mathcal{T}}_u$. Motivated by the modified mixup operator [2], the mixup operation is as follows:

$$\begin{aligned} \widehat{\mathbf{e}}_i &= \lambda' \mathbf{e}_i + (1 - \lambda') \mathbf{e}_j, & \widehat{y}_i &= \lambda' y_i + (1 - \lambda') y_j, \\ \lambda' &= \max(\lambda, 1 - \lambda), & \lambda &\sim \text{Beta}(\alpha, \alpha), \alpha \in (0, \infty), \end{aligned} \quad (12)$$

where λ' is a balance parameter. λ' can guarantee that the feature of each augmented entity $\widehat{\mathbf{e}}_i$ is closer to \mathbf{e}_i than the mixup partner \mathbf{e}_j .

3.3 Training the KGE Model

Because our framework is pluggable, we can combine other negative sampling methods to train the KGE model. Here we show the training objectives based on the uniform negative sampling [4] and self-adversarial sampling [18]. The loss function based on uniform sampling is as follows:

$$\mathcal{L} = \ell(f(h, r, t), 1) + \sum_{(\widehat{h}_i, r, \widehat{t}_i) \in \widehat{\mathcal{T}}_u} \ell(f(\widehat{h}_i, r, \widehat{t}_i), \widehat{y}_i). \quad (13)$$

The loss function based on self-adversarial sampling is as follows:

$$\begin{aligned} \mathcal{L} &= \ell(f(h, r, t), 1) \\ &+ \sum_{(\widehat{h}_i, r, \widehat{t}_i) \in \widehat{\mathcal{T}}_u} p(\widehat{h}_i, r, \widehat{t}_i) \ell(f(\widehat{h}_i, r, \widehat{t}_i), \widehat{y}_i), \end{aligned} \quad (14)$$

Algorithm 1 Training procedure of DeMix

Input: training set $\mathcal{P} = \{(h, r, t)\}$, entity set \mathcal{E} , relation set \mathcal{R} , embedding dimension d , scoring function f , the size of epoches E , the size of warm-up epoches W .

Output: embeddings for each $e \in \mathcal{E}$ and $r \in \mathcal{R}$

- 1: **Initialize** embeddings for each $e \in \mathcal{E}$ and $r \in \mathcal{R}$;
 - 2: Get $\mathcal{T}_{pattern}$ from \mathcal{T} ;
 - 3: **while** epoch $< W$ **do**
 - 4: Warm up model using uniform sampling with K negative samples;
 - 5: **end while**
 - 6: **while** epoch $< E$ **do**
 - 7: Sample a mini-batch $\mathcal{T}_{batch} \in \mathcal{T}$;
 - 8: **while** $(h, r, t) \in \mathcal{T}_{batch}$ **do**
 - 9: Uniformly sample M entities from \mathcal{E} to form unlabeled corrupted triplets $\mathcal{T}_u = \{(h_m, r, t_m), m = 1, 2, \dots, M\}$;
 - 10: Estimate marginal pseudo-negative triples \mathcal{T}_{mpn} using Eq.(9)(8);
 - 11: Select the mixup partners for each corrupted triples using Eq.(11);
 - 12: Construct $\hat{\mathcal{T}}_u$ applying Eq.(12) to corrupted triples and their mixup partners;
 - 13: calculate loss functions using Eq.(13), then update the embeddings of entites and relations via gradient descent;
 - 14: **end while**
 - 15: Update \mathcal{T}_{cap} using Eq.(10);
 - 16: **end while**
-

where $\hat{\mathcal{T}}_u = \text{AdaptiveMixup}(\mathcal{T}_u, \mathcal{T}_{pattern}, \alpha)$ and $\ell(\cdot, \cdot)$ is the cross-entropy loss. Movever, the probability distribution of sampling high-quality triples is as follows:

$$p(\hat{h}_j, r, \hat{t}_j | \{(\hat{h}_i, r, \hat{t}_i)\}) = \frac{\exp \alpha_t f(\hat{h}_j, r, \hat{t}_j)}{\sum_i \exp \alpha_t f(\hat{h}_i, r, \hat{t}_i)}, \quad (15)$$

where α_t is the temperature of sampling. The full training procedure is shown in Algorithm 1

4 Experiment

In this section, we perform detailed experiments to demonstrate the effectiveness of our proposed framework DeMix by answering the following questions. **Q1:** Does DeMix can mitigate the noisy triples issue? **Q2:** Whether DeMix can be effectively plugged into other negative sampling methods? **Q3:** How does each of designed modules influence the performance of DeMix?

4.1 Experiment Settings

Datasets Two public datasets are utilized for experiments, including WN18RR [10] and FB15K237 [20]. WN18RR is a subset of WN18 [4], where inverse relations are deleted. Similarly, FB15k237 is a subset of FB15K [4], which comes

Dataset	#Rel	#Ent	#Train	#Valid	#Test
WN18RR	11	40,943	86,835	3,034	3,134
FB15K237	237	14,541	272,115	17,535	20,466

Table 1. Statistics of two benchmarks. #Rel, #Ent, represent the number of relations, and entities of each dataset, respectively.

from FreeBase [3]. FB15K237 is denser than WN18RR, so it is more affected by false-negative triples. The statistics of the two datasets are given in Table 1.

Baseline methods We compare DeMix to seven negative sampling baselines. The details of baseline methods are presented as follows:

- *Uniform Sampling* [4]. The basic negative sampling method, which samples negative triples from a uniform distribution.
- *Bernoulli Sampling* [24]. which sample negative triples from a Bernoulli distribution considering false-negative triples.
- *NSCaching* [31]. The NSCaching introduces the cache strategy as a general negative sampling scheme.
- *Self-adversarial Sampling* [18]. It utilizes a self-scoring function and samples negative triples according to the current embedding model.
- *RW-SANS* [1]. It samples negative triples from the k-hop of the node neighborhood by utilizing the graph structure.
- *CANS* [16]. The CANS is a component of CAKE [16] responsible for solving the invalid negative sampling challenge. Considering our method focuses on negative sampling for KGEs, we mainly compare CANS instead of CAKE.
- *ESNS* [29]. It takes semantic similarities among entities into consideration to tackle the issue of false-negative samples.

Evaluation protocol Following the previous work [16], we calculate the score of triples in the test dataset by employing the learned KG embeddings and the score function. Then we can get the rank of the correct entity for each test triple based on the filtered setting, where all corrupted triples in the dataset are removed. The performance is evaluated by four metrics: mean reciprocal rank (MRR) and Hits at 1,3,10 (Hits@N). Higher MRR and Hits@N mean better performance.

Implementation details We firstly use the negative sampling method to warm-up² 8 epochs to give the KGE model discriminative ability, and then use our method to generate high-quality triples for further training. We use PyTorch [17] framework to implement our method and Adam [15] optimizer for model training. The mixup balance hyperparameter α is fixed to 1. In addition,

² We will discuss warm-up in the ablation study

Translational Distance-based Models	WN18RR		FB15K237	
	MRR	Hits@10	MRR	Hits@10
Neglecting False-negative Triples				
TransE+Uniform	0.175*	0.445*	0.171 [□]	0.323 [□]
TransE+NSCaching	0.200*	0.478*	0.205 [□]	0.353 [□]
TransE+Self-Adv	0.215 [△]	<u>0.516[△]</u>	0.268 [□]	0.454 [□]
TransE+RW-SANS [△]	<u>0.218</u>	0.510	0.295	0.483
Considering False-negative Triples				
TransE+Bernoulli*	0.178	0.451	0.256	0.419
TransE+CANS [□]	-	-	<u>0.298</u>	<u>0.490</u>
Ours				
TransE+ DeMix-Adv	0.220	0.521	0.318	0.510

Table 2. Link prediction results for TransE on two datasets. **DeMix-Adv** denotes DeMix based on self-adversarial sampling, **Bold** numbers are the best results for each type of model. Underlined numbers mean the best performances of baselines.

Translational Distance-based	WN18RR				FB15K237			
	MRR	Hits@10	Hits@3	Hits@1	MRR	Hits@10	Hits@3	Hits@1
RotatE+Uniform [△]	0.471	0.560	0.488	0.424	0.282	0.462	0.314	0.191
RotatE+Self-Adv	0.476 [△]	0.570 [△]	0.490 [△]	0.428 [△]	0.269 [□]	0.452 [□]	0.298 [□]	0.179 [□]
RotatE+RW-SANS [△]	<u>0.478</u>	<u>0.572</u>	0.494	0.430	0.295	0.481	0.327	0.202
RotatE+CANS [□]	-	-	-	-	<u>0.296</u>	<u>0.486</u>	<u>0.329</u>	<u>0.202</u>
RotatE+ DeMix-Adv	0.479	0.576	0.492	0.428	0.329	0.518	0.366	0.235
HAKE+Uniform [△]	0.493	0.580	0.510	0.450	0.304	0.482	0.333	0.216
HAKE+Self-Adv	<u>0.495[△]</u>	<u>0.580[△]</u>	<u>0.513[△]</u>	<u>0.450[△]</u>	0.306 [□]	0.486 [□]	0.337 [□]	0.216 [□]
HAKE+RW-SANS [△]	0.492	0.579	0.507	0.446	0.305	0.488	0.336	0.214
HAKE+CANS [□]	-	-	-	-	<u>0.315</u>	<u>0.501</u>	<u>0.344</u>	<u>0.221</u>
HAKE+ DeMix-Adv	0.498	0.584	0.514	0.451	0.337	0.533	0.374	0.239

Table 3. Link prediction results on RotatE and HAKE. **DeMix-Adv** denotes DeMix based on self-adversarial sampling, **Bold** numbers are the best results for each type of model. Underlined numbers mean the best performances of baselines.

we use grid search to tune hyper-parameters on the validation dataset. Specifically, we choose β from $\{1, 3\}$, δ from $\{0, 0.1, 1\}$, μ from $\{1, 3, 5\}$. The learning rate is tuned from 0.00001 to 0.01. The margin is chosen from $\{3, 6, 9, 12\}$. To make a fair comparison, the negative sampling size of each baseline is the same as generated negative triples size in our method. Specifically, for translational distance-based models, we follow the experimental setup in CANS [16]. The negative sampling size is 16. For semantic matching-based approaches, we set the negative sampling size as 50 following ESNS[29].

4.2 Experimental Results (Q1)

We compare the overall performance of DeMix applied to different KGE models in the link prediction task to answer Q1. X^* , X^\square , and X^Δ indicate the results

Semantic Matching Models	WN18RR		FB15K237	
	MRR	Hits@10	MRR	Hits@10
DistMult+Uniform	0.412	0.463	0.213	0.383
DistMult+Bernoulli	0.396	0.437	0.262	0.430
DistMult+NSCaching	0.413	0.455	0.288	0.458
DistMult+Self-Adv	0.416	0.463	0.215	0.395
DistMult+ESNS	<u>0.424</u>	<u>0.488</u>	<u>0.296</u>	<u>0.465</u>
DistMult+ DeMix-Adv	0.439	0.535	0.301	0.470
ComplEx+Uniform	0.429	0.478	0.214	0.387
ComplEx+Bernoulli	0.405	0.441	0.268	0.442
ComplEx+NSCaching	0.446	0.509	0.302	0.481
ComplEx+Self-Adv	0.435	0.493	0.211	0.395
ComplEx+ESNS	<u>0.450</u>	<u>0.512</u>	<u>0.303</u>	0.471
ComplEx+ DeMix-Adv	0.468	0.552	0.307	0.479

Table 4. Link prediction results for semantic matching KGE models on two datasets. **DeMix-Adv** denotes DeMix based on self-adversarial sampling, **Bold** numbers are the best results for each type of model. Underlined numbers mean the best performances of baselines. All baseline results are from ESNS[29].

are taken from [31], [16] and official code reproduction based on same setting respectively. First, we conduct experiments with translational distance-based KGE models. The link prediction results of TransE with different negative sampling methods on the two datasets are shown in Table 2. We can observe that DeMix can effectively improve the performance of the TransE model on each dataset. Compared with the best method without considering false-negative triples, our DeMix method improves MRR by 0.9%, 7.8% on WN18RR, and FB15K237. Even compared to CANS which reduces false-negative triples with commonsense, our method improves MRR by 6.7% on FB15K237. Besides, the improvement of our method on WN18RR is not as pronounced as on FB15K237, which is consistent with the observation that false-negative triples have a lower effect on WN18RR than the degrading effect on FB15K237 in Figure 1(b). Furthermore, we conduct experiments with recently proposed KGE models such as RotatE and HAKE on two datasets. From the results shown in Table 3, our DeMix outperforms all the other negative sampling methods on FB15K237 dataset and achieves competitive results on WN18RR. Second, we conduct experiments with semantic matching KGE models. As shown in Table 4, we can observe that our method DeMix outperforms ESNS, which specifically aim to tackle the issue of false-negative triples, on both datasets incorporating DistMult and ComplEx as backbones. These results demonstrate the superiority and effectiveness of our method. The results combined with different KGE models also illustrate the generalizability of our approach.

	MRR	Hits@10	Hits@3	Hits@1
HAKE+Uniform	0.304	0.482	0.333	0.216
HAKE+DeMix-Uni	0.332	0.524	0.368	0.236
HAKE+RW-SANS	0.305	0.488	0.336	0.214
HAKE+DeMix-RW-SANS	0.322	0.515	0.353	0.228

Table 5. DeMix upon different negative sampling methods on FB15K237 with HAKE as KGE.

Models	MRR	Hits@10	Hits@3	Hits@1
DeMix	0.337	0.533	0.374	0.239
-WARM	0.336	0.534	0.374	0.237
-MPNE	0.306	0.509	0.341	0.207
-AdaMix	0.288	0.471	0.318	0.198

Table 6. Ablation study of DeMix on FB15K237 with HAKE as KGE.

4.3 Combine with other NS (Q2)

To verify whether our approach is plug-and-play, we implement DeMix upon other negative sampling methods to answer Q2. From the results shown in Table 5, our proposed method can be effectively combined with different negative sampling methods, and all of them can obtain significant improvements. Such results demonstrate that DeMix can be a plug-and-play component for existing negative sampling methods.

4.4 Ablation Study (Q3)

We conduct ablation studies to show the contribution of different modules in DeMix to answer Q3. We choose HAKE [32] model as the backbone since it has high performance. Specifically, we integrate our framework into HAKE based on the following three ablation settings: 1) removing warm up the KGE model (-WARM); 2) neglecting marginal pseudo-negative triples (-MPNE); 3) directly using the label information of noisy triples (-AdaMix). The results of ablation studies using HAKE on FB15K237 are shown in Table 6. The results show that all ablation settings lead to degraded performance. In specific, we observe that the warm-up has a slight effect on the training of the KGE model. AdaMix module is essential for model performance, indicating the adaptive mixup mechanism can take more precise supervision to the KGE model. Moreover, we also find that the performance drops significantly after removing the MPNE module, which indicates the validity of estimating marginal pseudo-negative triples.

4.5 Further Analysis

Convergence Speed First, we demonstrate the generating manner of DeMix can help the KGE model converge quickly. To ensure a clear observation, we

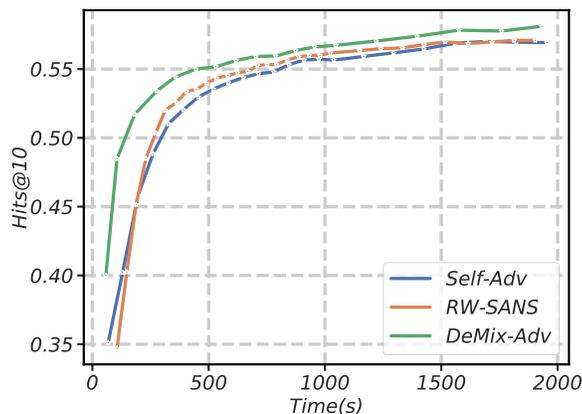


Fig. 3. Evaluating Hits@10 performance v.s. clock time (in seconds) of HAKE based on WN18RR.

compare the convergence speed of our method with other negative sampling methods using the same batch size and learning rate on WN18RR. Figure 3 shows the convergence of evaluating Hits@10 of HAKE based on WN18RR. We can observe our method help the KGE model converge more quickly compared with other searching-based methods. This means that even on WN18RR where the effect of false-negative triples is small, our method can help the model converge faster by generating high-quality triples, while not wasting training time on additional parameters.

Estimation Accuracy In order to investigate whether the MPNE module in DeMix can indeed identify false-negative triples, we inject triples from the validation set and test dataset into the MPNE module, then observe the estimation accuracy of the MPNE module for these false-negative triples as the training time increases. In specific, we calculate the estimation accuracy separately according to two patterns, i.e. (h, r) and (r, t) . Especially if a pattern does not exist in $\mathcal{T}_{pattern}$, we do not estimate this triple based on this pattern. Finally, the number of false-negative triples are 28926 and 34410 based on (h, r) pattern and (r, t) pattern respectively. As shown in Figure 4, After warming up HAKE 8 epochs, the estimation accuracy increases along the training, which implies our method can leverage the judgment power of the model itself to efficiently estimate false-negative triples.

Visualization of Entity Embeddings We visualize entities embeddings to verify the validity of our adaptive denoising mixup mechanism. In specific, we random sample an input (h, r) pattern from FB15K237, namely $(marriage, location)$, then retrieve 20 positive triples, in the training set, near the decision boundary, 10 false-negative triples, in the validation set and test set, with the lowest scores, and 10 true-negative triples with the highest scores. We visualize

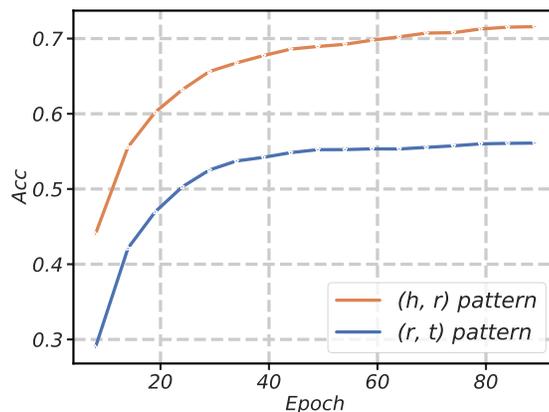


Fig. 4. The estimating accuracy of the MPNE module of HAKE based on FB15K237 with warming up 8 epochs.

tail entities embeddings of these triples. From Figure 5, we can notice that our method can push the false-negative triples closer to the positive triples and away from the true-negative triples, which indicates our adaptive denoising mixup mechanism can alleviate the noisy triples issue.

5 Related Work

KGE models The general learning approach of KGEs is to define a score function to measure the plausibility of triples. Traditional KGE Methods consist of translational distance-based models and semantic matching models. Translational distance-based models such as TransE [4] and SE[5] calculate the Euclidean distance between the relational projection of entities. Semantic matching models such as DistMult [27] and ComplEx [21] use similarity-based scoring function to measure the correctness of triples. The recent KGE approach attempt to model some of the properties present in KGs, such as RotatE [18] aims to model and infer various relation patterns by projecting entities into complex space, HAKE [32] aims to model semantic hierarchies in KGs by mapping entities into the polar coordinate system, and COMPGCN [22] uses deep neural networks to embed KGs. KGE approaches can be effectively applied to the task of knowledge graph complementation, also known as the link prediction task. However, the ability of KGE models to discriminate whether a triple is correct is often affected by the quality of the negative triples used in the training phase.

Negative Sampling for KGE The existing negative sampling techniques for KGE can be classified into two categories: 1) Negative sampling methods that ignore false-negative triples. This kind of approach is based on the closed-world assumption, where all unlabeled corrupted triples are considered as negative

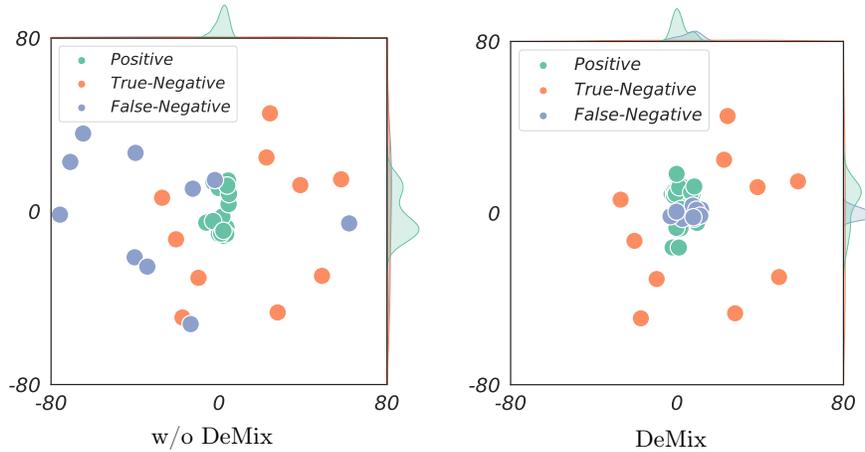


Fig. 5. 2D t-SNE visualisation of tail entities with their embeddings of positive triples near the decision boundary, false-negative triples, and true-negative triples with the same pattern (*marriage, location*).

triples, and search for high-quality negative triples. Uniform sampling[4] method samples negative triples from a uniform distribution. KBGAN [6] and IGAN [23] use the generative adversarial framework to feed the model with high-quality negative triples. Self-adversarial sampling [18] performs similarly to KBGAN, but it gives different training weights to negative triples. To achieve effectiveness and efficiency, NSCaching [31] maintains a cache containing candidates of negative triples. 2) Negative sampling considering false-negative triples. This type of method aims to minimize the chance of sampling false-negative triples, such as Bernoulli sampling [24], ESNS [29] and CANS [16]. For example, CANS leverages external commonsense information and the characteristics of complex relations to model false-negative triples’ distribution, then give lower training weights to false-negative triples following the self-adversarial negative sampling loss [18]. However, all previous negative sampling methods are searching-based methods, which inefficiently search for high-quality negative triples from a large set of unlabeled corrupted triples. Besides, CANS requires costly manual effort to gather valuable external information, and Bernoulli sampling does not use external information, but this method is a fixed sampling scheme. Our approach differs from these methods in that we do not avoid sampling false-negative triples, but rather refine negative triples to high-quality triples as much as possible.

Mixup Method Mixup [13] is a data augmentation method that generates a new instance by convex combinations of pairs of training instances. Despite its simplicity, it has shown that it can improve the generalization and Robustness of the model in many applications [19,7]. Further, a number of modified mixup versions have been proposed for supervised and unsupervised learning. For supervised learning, Mixup [13] for the first time linearly interpolates two samples and

their corresponding labels to generate virtual samples, and experimental results show that mixup is generally applicable to image, speech, and table datasets. In unsupervised scenarios, mixup method is mainly used to construct high-quality virtual negative samples to effectively improve the generalization and robustness of models. MixGCF [13] integrates multiple negative samples to synthesize a difficult negative sample by positive mixing and hop mixing to improve the performance of GNN-based recommendation models. More similar to us are methods that leverage off-the-shelf mixup methods to generate harder negative triples for KGE such as MixKG [8]. Different from existing methods, DeMix design an adaptive mixup mechanism to dynamically refine noisy negative triples.

6 Conclusion and Future Work

In this paper, we explore the denoising issue in sampling high-quality negative triples for KGEs and find these noisy triples have a significant impact on the performance of KGE. Further, we propose a novel and easily pluggable method to alleviate the denoising issue in negative sampling for KGEs. Empirical results on two benchmark datasets demonstrate the effectiveness of our approach. In the future, we plan to extend to recognize noisy triples with unseen patterns in the training set and apply active learning to our method.

Supplemental Material Statement: Source code and datasets are available for reproducing the results.³

7 Acknowledgments

This work has been supported in part by the Zhejiang NSF (LR21F020004), the NSFC (No. 62272411), Alibaba-Zhejiang University Joint Research Institute of Frontier Technologies, and Ant Group.

References

1. Ahrabian, K., Feizi, A., Salehi, Y., Hamilton, W.L., Bose, A.J.: Structure aware negative sampling in knowledge graphs. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing. pp. 6093–6101 (2020)
2. Berthelot, D., Carlini, N., Goodfellow, I., Papernot, N., Oliver, A., Raffel, C.A.: Mixmatch: A holistic approach to semi-supervised learning. In: Advances in Neural Information Processing Systems. vol. 32 (2019)
3. Bollacker, K., Evans, C., Paritosh, P., Sturge, T., Taylor, J.: Freebase: A collaboratively created graph database for structuring human knowledge. In: Proceedings of the ACM Conference on Management of Data. pp. 1247–1250 (2008)
4. Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J., Yakhnenko, O.: Translating embeddings for modeling multi-relational data. In: Proceedings of the Annual Conference on Neural Information Processing Systems. pp. 2787–2795 (2013)

³ <https://github.com/DeMix2023/Demix>

5. Bordes, A., Weston, J., Collobert, R., Bengio, Y.: Learning structured embeddings of knowledge bases. In: Proceedings of the AAAI Conference on Artificial Intelligence. p. 301–306 (2011)
6. Cai, L., Wang, W.Y.: KBGAN: adversarial learning for knowledge graph embeddings. In: NAACL-HLT. pp. 1470–1480 (2018)
7. Carratino, L., Cissé, M., Jenatton, R., Vert, J.P.: On mixup regularization. arXiv preprint arXiv:2006.06049 (2020)
8. Che, F., Yang, G., Shao, P., Zhang, D., Tao, J.: Mixkg: Mixing for harder negative samples in knowledge graph. arXiv preprint arXiv:2202.09606 (2022)
9. Chen, T.S., Hung, W.C., Tseng, H.Y., Chien, S.Y., Yang, M.H.: Incremental false negative detection for contrastive learning (2021)
10. Dettmers, T., Pasquale, M., Pontus, S., Riedel, S.: Convolutional 2d knowledge graph embeddings. In: Proceedings of the AAAI Conference on Artificial Intelligence. pp. 1811–1818 (2018)
11. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: Advances in Neural Information Processing Systems. pp. 2672–2680 (2014)
12. Hao, Y., Zhang, Y., Liu, K., He, S., Liu, Z., Wu, H., Zhao, J.: An end-to-end model for question answering over knowledge base with cross-attention combining global knowledge. In: Proceedings of Annual Meeting of the Association for Computational Linguistics. pp. 221–231 (2017)
13. Huang, T., Dong, Y., Ding, M., Yang, Z., Feng, W., Wang, X., Tang, J.: Mixgcf: An improved training method for graph neural network-based recommender systems. In: Proceedings of the ACM Knowledge Discovery and Data Mining (2021)
14. Ji, S., Pan, S., Cambria, E., Marttinen, P., Yu, P.S.: A survey on knowledge graphs: Representation, acquisition, and applications. *IEEE Transactions on Neural Networks and Learning Systems* **33**(2), 494–514 (2022)
15. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: Proceedings of the International Conference on Learning Representations (2015)
16. Niu, G., Li, B., Zhang, Y., Pu, S.: Cake: A scalable commonsense-aware framework for multi-view knowledge graph completion. In: Proceedings of the Annual Meeting of the Association for Computational Linguistics (2022)
17. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E.Z., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: Proceedings of the Annual Conference on Neural Information Processing Systems. pp. 8024–8035 (2019)
18. Sun, Z., Deng, Z., Nie, J., Tang, J.: Rotate: Knowledge graph embedding by relational rotation in complex space. In: Proceedings of the International Conference on Learning Representations (2019)
19. Thulasidasan, S., Chennupati, G., Bilmes, J.A., Bhattacharya, T., Michalak, S.: On mixup training: Improved calibration and predictive uncertainty for deep neural networks. In: Proceedings of the Annual Conference on Neural Information Processing Systems. pp. 13888–13899 (2019)
20. Toutanova, K., Chen, D., Pantel, P., Poon, H., Choudhury, P., Gamon, M.: Representing text for joint embedding of text and knowledge bases. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing. pp. 1499–1509 (2015)

21. Trouillon, T., Welbl, J., Riedel, S., Gaussier, É., Bouchard, G.: Complex embeddings for simple link prediction. In: Proceedings of the International Conference on Machine Learning. pp. 2071–2080 (2016)
22. Vashishth, S., Sanyal, S., Nitin, V., Talukdar, P.: Composition-based multi-relational graph convolutional networks. In: International Conference on Learning Representations (2020)
23. Wang, P., Li, S., Pan, R.: Incorporating GAN for negative sampling in knowledge representation learning. In: Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence. pp. 2005–2012 (2018)
24. Wang, Z., Zhang, J., Feng, J., Chen, Z.: Knowledge graph embedding by translating on hyperplanes. In: Proceedings of the Twenty-Eighth AAAI Conference. pp. 1112–1119 (2014)
25. Wang, Z., Lai, K., Li, P., Bing, L., Lam, W.: Tackling long-tailed relations and uncommon entities in knowledge graph completion. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). pp. 250–260 (2019)
26. Xiong, C., Power, R., Callan, J.: Explicit semantic ranking for academic search via knowledge graph embedding. In: Proceedings of the International World Wide Web Conferences. pp. 1271–1279 (2017)
27. Yang, B., Yih, W., He, X., Gao, J., Deng, L.: Embedding entities and relations for learning and inference in knowledge bases. In: Proceedings of the International Conference on Learning Representations (2015)
28. Yang, S., Zhang, R., Erfani, S.: GraphDialog: Integrating graph knowledge into end-to-end task-oriented dialogue systems. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 1878–1888 (2020)
29. Yao, N., Liu, Q., Li, X., Yang, Y., Bai, Q.: Entity similarity-based negative sampling for knowledge graph embedding. In: Proceedings of the 19th Pacific Rim International Conference on Artificial Intelligence, PRICAI. pp. 73–87 (2022)
30. Zhang, W., Deng, S., Wang, H., Chen, Q., Zhang, W., Chen, H.: Xtranse: Explainable knowledge graph embedding for link prediction with lifestyles in e-commerce. In: Proceedings of the Joint International Semantic Technology Conference. vol. 1157, pp. 78–87 (2019)
31. Zhang, Y., Yao, Q., Shao, Y., Chen, L.: Nscaching: Simple and efficient negative sampling for knowledge graph embedding. In: Proceedings of the IEEE International Conference on Data Engineering. pp. 614–625 (2019)
32. Zhang, Z., Cai, J., Zhang, Y., Wang, J.: Learning hierarchy-aware knowledge graph embeddings for link prediction. In: Thirty-Fourth AAAI Conference on Artificial Intelligence. pp. 3065–3072 (2020)