

Noisy Positive-Unlabeled Learning with Self-Training for Speculative Knowledge Graph Reasoning

Anonymous ACL submission

Abstract

This paper studies speculative reasoning task on real-world knowledge graphs (KG) that contain both *false negative issue* (i.e., potential true facts being excluded) and *false positive issue* (i.e., unreliable or outdated facts being included). State-of-the-art methods fall short in the speculative reasoning ability, as they assume the correctness of a fact is solely determined by its presence in KG, making them vulnerable to false negative/positive issues. The new reasoning task is formulated as a noisy Positive-Unlabeled learning problem. We propose a variational framework, namely nPUGraph, that jointly estimates the correctness of both collected and uncollected facts (which we call *label posterior*) and updates model parameters during training. The label posterior estimation facilitates speculative reasoning from two perspectives. First, it improves the robustness of a label posterior-aware graph encoder against false positive links. Second, it identifies missing facts to provide high-quality grounds of reasoning. They are unified in a simple yet effective self-training procedure. Empirically, extensive experiments on three benchmark KG and one Twitter dataset with various degrees of false negative/positive cases demonstrate the effectiveness of nPUGraph.

1 Introduction

Knowledge graphs (KG), which store real-world facts in triples (*head entity, relation, tail entity*), have facilitated a wide spectrum of knowledge-intensive applications (Saxena et al., 2021; Qian et al., 2019; Wang et al., 2018, 2022). Automatically reasoning facts based on observed ones, a.k.a. Knowledge Graph Reasoning (KGR) (Bordes et al., 2013), becomes increasingly vital since it allows for expansion of the existing KG at a low cost.

Numerous efforts have been devoted to KGR task (Bordes et al., 2013; Lin et al., 2015; Trouillon et al., 2017; Sun et al., 2019; Li et al., 2021), which assume the correctness of a fact is solely

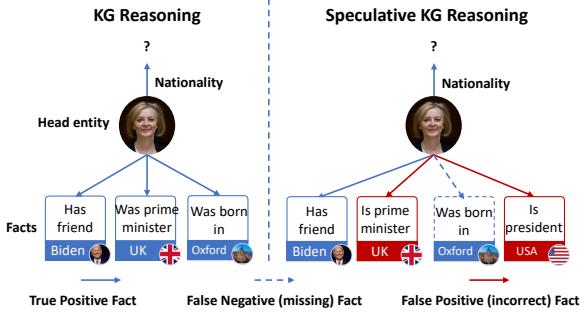


Figure 1: An illustrative example of speculative KG reasoning. Blue solid lines denote the true positive fact, blue dashed lines denote the *false negative* (missing) fact, and red solid lines denote the *false positive* (incorrect) fact. We aim to mitigate the false negative/positive issues to enable speculative KG reasoning.

determined by its presence in KG. They ideally view facts included in KG as positive samples and excluded facts as negative samples. However, most real-world reasoning has to be performed based on sparse and unreliable observations, where there may be true facts excluded or false facts included. Reasoning facts based on sparse and unreliable observations (which we call *speculative KG reasoning*) are still underexplored.

In this paper, we aim to enable the speculative reasoning ability on real-world KG. The fulfillment of the goal needs to address two commonly existing issues, as shown in Figure 1: 1) **The false negative issue** (i.e., sparse observation): Due to the graph incompleteness, facts excluded from the KG can be used as implicit grounds of reasoning. This is particularly applicable to non-obvious facts. For example, personal information such as the birthplace of politicians may be missing when constructing a political KG, as they are not explicitly stated in the political corpus (Tang et al., 2022). However, it can be critical while reasoning personal facts like nationality. 2) **The false positive issue** (i.e., noisy observation): Facts included in the KG may be unreliable and should not be directly grounded without

068
069
070
071
072
073
074
075
076
077
078
079

inspection. It can happen when relations between entities are incorrectly collected or when facts are extracted from outdated or unreliable sources. For example, Mary Elizabeth is no longer the Prime Minister of the United Kingdom, which may affect the reasoning accuracy of her current workplace. These issues generally affect both one-hop reasoning (Bordes et al., 2013) and multi-hop reasoning (Saxena et al., 2021). The main focus of this paper is investigating the one-hop speculative reasoning task as it lays the basis for complicated multi-hop reasoning capability.

080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102

Speculative KG reasoning differs from conventional KG reasoning in that the correctness of each collected/uncollected fact needs to be dynamically estimated as part of the learning process, such that the grounds of reasoning can be accordingly calibrated. Unfortunately, most existing work, if not all, lacks such inspection capability. Knowledge graph embedding methods (Bordes et al., 2013; Lin et al., 2015; Yang et al., 2014; Trouillon et al., 2017; Sun et al., 2019) and graph neural network (GNN) methods (Schlichtkrull et al., 2017; Dettmers et al., 2018; Nguyen et al., 2018; Vashishth et al., 2020; Li et al., 2021) can easily overfit the false negative/positive cases because of their training objective that ranks the collected facts higher than other uncollected facts in terms of plausibility. Recent attempts on uncertain KG (Chen et al., 2019; Kertkeidkachorn et al., 2019) measure the uncertainty scores for facts, which can be utilized to detect false negative/positive samples. However, they explicitly require the ground truth uncertainty scores as supervision for reasoning model training, which are usually unavailable in practice.

103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118

Motivated by these observations, we formulate the speculative KG reasoning task as a noisy Positive-Unlabeled learning problem. The facts contained in the KG are seen as noisy positive samples with a certain level of label noise, and the facts excluded from the KG are treated as unlabeled samples, which include both negative ones and possible factual ones. Instead of determining the correctness of facts before training the reasoning model without inspection, we learn the two perspectives in an end-to-end training process. To this end, we propose nPUGraph, a novel variational framework that regards the underlying correctness of collected/uncollected facts in the KG as latent variables for the reasoning process. We jointly update model parameters and estimate the

119
120
121
122
123

posterior likelihood of the correctness of each collected/uncollected fact (referred to as *label posterior*), through maximizing a theoretical lower bound of the log-likelihood of each fact being collected or uncollected.

124
125
126
127
128
129
130

The estimated label posterior further facilitates the speculative KG reasoning from two aspects: 1) It removes false positive facts contained in KG and improves the representation quality. We accordingly propose a label posterior-aware encoder to incorporate information only from entity neighbors induced by facts with a high posterior probability, under the assumption that the true positive facts from the collected facts provide more reliable information for reasoning. 2) It complements the grounds of reasoning by selecting missing but possibly plausible facts with high label posterior, which are iteratively added to acquire more informative samples for model training. These two procedures are ultimately unified in a simple yet effective self-training strategy that alternates between the *data sampling based on latest label posteriors* and the *model training based on latest data samples*. Empirically, nPUGraph outperforms eleven state-of-the-art baselines on three benchmark KG data and one Twitter data we collected by large margins. Additionally, its robustness is demonstrated in speculative reasoning on data with multiple ratios of false negative/positive cases.

131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147

Our contributions are summarized as follows: (1) We open up a practical but underexplored problem space of speculative KG reasoning, and formulate it as a noisy Positive-Unlabeled learning task; (2) We take the first step in tackling this problem by proposing a variational framework nPUGraph to jointly optimize reasoning model parameters and estimate fact label posteriors; (3) We propose a simple yet effective self-training strategy for nPUGraph to simultaneously deal with false negative/positive issues; (4) We perform extensive evaluations to verify the effectiveness of nPUGraph on both benchmark KG and Twitter interaction data with a wide range of data perturbations.

2 Preliminaries

2.1 Speculative Knowledge Graph Reasoning

163
164
165
166
167
168

A knowledge graph (KG) is denoted as $\mathcal{G} = \{(e_h, r, e_t)\} \subseteq \mathcal{S}$, where $\mathcal{S} = \mathcal{E} \times \mathcal{R} \times \mathcal{E}$ denotes triple space, \mathcal{E} denotes the entity set, \mathcal{R} denotes the relation set. Each triple $s = (e_h, r, e_t)$ refers to that a head entity $e_h \in \mathcal{E}$ has a relation $r \in \mathcal{R}$

169
170
171
172
173
174
with a tail entity $e_t \in \mathcal{E}$. The goal of speculative
KG reasoning is to infer the most plausible triple
for each incomplete triple $(e_h, r, e_?)$ or $(e_?, r, e_t)$
given by sparse and unreliable observations in \mathcal{G} .
In addition, it requires correctness estimation for
each potential fact collected or uncollected by \mathcal{G} .

175 2.2 Noisy Positive-Unlabeled Learning

176 Positive-Unlabeled (PU) learning is a learning
177 paradigm for training a model when only positive
178 and unlabeled data is available (Plessis et al., 2015).
179 We formulate the speculative KG reasoning task
180 as a noisy Positive-Unlabeled learning problem,
181 where the positive set contains potentially label
182 noise from false facts (Jain et al., 2016).

183 **PU Triple Distribution.** For the speculative KG
184 reasoning task, we aim to learn a binary classi-
185 fier that maps a triple space \mathcal{S} to a label space
186 $\mathcal{Y} = \{0, 1\}$. Data are split as labeled (collected)¹
187 triples $s^l \in \mathcal{S}^L$ and unlabeled (uncollected) triples
188 $s^u \in \mathcal{S}^U$. The labeled triples are considered noisy
189 positive samples with a certain level of label noise.
190 The distribution of labeled triples can be re-
191 presented as follows:

$$192 s^l \sim \beta \phi_1^l(s^l) + (1 - \beta) \phi_0^l(s^l), \quad (1)$$

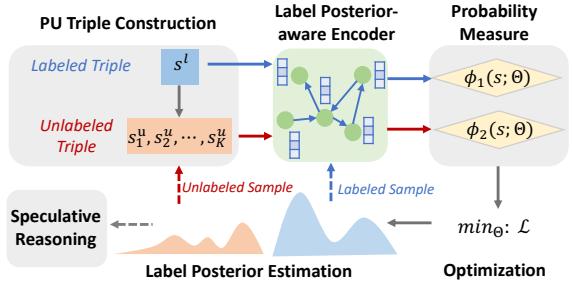
193 where ϕ_y^l denotes the probability of being collected
194 over triple space \mathcal{S} for the positive class ($y = 1$)
195 and negative class ($y = 0$), and $\beta \in [0, 1]$ denotes
196 the proportion of true positive samples in labeled
197 data. Unlabeled triples include both negative sam-
198 ples and possible factual samples. The distribution
199 of unlabeled samples can be represented as follows:

$$200 s^u \sim \alpha \phi_1^u(s^u) + (1 - \alpha) \phi_0^u(s^u), \quad (2)$$

201 where $\phi_y^u = 1 - \phi_y^l$ denotes the probability of being
202 uncollected, $\alpha \in [0, 1]$ denotes the positive class
203 prior, i.e., the proportion of positive samples in
204 unlabeled data.

205 **PU Triple Construction.** We then discuss the con-
206 struction of \mathcal{S}^L and \mathcal{S}^U based on the collected KG
207 \mathcal{G} . Triples in \mathcal{G} naturally serve as labeled sam-
208 ples with a ratio of noise, i.e., $\mathcal{S}^L = \mathcal{G}$. For unlabeled set \mathcal{S}^U , However, directly using $\mathcal{S} \setminus \mathcal{G}$ as
209 the unlabeled set \mathcal{S}^U would result in too many
210 unlabeled samples for training due to the large
211 number of possible triples in triple space \mathcal{S} . Fol-
212 lowing (Tang et al., 2022), we construct \mathcal{S}^U as
213 follows: For each labeled triple $s_i^l = (e_h, r, e_t)$,

214 ¹In this paper, we interchangeably use the term *la-
215 beled/unlabeled* and *collected/uncollected* with no distinction.



216
217
218
219
220
221
222
223
Figure 2: nPUGraph overview. It jointly optimizes
parameters and estimates label posterior, to detect false
negative/positive cases for the encoder and self-training.

224 we construct K unlabeled triples s_{ik}^u by replacing
225 the head and tail respectively with other entities:
226 $s_{ik}^u = (e_h, r, e_k^-)$ or (e_k^-, r, e_t) , where e_k^- is the
227 selected entity that ensures $s_{ik}^u \notin \mathcal{S}^L$. Initially,
228 the construction can be randomized. During the
229 training process, it is further improved by select-
230 ing unlabeled samples with high label posterior
231 in a self-training scheme, so as to cover positive
232 samples in the unlabeled set to the greatest extent.

233 3 Methodology

234 3.1 Overview

235 Our approach views underlying triple labels (pos-
236 itive/negative) as latent variables, influencing the
237 collection probability. Unlike the common objec-
238 tive of reasoning training that ranks the plausibil-
239 ity of the collected triples higher than uncollected
240 ones, we instead maximize the data likelihood of
241 each potential triple being collected or not. To
242 this end, as shown in Figure 2, we propose nPUGraph
243 framework to jointly optimize parameters and
244 infer the label posterior. During the training
245 process, the latest label posterior estimation can be
246 utilized by a label posterior-aware encoder, which
247 improves the quality of representation learning by
248 only integrating information from the entity neigh-
249 bors induced by true facts. Finally, a simple yet
250 effective self-training strategy based on label pos-
251 terior is proposed, which can dynamically update
252 neighbor sets for the encoder and sample unlabeled
253 triplets to cover positive samples in the unlabeled
254 set to the greatest extent for model training.

255 The remaining of this section is structured as
256 follows: Section 3.2 first formalizes the learning
257 objective and the variational framework for likeli-
258 hood maximization. Section 3.3 details the label
259 posterior-aware encoder for representation learn-
260 ing, followed by Section 3.4 that introduces the

252 self-training strategy.

253 3.2 Noisy PU Learning on KG

254 Due to the false negative/positive issues, the correctness of a fact (y) is not solely determined by
 255 its presence in a knowledge graph. nPUGraph
 256 addresses the issue by treating the underlying label as
 257 a latent variable that influences the probability of
 258 being collected or not. We, therefore, set maximizing
 259 the data collection likelihood as our objective.
 260 In such a learning paradigm, the assumptions that
 261 collected triples are correct $p(y = 1|s^l) = 1$ and
 262 uncollected triples are incorrect $p(y = 0|s^u) = 1$
 263 are removed. We aim to train a model on labeled
 264 triples \mathcal{S}^L and unlabeled ones \mathcal{S}^U , and infer the label
 265 posterior $p(y|s^u)$ and $p(y|s^l)$ at the same time
 266 by data likelihood maximization. The latest label
 267 posterior can help to detect false negative/positive
 268 cases during model training.

269 We first derive our training objective. To be
 270 more formal, the log-likelihood of each potential
 271 fact being collected or not is lower bounded by
 272 Eq. (3), which is given by Theorem 1.

273 **Theorem 1.** *The log-likelihood of the complete
 274 data $\log p(\mathbf{S})$ is lower bounded as follows:*

$$\begin{aligned} \log p(\mathbf{S}) &\geq \mathbb{E}_{q(\mathbf{Y})} [\log p(\mathbf{S}|\mathbf{Y})] - \text{KL}(q(\mathbf{Y})\|p(\mathbf{Y})) \\ &= \mathbb{E}_{s^l \in \mathcal{S}^L} [w^l \log[\phi_1^l(s^l)] + (1-w^l) \log[\phi_0^l(s^l)]] \\ &+ \mathbb{E}_{s^u \in \mathcal{S}^U} [(w^u \log[\phi_1^u(s^u)] + (1-w^u) \log[\phi_0^u(s^u)])] \\ &- \text{KL}(\mathbf{W}^U\|\tilde{\mathbf{W}}^U) - \text{KL}(\mathbf{W}^L\|\tilde{\mathbf{W}}^L) - \frac{\|\mathbf{W}^L\|_1}{|\mathcal{S}^L|} - \frac{\|\mathbf{W}^U\|_1}{|\mathcal{S}^U|}, \end{aligned} \quad (3)$$

275 where \mathbf{S} denotes all labeled/unlabeled triples, \mathbf{Y}
 276 is the corresponding latent variable indicating the
 277 positive/negative labels for triples, $\mathbf{W}^U = \{w_i^u\}$
 278 denotes the point-wise probability for the uncollected
 279 triples being positive, $\mathbf{W}^L = \{w_i^l\}$ denotes
 280 the probability for the collected triples being positive.
 281 $\tilde{\mathbf{W}}^U$ and $\tilde{\mathbf{W}}^L$ are the approximation of
 282 the collection probability for uncollected/collected
 283 triples respectively, produced by nPUGraph based
 284 on the latest parameters.

285 *Proof.* Refer to Appendix A.1 for proof. \square

286 We treat label \mathbf{Y} as a latent variable and derive
 287 the lower bound for the log-likelihood, which is
 288 influenced by the prior knowledge of positive class
 289 prior α and true positive ratio β . Thus, maximizing
 290 the lower bound can jointly optimize model param-
 291 eters and infer the posterior label distribution, \mathbf{W}^U

292 and \mathbf{W}^L . Such a learning process enables us to
 293 avoid false negative/positive issues during model
 294 training since it considers ϕ_0^l (one negative triple
 295 is collected) and ϕ_1^u (one positive triple is missing)
 296 as non-zero probability, which are determined by
 297 the latest label posterior during model training.

298 **Probability Measure.** We then specify the proba-
 299 bility measures for positive/negative triples being
 300 collected, i.e., $\phi_1^l(\cdot)$ and $\phi_0^l(\cdot)$ ($\phi_y^u(\cdot) = 1 - \phi_y^l(\cdot)$
 301 for $y = 1/0$). To better connect to other meth-
 302 ods utilizing score functions for KGR, we hereby
 303 utilize the sigmoid function $\sigma(\cdot)$ to directly trans-
 304 form the score function $\psi(s; \Theta)$, parameterized by
 305 model parameters Θ , to probability:

$$\phi_1^l(s) = \sigma(\psi_1(s; \Theta)), \quad \phi_0^l(s) = \sigma(\psi_0(s; \Theta)), \quad (4)$$

340 we hereby utilize two score functions $\psi_1(s; \Theta)$ and
 341 $\psi_0(s; \Theta)$ to measure the positive/negative triples
 342 being collected, as the influencing factors based
 343 on triple information can be different. We utilize
 344 two neural networks to approximate the probability
 345 measure, which will be detailed in Section 3.3.

346 Since we aim to detect the potential existence of
 347 positive triples in an unlabeled set, it is unnecessary
 348 to push the collection probability of all uncollected
 349 triple $\phi_y^l(s^u)$ to 0 ($\phi_y^u(s^u)$ to 1). A loose constraint
 350 is that we force the uncollection probability of a
 351 collected triple s^l lower than its corresponding un-
 352 collected triples s^u : $\phi_y^u(s^l) < \phi_y^u(s^u)$. Therefore,
 353 we adopt the pair-wise ranking measure $\phi_y^*(s^l, s^u)$
 354 to replace $\phi_y^u(s^u)$ as follows:

$$\phi_y^u(s^u) \rightarrow \phi_y^*(s^l, s^u) = \sigma(\psi_y(s^u; \Theta) - \psi_y(s^l; \Theta)). \quad (5)$$

355 **Maximum Probability Training.** We then derive
 356 the training objective based on Eq. (3). The first
 357 part of Eq. (3) measures the probability of data be-
 358 ing collected/uncollected. Concretely, given each
 359 collected triple $s_i^l \in \mathcal{S}^L$ and its corresponding K
 360 uncollected triples $s_{ik}^u \in \mathcal{S}^U$, We denote the loss
 361 function measuring the probability as \mathcal{L}_{triple} :

$$\begin{aligned} \mathcal{L}_{triple} &= -\frac{1}{K|\mathcal{S}^L|} \sum_i \sum_k (w_i^l \log[\phi_1^l(s_i^l)] \\ &+ (1-w_i^l) \log[\phi_0^l(s_i^l)] + w_{ik}^u \log[\phi_1^*(s_i^l, s_{ik}^u)] \\ &+ (1-w_{ik}^u) \log[\phi_0^*(s_i^l, s_{ik}^u)]), \end{aligned} \quad (6)$$

362 where w_i^l denotes the point-wise probability for
 363 the collected triple s_i^l being positive, w_{ik}^u denotes
 364 the probability for the uncollected triple s_{ik}^u being

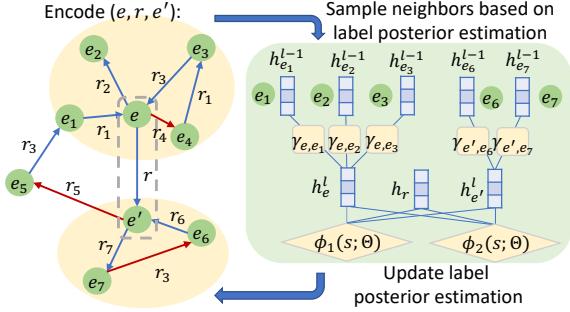


Figure 3: The label posterior-aware encoder. Red line denotes the detected false positive facts based on posterior, which are excluded during neighbor sampling.

positive. Based on the definition, the posterior probability of each collected/uncollected triple being positive can be computed as:

$$\tilde{w}_i^l = \frac{\beta\phi_1^l(s_i^l)}{\beta\phi_1^l(s_i^l) + (1-\beta)\phi_0^l(s_i^l)}, \quad (7)$$

$$\tilde{w}_{ik}^u = \frac{\alpha\phi_1^u(s_{ik}^u)}{\alpha\phi_1^u(s_{ik}^u) + (1-\alpha)\phi_0^u(s_{ik}^u)}. \quad (8)$$

To increase model expression ability, instead of forcing $\mathbf{W}^L = \tilde{\mathbf{W}}^L$ and $\mathbf{W}^U = \tilde{\mathbf{W}}^U$, we set \mathbf{W}^L and \mathbf{W}^U as free parameters and utilize the term $\mathcal{L}_{KL} = \text{KL}(\mathbf{W}^L \parallel \tilde{\mathbf{W}}^L) + \text{KL}(\mathbf{W}^U \parallel \tilde{\mathbf{W}}^U)$ to regularize the difference. Finally, based on Eq. (3), the training objective is formalized as follows:

$$\min_{\Theta} \mathcal{L} = \min_{\Theta} \mathcal{L}_{\text{triple}} + \mathcal{L}_{KL} + \mathcal{L}_{reg}, \quad (9)$$

where $\mathcal{L}_{reg} = \|\mathbf{W}^L\|_1 + \|\mathbf{W}^U\|_1$ can be viewed as a normalization term. Considering the sparsity property of real-world graphs, \mathcal{L}_{reg} penalizes the posterior estimation that there are too many true positive facts on KG.

3.3 Label Posterior-aware Encoder

We then introduce the encoder and the score functions $\psi_1(s; \Theta)$ and $\psi_0(s; \Theta)$ to measure the probability of positive/negative triples being collected, as shown in Figure 3. Recent work (Schlichtkrull et al., 2017; Dettmers et al., 2018; Nguyen et al., 2018; Vashishth et al., 2020) has shown that integrating information from neighbors to represent entities engenders better reasoning performance. However, the message-passing mechanism is vulnerable to the false positive issue, as noise can be integrated via a link induced by a false positive fact. In light of this, we propose a label posterior-aware encoder to improve the quality of representations.

We represent each entity $e \in \mathcal{E}$ and each relation $r \in \mathcal{R}$ into a d -dimensional latent space:

$\mathbf{h}_e, \mathbf{h}_r \in \mathbb{R}^d$. To encode more information in \mathbf{h}_e , we first construct a neighbor set \mathcal{N}_e induced by the positive facts related to entity e . The latest label posterior for collected facts $\tilde{\mathbf{W}}^L$ naturally serves this purpose, as it indicates the underlying correctness for each collected fact.

Therefore, for each entity e , we first sort the related facts by label posterior $\tilde{\mathbf{W}}^L$ and construct the neighbor set $\mathcal{N}_e(\tilde{\mathbf{W}}^L) = \{(e_i, r_i)\}$ from the top facts. Then the encoder attentively aggregates information from the collected neighbors, where the attention weights take neighbor features, relation features into account. Specifically:

$$\mathbf{h}_e^l = \mathbf{h}_e^{l-1} + \sigma \left(\sum_{(e_i, r_i) \in \mathcal{N}_e(\tilde{\mathbf{W}}^L)} \gamma_{e,e_i}^l (\mathbf{h}_{e_i}^{l-1} \mathbf{M}) \right), \quad (10)$$

where l denotes the layer number, $\sigma(\cdot)$ denotes the activation function, γ_{e,e_i}^l denotes the attention weight of entity e_i to the represented entity e , and \mathbf{M} is the trainable transformation matrix. The attention weight γ_{e,e_i}^l is supposed to be aware of entity feature and topology feature induced by relations. We design the attention weight γ_{e,e_i}^l as follows:

$$\gamma_{e,e_i}^l = \frac{\exp(q_{e,e_i}^l)}{\sum_{\mathcal{N}_e(\tilde{\mathbf{W}}^L)} \exp(q_{e,e_k}^l)}, \quad q_{e,e_k}^l = \mathbf{a} (\mathbf{h}_e^{l-1} \parallel \mathbf{h}_{e_k}^{l-1} \parallel \mathbf{h}_{r_k}), \quad (11)$$

where q_{e,e_k}^l measures the pairwise importance from neighbor e_k by considering the entity embedding, neighbor embedding, and relation embedding, $\mathbf{a} \in \mathbb{R}^{3d}$ is a shared parameter in the attention.

To measure the collection probability for positive/negative triples, we utilize two multilayer perception (MLP) to approximate score function $\psi_1(s; \Theta)$ and $\psi_0(s; \Theta)$. Specifically, for each triple $s = (e_h, r, e_t)$:

$$\psi_1(s; \Theta) = \text{MLP}_1(\mathbf{h}_s), \quad \psi_0(s; \Theta) = \text{MLP}_0(\mathbf{h}_s), \quad (12)$$

where the MLP input $\mathbf{h}_s = [\mathbf{h}_{e_h}^l \parallel \mathbf{h}_r \parallel \mathbf{h}_{e_t}^l]$ concatenates entity embeddings and relation embedding.

3.4 Self-Training Strategy

The latest label posterior $\tilde{\mathbf{W}}^L$ and $\tilde{\mathbf{W}}^U$ is further utilized in a self-training strategy to enhance speculative reasoning. First, the latest posterior estimation $\tilde{\mathbf{W}}^L$ for collected links updates neighbor sets to gradually prevent the encoder effects by false positive links. Moreover, the latest estimation $\tilde{\mathbf{W}}^L$ for uncollected facts enables us to continuously sample unlabeled triplets with high label posterior to cover positive samples in the unlabeled set to

Algorithm 1: Summary of nPUGraph.

Input: The collected triple set \mathcal{S}^L .
Output: The model parameter Θ , predicted triples.
1 Construct the uncollected triple set \mathcal{S}^U randomly;
2 Initialize the model parameter Θ and the label posterior $\tilde{\mathbf{W}}^L$ and $\tilde{\mathbf{W}}^U$ randomly;
3 **for** each training epoch **do**
4 Construct neighbor set $\mathcal{N}_e(\tilde{\mathbf{W}}^L)$ by $\tilde{\mathbf{W}}^L$;
5 Construct uncollected triple set \mathcal{S}^U by $\tilde{\mathbf{W}}^U$;
6 **for** each collected triple $s_i^l \in \mathcal{S}^L$ **do**
7 Collect unlabeled triples $\{s_{ik}^u\}_{k=1}^K$;
8 Calculate $\phi_y^l(s_i^l)$ by Eq. (4);
9 Calculate each $\phi_y^*(s_i^l, s_{ik}^u)$ by Eq. (5);
10 **end**
11 Calculate the total loss \mathcal{L} by Eq. (9);
12 Optimize model parameter: $\Theta = \Theta - \frac{\partial \mathcal{L}}{\partial \Theta}$;
13 Update label posterior $\tilde{\mathbf{W}}^L$ and $\tilde{\mathbf{W}}^U$;
14 **end**

414 the greatest extent. For each labeled triple $s_i^l = (e_h, r, e_t)$, we construct K unlabeled triples s_{ik}^u by
415 replacing the head and tail respectively with other
416 entities: $s_{ik}^u = (e_h, r, e_k^-)$ or (e_k^-, r, e_t) , where
417 e_k^- is the selected entity that ensures $s_{ik}^u \notin \mathcal{S}^L$.
418 Such selection is performed by ranking the corresponding
419 label posterior \tilde{w}_{ik}^u . The updates of
420 neighbor sets and unlabeled triples based on label
421 posterior are nested with parameter optimization
422 during model training alternatively. The training of
423 nPUGraph is summarized in Algorithm 1.
424

4 Experiment

4.1 Experimental Setup

427 **Dataset.** We evaluate nPUGraph mainly on
428 three benchmark datasets: FB15K (Bordes et al.,
429 2013), FB15k-237 (Toutanova et al., 2015), and
430 WN18 (Bordes et al., 2013) and one Twitter data
431 we collected, which describes user interaction
432 information towards tweets and hashtags. Table 3 in
433 Appendix A.2 summarizes the dataset statistics.
434

435 To better fit the real scenario for speculative
436 reasoning, we randomly modify links on KG to simu-
437 late more false negative/positive cases. We modify
438 a specific amount of positive/negative links (the
439 ratio of the modified links is defined as perturbation
440 rate, i.e., ptb_rate) by flipping. 90% of them are the
441 removed positive links to simulate false negative
442 cases and the remaining 10% are the added neg-
443 ative links to simulate false positive cases. More
444 details about datasets and the data perturbation
445 process can be found in Appendix A.2.
446

Baselines. We compare to eleven state-of-the-art baselines: 1) KG embedding methods:

447 **TransE** (Bordes et al., 2013), **TransR** (Lin
448 et al., 2015), **DistMult** (Yang et al., 2014),
449 **ComplEX** (Trouillon et al., 2017), and **Rot-
450 ate** (Sun et al., 2019); 2) GNN methods on
451 KG: **RGCN** (Schlichtkrull et al., 2017) and
452 **CompGCN** (Vashisht et al., 2020); 3) Uncertain
453 KG reasoning: **UKGE** (Chen et al., 2019); 4) Neg-
454 ative sampling methods: **NSCaching** (Zhang et al.,
455 2019) and **SANS** (Ahrabian et al., 2020); 5) PU
456 learning on KG: **PUDA** (Tang et al., 2022). More
457 details can be found in Appendix A.3.

458 **Evaluation and Implementation.** For each
459 (e_h, r, e_t) or (e_t, r, e_t) , we rank all entities at
460 the missing position in triples, and adopt filtered
461 mean reciprocal rank (*MRR*) and filtered Hits at
462 $\{1, 3, 10\}$ as evaluation metrics (Bordes et al.,
463 2013). More implementation details of baselines
464 and nPUGraph can be found in Appendix A.4.

4.2 Main Results

465 We first discuss the model performance on noisy
466 and incomplete graphs, with $ptb_rate = 0.3$, as
467 shown in Table 1. nPUGraph achieves consistently
468 better results than all baseline models, with 10.3%
469 relative improvement on average. Specifically,
470 conventional KGE and GNN-based methods produce
471 unsatisfying performance, as they ignore the false
472 negative/positive issues during model training. In
473 some cases, GNN-based ways are worse, as the
474 message-passing mechanism is more vulnerable
475 to false positive links. As expected, the perfor-
476 mance of uncertain knowledge graph embedding
477 model (Chen et al., 2019) is much worse when there
478 are no available uncertainty scores for model train-
479 ing. SANS and PUDA generate competitive results
480 in some cases, as their negative sampling strategy
481 and PU learning objective can respectively mitigate
482 the false negative/positive issues to some extent.
483 Table 1 demonstrates the superiority of nPUGraph
484 which addresses the false negative/positive issues
485 simultaneously. For space limitations, we report
486 and discuss the model performance on Twitter data
487 in Table 4 in Appendix A.5.1.

4.3 Experiments under Various Degrees of Noise and Incompleteness

491 We investigate the performance of baseline mod-
492 els and nPUGraph under different degrees of noise
493 and incompleteness. Figure 4 reports the per-
494 formance under various ptb_rate , from 0.1 to 0.7,
495 where higher ptb_rate means more links are per-
496 turbed as false positive/negative cases. Full re-

Table 1: Overall performance on noisy and incomplete graphs, with $ptb_rate = 0.3$. Average results on 5 independent runs are reported. * indicates the statistically significant results over baselines, with $p\text{-value} < 0.01$. The best results are in boldface, and the strongest baseline performance is underlined.

Dataset	FB15K				FB15K-237				WN18			
	Metrics	MRR	H@10	H@3	H@1	MRR	H@10	H@3	H@1	MRR	H@10	H@3
<i>Knowledge graph embedding methods</i>												
TransE	0.336	0.603	0.425	0.189	0.196	0.394	0.236	0.094	0.229	0.481	0.416	0.030
TransR	0.314	0.579	0.397	0.170	0.184	0.359	0.211	0.098	0.229	0.480	0.408	0.035
DistMult	0.408	0.627	0.463	0.296	0.240	0.407	0.262	0.158	0.397	0.518	0.453	0.320
ComplEx	0.396	0.616	0.451	0.284	0.238	0.411	0.262	0.154	0.448	0.526	0.475	0.403
RotatE	0.431	0.636	0.489	0.323	0.255	0.433	0.280	0.169	0.446	0.524	0.474	0.400
<i>Graph neural network methods on KG</i>												
RGCN	0.154	0.307	0.164	0.078	0.141	0.276	0.145	0.075	0.362	0.464	0.412	0.300
CompGCN	0.409	0.631	0.465	0.294	0.253	0.422	0.275	0.171	0.445	0.522	0.471	0.400
<i>Uncertain knowledge graph embedding method</i>												
UKGE	0.311	0.556	0.337	0.189	0.172	0.233	0.128	0.081	0.241	0.447	0.309	0.119
<i>Negative sampling methods</i>												
NSCaching	0.371	0.576	0.424	0.265	0.190	0.329	0.208	0.121	0.306	0.401	0.334	0.255
SANS	0.372	0.599	0.434	0.252	0.243	0.416	0.267	0.158	0.453	0.528	0.479	0.409
<i>Positive-Unlabeled learning methods on KG</i>												
PUDA	0.403	0.623	0.458	0.291	0.234	0.394	0.255	0.156	0.382	0.499	0.444	0.306
nPUGraph	0.486*	0.718*	0.534*	0.342*	0.287*	0.481*	0.315*	0.191*	0.493*	0.582*	0.519*	0.442*
Gains (%)	12.7	12.8	9.2	5.9	12.6	11.2	12.5	11.4	8.9	10.3	8.3	8.0

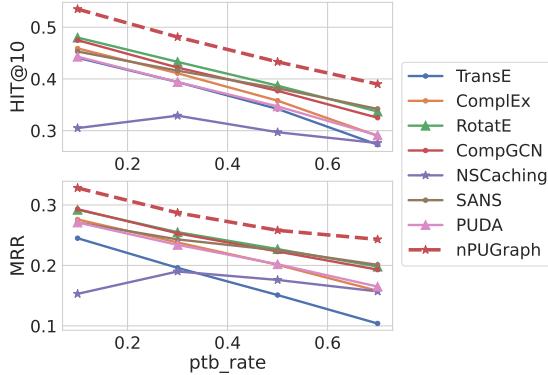


Figure 4: Performance with respect to various ptb_rate , i.e., different degrees of noise and incompleteness, on FB15K-237. nPUGraph exhibits impressive robustness against false negative/positive issues.

sults are included in Appendix A.5.2. The performance degrades as the ptb_rate increases for all in most cases, demonstrating that the false negative/positive issues significantly affect the reasoning performance. However, nPUGraph manages to achieve the best performance in all cases. Notably, the relative improvements are more significant under higher ptb_rate .

4.4 Model Analysis

Ablation Study. We evaluate performance improvements brought by the nPUGraph framework by following ablations: 1) **nPUGraph w/o nPU** is trained without the noisy Positive-Unlabeled framework, which instead utilizes the margin loss for model training; 2) **nPUGraph w/o LP-Encoder**

Table 2: Ablation Studies.

Dataset	FB15K		FB15K-237		Gains	
	Ablations	MRR	H@10	MRR	H@10	
nPUGraph w/o nPU		0.401	0.619	0.230	0.407	-20.0
nPUGraph w/o LP-Encoder		0.457	0.681	0.261	0.459	-6.6
nPUGraph w/o Self-Training		0.471	0.704	0.276	0.461	-3.4
nPUGraph		0.486	0.718	0.287	0.481	-

eliminates the label posterior-aware encoder (LP-Encoder), which aggregates information from all neighbors instead of the sampled neighbors; 3) **nPUGraph w/o Self-Training** is trained without the proposed self-training algorithm.

We report *MRR* and *Hit@10* over FB15K and FB15K-237 data, as shown in Table 2. As we can see, training the encoder without the proposed noisy Positive-Unlabeled framework will cause the performance drop, as this variant ignores the false negative/positive issues. Removing the label posterior-based neighbor sampling in the encoder will also cause performance degradation, as the information aggregation no longer distinguishes between true and false links. Such a variant can be easily influenced by the existence of false positive facts. Moreover, the last ablation result shows if the training process is further equipped with the self-training strategy, the performance will be enhanced, which verifies its effectiveness to select informative unlabeled samples for model training.

The Effect of PU Triple Construction. We then investigate the effect of PU Triple Construction on model performance, by varying different sizes of unlabeled samples from 10 to 50. Figure 5a shows that the performance improves as the num-

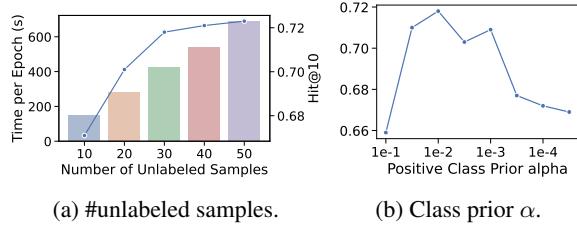


Figure 5: Model analysis w.r.t. the number of unlabeled samples and positive class prior α .

ber of unlabeled samples increases. Because more unlabeled samples can cover more false negative cases for model training. The training time grows linearly.

The Effect of Positive Class Prior α . The positive class prior α and true positive ratio β are two important hyperparameters. While β has a clear definition from real-world data, the specific value of α is unknown in advance. Figure 5b shows the model performance w.r.t. different values of α by grid search. Reasoning performance fluctuates a bit with different values of α since incorrect prior knowledge of α can bias the label posterior estimation and thus hurt the performance.

5 Related Work

Knowledge Graph Reasoning. Knowledge graph reasoning (KGR) aims to predict missing facts to automatically complete KG, including one-hop reasoning (Bordes et al., 2013) and multi-hop reasoning (Saxena et al., 2021). To set the scope, we primarily focus on one-hop reasoning and are particularly interested in predicting missing entities in a partial fact. Knowledge graph embeddings achieve state-of-the-art performance (Bordes et al., 2013; Lin et al., 2015; Yang et al., 2014; Trouillon et al., 2017; Sun et al., 2019; García-Durán et al., 2018). Recently, graph neural networks (GNN) have been incorporated to enhance representation learning by aggregating neighborhood information (Schlichtkrull et al., 2017; Dettmers et al., 2018; Nguyen et al., 2018; Vashishth et al., 2020; Li et al., 2021). However, most approaches significantly degrade when KG are largely incomplete and contain certain errors (Pujara et al., 2017), as they ignore the false negative/positive issues. Recent attempts on uncertain KG (Chen et al., 2019; Kertkeidkachorn et al., 2019; Chen et al., 2021) measure the uncertainty score for facts, which can detect false negative/positive samples. However, they explicitly require the ground truth uncertainty

scores for model training, which are usually unavailable in practice. Various negative sampling strategies have been explored to sample informative negative triples to facilitate model training (Cai and Wang, 2018; Zhang et al., 2019; Ahrabian et al., 2020). However, they cannot detect false negative/positive facts. We aim to mitigate the false negative/positive issues and enable the automatic detection of false negative/positive facts during model training.

Positive-Unlabeled Learning. Positive-Unlabeled (PU) learning is a learning paradigm for training a model that only has access to positive and unlabeled data, where unlabeled data includes both positive and negative samples (Plessis et al., 2015; Bekker and Davis, 2018). PU learning roughly includes (1) two-step solutions (He et al., 2018; Jain et al., 2016); (2) methods that consider the unlabeled samples as negative samples with label noise (Shi et al., 2018); (3) unbiased risk estimation methods (Plessis et al., 2015; Tang et al., 2022). Recent work further studies the setting that there exists label noise in the observed positive samples (Jain et al., 2016). We formulate the KGR task on noisy and incomplete KG as a noisy Positive-Unlabeled learning problem and propose a variational framework for it, which relates to two-step solutions and unbiased risk estimation methods.

6 Conclusion

We studied speculative KG reasoning based on sparse and unreliable observations, which contains both *false negative issue* and *false positive issue*. We formulated the task as a noisy Positive-Unlabeled learning problem and proposed a variational framework nPUGraph to jointly update model parameters and estimate the posterior likelihood of collected/uncollected facts being true or false, where the underlying correctness is viewed as latent variables. During the training process, a label posterior-aware encoder and a self-training strategy were proposed to further address the false positive/negative issues. We found label posterior estimation plays an important role in moving toward speculative KG reasoning in reality, and the estimation can be fulfilled by optimizing an alternative objective without additional cost. Extensive experiments verified the effectiveness of nPUGraph on both benchmark KGs and Twitter interaction data with various degrees of data perturbations.

627 7 Limitations

628 There are certain limitations that can be concerned
629 for further improvements. First, the posterior inference
630 relies on the prior estimation of positive class prior α and true positive ratio β . Our experiments
631 show that a data-driven estimation based on end-to-end model training produces worse results
632 than a hyperparameter grid search. An automatic prior estimation is desirable for real-world applica-
633 tions. Moreover, in nPUGraph, we approximate the probability of negative/positive facts being col-
634 lected/uncollected via neural networks, which lacks
635 a degree of interpretability. In the future, we plan
636 to utilize a more explainable random process de-
637 pending on entity/relation features to model the
638 collection probability distribution.

643 8 Ethical Impact

644 nPUGraph neither introduces any social/ethical
645 bias to the model nor amplifies any bias in data.
646 Benchmark KG are publicly available. For Twitter
647 interaction data, we mask all identity and privacy
648 information for users, where only information re-
649 lated to user interactions with tweets and hashtags
650 is presented. Our model is built upon public li-
651 braries in PyTorch. We do not foresee any direct
652 social consequences or ethical issues.

653 References

- 654 Kian Ahrabian, Aarash Feizi, Yasmin Salehi, William L.
655 Hamilton, and Avishek Joey Bose. 2020. **Structure**
656 **aware negative sampling in knowledge graphs**. In *Proceedings of the 2020 Conference on Empirical*
657 *Methods in Natural Language Processing (EMNLP)*, pages 6093–6101, Online. Association for Compu-
658 *tational Linguistics*.
- 661 Jessa Bekker and Jesse Davis. 2018. **Learning from**
662 **positive and unlabeled data: A survey**. *CoRR*,
663 abs/1811.04820.
- 664 Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim
665 Sturge, and Jamie Taylor. 2008. Freebase: a col-
666 laboratively created graph database for structuring
667 human knowledge. In *SIGMOD ’08: Proceedings of*
668 *the 2008 ACM SIGMOD international conference on*
669 *Management of data*, pages 1247–1250.
- 670 Antoine Bordes, Nicolas Usunier, Alberto Garcia-
671 Duran, Jason Weston, and Oksana Yakhnenko.
672 2013. Translating embeddings for modeling multi-
673 relational data. In *Advances in Neural Information*
674 *Processing Systems*.
- 675 Liwei Cai and William Yang Wang. 2018. **KBGAN:**
676 **Adversarial learning for knowledge graph embed-**
677 **dings**. In *Proceedings of the 2018 Conference of*
678 *the North American Chapter of the Association for*
679 *Computational Linguistics: Human Language Tech-*
680 *nologies, Volume 1 (Long Papers)*, pages 1470–1480,
681 New Orleans, Louisiana. Association for Compu-
682 *tational Linguistics*.
- 683 Xuelu Chen, Michael Boratko, Muhan Chen,
684 Shib Sankar Dasgupta, Xiang Lorraine Li, and
685 Andrew McCallum. 2021. Probabilistic box em-
686 beddings for uncertain knowledge graph reasoning.
687 In *Proceedings of the 2021 Conference of the*
688 *North American Chapter of the Association for*
689 *Computational Linguistics: Human Language*
690 *Technologies*.
- 691 Xuelu Chen, Muhan Chen, Weijia Shi, Yizhou Sun,
692 and Carlo Zaniolo. 2019. Embedding uncertain
693 knowledge graphs. In *Proceedings of the Thirty-*
694 *Third AAAI Conference on Artificial Intelligence and*
695 *Thirty-First Innovative Applications of Artificial In-*
696 *telligence Conference and Ninth AAAI Symposium*
697 *on Educational Advances in Artificial Intelligence,*
698 *AAAI’19/IAAI’19/EAAI’19*.
- 699 Tim Dettmers, Pasquale Minervini, Pontus Stenetorp,
700 and Sebastian Riedel. 2018. Convolutional 2d knowl-
701 edge graph embeddings. In *Proceedings of the Thirty-*
702 *Second AAAI Conference on Artificial Intelligence and*
703 *Thirtieth Innovative Applications of Artificial In-*
704 *telligence Conference and Eighth AAAI Symposium*
705 *on Educational Advances in Artificial Intelligence,*
706 *AAAI’18/IAAI’18/EAAI’18*.
- 707 Christiane Fellbaum. 1998. *WordNet: An Electronic*
708 *Lexical Database*. Bradford Books.
- 709 Alberto García-Durán, Sebastijan Dumančić, and Math-
710 ias Niepert. 2018. Learning sequence encoders for
711 temporal knowledge graph completion. In *Proceed-
712 ings of the 2018 Conference on Empirical Methods*
713 *in Natural Language Processing*.
- 714 Fengxiang He, Tongliang Liu, Geoffrey I. Webb,
715 and Dacheng Tao. 2018. **Instance-dependent PU**
716 **learning by bayesian optimal relabeling**. *CoRR*,
717 abs/1808.02180.
- 718 Shantanu Jain, Martha White, and Predrag Radivojac.
719 2016. Estimating the class prior and posterior from
720 noisy positives and unlabeled data. In *Proceedings of*
721 *the 30th International Conference on Neural Infor-*
722 *mation Processing Systems, NIPS’16*, page 2693–2701,
723 Red Hook, NY, USA. Curran Associates Inc.
- 724 Natthawut Kertkeidkachorn, Xin Liu, and Ryutaro
725 Ichise. 2019. Gtranse: Generalizing translation-
726 based model on uncertain knowledge graph embed-
727 ding. In *Advances in Artificial Intelligence - Selected*
728 *Papers from the Annual Conference of Japanese So-*
729 *cieties of Artificial Intelligence (JSAI 2019), Niigata,*
730 *Japan, 4-7 June 2019*, volume 1128 of *Advances in*
731 *Intelligent Systems and Computing*, pages 170–178.

732	Zixuan Li, Xiaolong Jin, Wei Li, Saiping Guan, Jiafeng	786
733	Guo, Huawei Shen, Yuanzhuo Wang, and Xueqi	787
734	Cheng. 2021. Temporal knowledge graph reason-	788
735	ing based on evolutional representation learning.	789
736	In <i>SIGIR</i> .	790
737	Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and	791
738	Xuan Zhu. 2015. Learning entity and relation embed-	792
739	dings for knowledge graph completion. In <i>AAAI'15</i> .	793
740	Dai Quoc Nguyen, Tu Dinh Nguyen, Dat Quoc Nguyen,	794
741	and Dinh Phung. 2018. A novel embedding model	795
742	for knowledge base completion based on convolu-	796
743	tional neural network. In <i>Proceedings of the 16th</i>	797
744	<i>Annual Conference of the North American Chapter</i>	798
745	<i>of the Association for Computational Linguistics: Hu-</i>	799
746	<i>man Language Technologies (NAACL-HLT)</i> , pages	800
747	327–333.	801
748	Marthinus Du Plessis, Gang Niu, and Masashi	802
749	Sugiyama. 2015. Convex formulation for learning	803
750	from positive and unlabeled data. In <i>Proceedings</i>	804
751	<i>of the 32nd International Conference on Machine</i>	805
752	<i>Learning</i> , volume 37 of <i>Proceedings of Machine</i>	806
753	<i>Learning Research</i> , pages 1386–1394, Lille, France.	807
754	Jay Pujara, Eriq Augustine, and Lise Getoor. 2017. <i>Spar-</i>	808
755	<i>sity and noise: Where knowledge graph embeddings</i>	809
756	<i>fall short</i> . In <i>Proceedings of the 2017 Conference on</i>	809
757	<i>Empirical Methods in Natural Language Processing</i> ,	809
758	pages 1751–1756, Copenhagen, Denmark. Associa-	809
759	tion for Computational Linguistics.	809
760	Jianwei Qian, Xiang-Yang Li, Chunhong Zhang, Linlin	810
761	Chen, Taeho Jung, and Junze Han. 2019. Social net-	811
762	work de-anonymization and privacy inference with	812
763	knowledge graph model. <i>IEEE Transactions on De-</i>	813
764	<i>pendable and Secure Computing</i> , pages 679–692.	814
765	Apoorv Saxena, Soumen Chakrabarti, and Partha Taluk-	815
766	dar. 2021. Question answering over temporal knowl-	816
767	edge graphs. In <i>Proceedings of the 59th Annual</i>	817
768	<i>Meeting of the Association for Computational Lin-</i>	818
769	<i>guistics and the 11th International Joint Conference</i>	819
770	<i>on Natural Language Processing (Volume 1: Long</i>	819
771	<i>Papers</i>).	819
772	Michael Sejr Schlichtkrull, Thomas N. Kipf, Peter	820
773	Bloem, Rianne van den Berg, Ivan Titov, and Max	821
774	Welling. 2017. Modeling relational data with graph	822
775	convolutional networks. In <i>ESWC</i> , pages 593–607.	823
776	Hong Shi, Shaojun Pan, Jian Yang, and Chen Gong.	824
777	2018. Positive and unlabeled learning via loss de-	825
778	composition and centroid estimation. In <i>Proceedings</i>	826
779	<i>of the 27th International Joint Conference on Artifi-</i>	827
780	<i>cial Intelligence</i> , IJCAI'18, page 2689–2695. AAAI	827
781	Press.	828
782	Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian	828
783	Tang. 2019. Rotate: Knowledge graph embedding by	828
784	relational rotation in complex space. In <i>International</i>	828
785	<i>Conference on Learning Representations</i> .	828

A Appendix

A.1 Theorem Proof

Theorem 2. *The log-likelihood of the complete data $\log p(\mathbf{S})$ is lower bounded as follows:*

$$\begin{aligned} \log p(\mathbf{S}) &\geq \mathbb{E}_{q(\mathbf{Y})} [\log p(\mathbf{S}|\mathbf{Y})] - \text{KL}(q(\mathbf{Y})\|p(\mathbf{Y})) \\ &= \mathbb{E}_{s^l \in \mathcal{S}^L} [w^l \log[\phi_1^l(s^l)] + (1-w^l) \log[\phi_0^l(s^l)]] \\ &+ \mathbb{E}_{s^u \in \mathcal{S}^U} [(w^u \log[\phi_1^u(s^u)] + (1-w^u) \log[\phi_0^u(s^u)])] \\ &- \text{KL}(\mathbf{W}^U\|\tilde{\mathbf{W}}^U) - \text{KL}(\mathbf{W}^L\|\tilde{\mathbf{W}}^L) - \frac{\|\mathbf{W}^L\|_1}{|\mathcal{S}^L|} - \frac{\|\mathbf{W}^U\|_1}{|\mathcal{S}^U|}, \end{aligned} \quad (13)$$

where \mathbf{S} denotes all labeled/unlabeled triples, \mathbf{Y} is the corresponding latent variable indicating the positive/negative labels for triples, $\mathbf{W}^U = \{w_i^u\}$ denotes the point-wise probability for the uncollected triples being positive, $\mathbf{W}^L = \{w_i^l\}$ denotes the probability for the collected triples being positive.

Proof. Let $\log p(\mathbf{S})$ denote the log-likelihood of all potential triples being collected in the KG or not, \mathbf{Y} denote the corresponding latent variable indicating the positive/negative labels. We aim to infer the label posterior $p(\mathbf{Y}|\mathbf{S})$, which can be approximated by $q(\mathbf{Y}|\mathbf{S})$. We therefore are interested at the difference between the two, measured by the Kullback–Leibler (KL) divergence as follows:

$$\text{KL}(q(\mathbf{Y}|\mathbf{S})\|p(\mathbf{Y}|\mathbf{S})) = -\mathbb{E}_{q(\mathbf{Y}|\mathbf{S})} \log \left(\frac{p(\mathbf{S}|\mathbf{Y})p(\mathbf{Y})}{q(\mathbf{Y}|\mathbf{S})} \right) + \log p(\mathbf{S}), \quad (14)$$

as KL divergence is positive, we derive the lower bound of the log-likelihood as follows:

$$\begin{aligned} \log p(\mathbf{S}) &\geq \mathbb{E}_{q(\mathbf{Y}|\mathbf{S})} \log \left(\frac{p(\mathbf{S}|\mathbf{Y})p(\mathbf{Y})}{q(\mathbf{Y}|\mathbf{S})} \right) \\ &\geq \mathbb{E}_{q(\mathbf{Y}|\mathbf{S})} \log p(\mathbf{S}|\mathbf{Y}) - \text{KL}(q(\mathbf{Y}|\mathbf{S})\|p(\mathbf{Y})) - \mathbb{E}_{p(\mathbf{S})} q(\mathbf{Y}|\mathbf{S}), \end{aligned} \quad (15)$$

which consists of three terms: triple collection probability measure, KL term and regularization of label posterior (positive). We discuss each term in detail.

Recall that the distribution of labeled triples can be represented as follows:

$$s^l \sim \beta\phi_1^l(s^l) + (1-\beta)\phi_0^l(s^l), \quad (16)$$

where ϕ_y^l denotes the probability of being collected over triple space \mathcal{S} for the positive class ($y = 1$) and negative class ($y = 0$), and $\beta \in [0, 1]$ denotes the proportion of true positive samples in labeled

data. Similarly, considering the existence of unlabeled positive triples, the distribution of unlabeled samples can be represented as follows:

$$s^u \sim \alpha\phi_1^u(s^u) + (1-\alpha)\phi_0^u(s^u), \quad (17)$$

where $\phi_y^u = 1 - \phi_y^l$ denotes the probability of being uncollected, $\alpha \in [0, 1]$ is the class prior or the proportion of positive samples in unlabeled data. Based on that, the first term can be detailed as follows:

$$\begin{aligned} \mathbb{E}_{q(\mathbf{Y}|\mathbf{S})} \log p(\mathbf{S}|\mathbf{Y}) &= \mathbb{E}_{s^l \in \mathcal{S}^L} \mathbb{E}_{y \in \{0,1\}} q(y|s^l) \log p(s^l|y) \\ &+ \mathbb{E}_{s^u \in \mathcal{S}^U} \mathbb{E}_{y \in \{0,1\}} q(y|s^u) \log p(s^u|y) \\ &= \mathbb{E}_{s^l \in \mathcal{S}^L} [q(y=1|s^l) \log[p(s^l|y=1)] + q(y=0|s^l) \log[p(s^l|y=0)]] \\ &+ \mathbb{E}_{s^u \in \mathcal{S}^U} [q(y=1|s^u) \log[p(s^u|y=1)] + q(y=0|s^u) \log[p(s^u|y=0)]] \\ &= \mathbb{E}_{s^l \in \mathcal{S}^L} [w^l \log[\phi_1^l(s^l)] + (1-w^l) \log[\phi_0^l(s^l)]] \\ &+ \mathbb{E}_{s^u \in \mathcal{S}^U} [(w^u \log[\phi_1^u(s^u)] + (1-w^u) \log[\phi_0^u(s^u)])], \end{aligned} \quad (18)$$

where $\mathbf{W}^U = \{w^u\}$ denotes the point-wise probability for the uncollected triples being positive $q(y=1|s^u)$, $\mathbf{W}^L = \{w^l\}$ denotes the probability for the collected triples being positive $q(y=1|s^l)$.

We view \mathbf{W}^U and \mathbf{W}^L as free parameters and regularize them by $\tilde{\mathbf{W}}^U$ and $\tilde{\mathbf{W}}^L$, which are estimated posterior probability as follows:

$$\tilde{w}_i^l = \frac{\beta\phi_1^l(s_i^l)}{\beta\phi_1^l(s_i^l) + (1-\beta)\phi_0^l(s_i^l)}, \quad (19)$$

$$\tilde{w}_{ik}^u = \frac{\alpha\phi_1^u(s_{ik}^u)}{\alpha\phi_1^u(s_{ik}^u) + (1-\alpha)\phi_0^u(s_{ik}^u)}. \quad (20)$$

Therefore, the KL term in Eq. (15) becomes $\text{KL}(\mathbf{W}^U\|\tilde{\mathbf{W}}^U) + \text{KL}(\mathbf{W}^L\|\tilde{\mathbf{W}}^L)$.

Last but not least, the third term $\mathbb{E}_{p(\mathbf{S})} q(\mathbf{Y}|\mathbf{S}) = \|\mathbf{W}^U\|_1/|\mathcal{S}^U| + \|\mathbf{W}^L\|_1/|\mathcal{S}^L|$ regulates the total number of potential triples (including both collected ones and uncollected ones), because of the sparsity nature of graphs. Finally, we derive the lower bound of the log-likelihood, as shown in Eq (13). \square

A.2 Datasets

A.2.1 Dataset Information

We evaluate our proposed model based on three widely used knowledge graphs and one Twitter interaction graph:

Table 3: The statistics of the datasets.

#Ptb Rate	Dataset	#Entities	#Relation	#Train	#Valid	#Test
0.1	FB15K	14,951	1,345	340,968	146,129	59,071
	FB15K-237	14,541	237	184,803	79,201	20,466
	WN18	40,943	18	92,428	39,612	5,000
	Twitter	17,839	2	282,233	120,956	110,456
0.3	FB15K	14,951	1,345	276,940	118,688	59,071
	FB15K-237	14,541	237	149,229	63,954	20,466
	WN18	40,943	18	72,462	31,055	5,000
	Twitter	17,839	2	232,748	99,749	110,456
0.5	FB15K	14,951	1,345	213,380	91,448	59,071
	FB15K-237	14,541	237	113,772	48,759	20,466
	WN18	40,943	18	52,707	22,588	5,000
	Twitter	17,839	2	183,263	78,540	110,456
0.7	FB15K	14,951	1,345	150,485	64,493	59,071
	FB15K-237	14,541	237	78,531	33,656	20,466
	WN18	40,943	18	34,984	14,993	5,000
	Twitter	17,839	2	133,778	57,333	110,456

- **FB15K** (Bordes et al., 2013) is a subset of Freebase (Bollacker et al., 2008), a large database containing general knowledge facts with a variety of relation types;
- **FB15K-237** (Toutanova et al., 2015) is a reduced version of FB15K, where inverse relations are removed;
- **WN18** (Bordes et al., 2013) is a subset of WorldNet (Fellbaum, 1998), a massive lexical English database that captures semantic relations between words;
- **Twitter** is an interaction graph relevant with *Russo-Ukrainian War*. Data is collected in the Twitter platform from May 1, 2022, to December 25, 2022, which records user-tweet interactions and user-hashtag interactions. Thus, the graph is formed by two relations (user-tweet and user-hashtag) and multiple entities, which can be categorized into three types (user, tweet, and hashtag). When constructing the graph, thresholds will be selected to remove inactive users, tweets, and hashtags according to their occurrence frequency. We set the thresholds for the user, tweet, and hashtag as 30, 30, and 10, respectively, i.e., if a tweet has fewer than 30 interactions with users, it will be regarded as inactive and removed from the graph. Besides, the extremely frequent users and tweets are deleted as they may be generated by bots.

For all datasets, we first merge the training set and validation set as a whole. Then we simulate noisy and incomplete graphs for the training process by adding various proportions of false negative/positive cases in the merged set. After that, we partition the simulated graphs into new train/valid sets with a ratio of 7 : 3 and the test set remains the

same. Table 3 provides an overview of the statistics of the simulated datasets corresponding to various perturbation rates and based graphs.

A.2.2 Dataset Perturbation

Data perturbation aims to simulate noisy and incomplete graphs from clean benchmark knowledge graphs. It consists of two aspects: First, to simulate the *false negative issue*, it randomly removes some existing links in a graph, considering the removed links as missing but potentially plausible facts. Second, to simulate the *false positive issue*, it randomly adds spurious links to the graph as unreliable or outdated facts.

We define perturbation rate, i.e., ptb_rate , as a proportion of modified edges in a graph to control the amount of removing positive links and adding negative links. For example, if a graph has 100 links and the perturbation rate is 0.5, then we will randomly convert the positivity or negativity of 50 links. Among these 50 modified links, 10% of them will be added and the rest of them will be removed, i.e., we will randomly add 5 negative links and remove 45 positive links to generate a perturbed graph. The perturbed graph can be regarded as noisy and incomplete, leading to significant performance degradation. In our experiments, we set ptb_rate in a range of {0.1, 0.3, 0.5, 0.7}. The detailed perturbation process is summarized in Figure 6.

A.3 Baselines

We describe the baseline models utilized in the experiments in detail:

- **TransE**² (Bordes et al., 2013) is a translation-based embedding model, where both entities and relations are represented as vectors in the latent space. The relation is utilized as a translation operation between the subject and the object entity;
- **TransR** (Lin et al., 2015) advances TransE by optimizing modeling of n-n relations, where each entity embedding can be projected to hyperplanes defined by relations;
- **DistMult**³ (Yang et al., 2014) is a general framework with the bilinear objective for multi-

²The experiments of TransE and TransR are implemented with <https://github.com/thunlp/OpenKE>.

³The experiments of DistMult, ComplEx, and RotatE are implemented with <https://github.com/DeepGraphLearning/KnowledgeGraphEmbedding>.

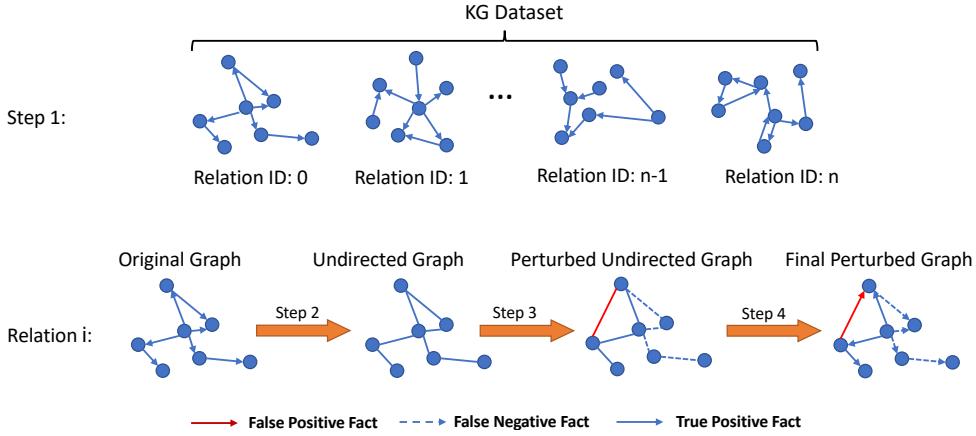


Figure 6: Summary of graph perturbation to construct noisy and incomplete KGs.

977 relational learning that unifies most multi-
978 relational embedding models;

- 979 • **ComplEx** (Trouillon et al., 2017) introduces
980 complex embeddings, which can effectively
981 capture asymmetric relations while retaining
982 the efficiency benefits of the dot product;
- 983 • **RotatE** (Sun et al., 2019) extends ComplEx
984 by representing entities as complex vectors and
985 relations as rotation operations in a complex
986 vector space;
- 987 • **R-GCN⁴** (Schlichtkrull et al., 2017) uses
988 relation-specific weight matrices that are de-
989 fined as linear combinations of a set of basis
990 matrices;
- 991 • **CompGCN⁵** (Vashishth et al., 2020) is a frame-
992 work for incorporating multi-relational informa-
993 tion in graph convolutional networks to jointly
994 embeds both nodes and relations in a graph;
- 995 • **UKGE⁶** (Chen et al., 2019) learns embeddings
996 according to the confidence scores of uncertain
997 relation facts to preserve both structural and un-
998 certainty information of facts in the embedding
999 space;
- 1000 • **NSCaching⁷** (Zhang et al., 2019) is an inex-
1001 pensive negative sampling approach by using cache
1002 to keep track of high-quality negative triplets,
1003 which have high scores and rare;

- 1004 • **SANS⁸** (Ahrabian et al., 2020) utilizes the rich
1005 graph structure by selecting negative samples
1006 from a node’s k-hop neighborhood for negative
1007 sampling without additional parameters and dif-
1008 ficult adversarial optimization;
- 1009 • **PUDA⁹** (Tang et al., 2022) is a KGC method
1010 to circumvent the impact of the false negative
1011 issue by tailoring positive unlabeled risk esti-
1012 mator and address the data sparsity issue by
1013 unifying adversarial training and PU learning
1014 under the positive-unlabeled minimax game.

A.4 Reproducibility

A.4.1 Baseline Setup

All baseline models and nPUGraph are trained on the perturbed training set and validated on the perturbed valid set. We utilize *MRR* on the valid set to determine the best models and evaluate them on the clean test set. For uncertain knowledge graph embedding methods UKGE (Chen et al., 2019), since the required uncertainty scores are unavailable, we set the scores for triples in training set as 1 and 0 otherwise. The predicted uncertainty scores produced by UKGE are utilized to rank the potential triples for ranking evaluation. We train all baseline models and nPUGraph on the same GPUs (GeForce RTX 3090) and CPUs (AMD Ryzen Threadripper 3970X 32-Core Processor).

A.4.2 nPUGraph Setup

For model training, we utilize Adam optimizer and set the maximum number of epochs as 200. Within the first 50 epochs, we disable self-training and focus on learning suboptimal model parameters on

⁴<https://github.com/JinheonBaek/RGCN>

⁵<https://github.com/mallabiisc/CompGCN>

⁶<https://github.com/stas10217/UKGE>

⁷<https://github.com/AutoML-Research/NSCaching>

⁸<https://github.com/kahrabian/SANS>

⁹<https://github.com/lilv98/PUDA>

Table 4: Overall performance on noisy and incomplete Twitter data, with $ptb_rate = 0.3$. Average results on 5 independent runs are reported. * indicates the statistically significant results over baselines, with p -value < 0.01 . The best results are in boldface, and the strongest baseline performance is underlined.

Dataset		Twitter			
Metrics		MRR	HIT@100	HIT@50	HIT@30
Random		0.001	0.006	0.003	0.002
<i>Knowledge graph embedding methods</i>					
TransE		0.010	0.091	0.058	0.041
TransR		0.009	0.078	0.048	0.033
DistMult		0.021	0.091	0.065	0.052
ComplEx		0.022	0.089	0.064	0.051
RotatE		0.022	0.1115	0.077	0.059
<i>Graph neural network methods on KG</i>					
RGCN		0.005	0.054	0.029	0.019
CompGCN		0.014	0.089	0.059	0.044
<i>Uncertain knowledge graph embedding method</i>					
UKGE		0.011	0.072	0.053	0.033
<i>Negative sampling methods</i>					
NSCaching		0.012	0.095	0.060	0.043
SANS		0.019	0.104	0.070	0.054
<i>Positive-Unlabeled learning method</i>					
PUDA		0.013	0.082	0.057	0.044
nPUGraph		0.030*	0.127*	0.096*	0.074*
Gains (%)		38.2	13.9	25.3	26.5

noisy and incomplete data. After that, we start the self-training strategy, where the latest label posterior estimation is utilized to sample neighbors for the encoder and select informative unlabeled samples for model training. We set batch size as 256, the dimensions of all embeddings as 128, and the dropout rate as 0.5. For the sake of efficiency, we employ 1 neighborhood aggregation layer in the encoder.

For the setting of hyperparameter, we mainly tune positive class prior α in the range of $\{1e - 1, 5e - 2, 1e - 2, 5e - 3, 1e - 3, 5e - 4, 1e - 4, 5e - 5\}$; true positive ratio β in the range of $\{0.3, 0.2, 0.1, 0.005, 0.001\}$; learning rate in the range of $\{0.02, 0.01, 0.005, 0.001, 0.0005\}$; the number of sampled unlabeled triples for each labeled one in the range of $\{50, 40, 30, 20, 10\}$. We will publicly release our code and data upon acceptance.

A.5 Experiments

A.5.1 Experimental Results on Twitter Data

We discuss the model performance on noisy and incomplete Twitter data with $ptb_rate = 0.3$ in this section, which is shown in Table 4. According to the result of Random, we can infer that all relations on Twitter are $n - n$, where relations can be $1 - n$, $n - 1$, and $n - n$ for benchmark KG. Therefore, link prediction is more challenging for Twitter data and

we adopt Hits at 30, 50, 100 as evaluation metrics, instead.

For Twitter data, nPUGraph achieves impressive performance compared with the baseline models, with 25.98% relative improvement on average. The results of Twitter data support the robustness of nPUGraph, which can mitigate false negative/positive issues not only in benchmark KG but also in the real-world social graph.

A.5.2 Experimental Results under Different Perturbation Rates

The experimental results under perturbation rates 0.1, 0.5, and 0.7 are shown in Table 5, Table 6, and Table 7, respectively. nPUGraph outperforms all baseline models for various perturbation rates, demonstrating that nPUGraph can mitigate false negative/positive issues on knowledge graphs with different degrees of noise and incompleteness. Notably, comparing these three tables, the relative improvements are more significant under higher ptb_rate in most cases, showing stronger robustness for nPUGraph on graphs with more false negative/positive facts.

1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086

Table 5: Overall performance on noisy and incomplete graphs, with $ptb_rate = 0.1$. Average results on 5 independent runs are reported. * indicates the statistically significant results over baselines, with $p\text{-value} < 0.01$. The best results are in boldface, and the strongest baseline performance is underlined.

Dataset	FB15K				FB15K-237				WN18			
	Metrics	MRR	H@10	H@3	H@1	MRR	H@10	H@3	H@1	MRR	H@10	H@3
<i>Knowledge graph embedding methods</i>												
TransE	0.419	0.666	0.507	0.280	0.245	0.441	0.284	0.144	0.318	0.646	0.579	0.044
TransR	0.398	0.658	0.494	0.250	0.246	0.434	0.280	0.153	0.319	0.646	0.575	0.051
DistMult	0.516	0.719	0.578	0.407	0.271	0.446	0.297	0.185	0.524	0.686	0.610	0.418
ComplEx	0.503	0.711	0.570	0.390	0.276	0.459	0.305	0.186	0.612	0.702	0.643	0.560
Rotate	0.544	0.732	0.608	0.441	0.292	0.480	0.324	0.199	0.613	0.691	0.642	0.567
<i>Graph neural network methods on KG</i>												
RGCN	0.196	0.372	0.209	0.110	0.169	0.317	0.177	0.097	0.483	0.625	0.554	0.396
CompGCN	0.460	0.677	0.525	0.343	0.293	0.475	0.324	0.203	0.608	0.686	0.636	0.564
<i>Uncertain knowledge graph embedding method</i>												
UKGE	0.338	0.425	0.321	0.233	0.231	0.411	0.204	0.110	0.381	0.541	0.407	0.331
<i>Negative sampling method</i>												
NSCaching	0.495	0.689	0.557	0.390	0.153	0.305	0.167	0.080	0.434	0.542	0.470	0.374
SANS	0.422	0.649	0.493	0.298	0.271	0.453	0.301	0.182	0.619	0.702	0.644	0.574
<i>Positive-Unlabeled learning method</i>												
PUDA	0.493	0.713	0.559	0.377	0.271	0.443	0.298	0.185	0.520	0.667	0.608	0.419
nPUGraph	0.561*	0.791*	0.621*	0.449*	0.328*	0.535*	0.343*	0.221*	0.630*	0.754*	0.671*	0.599*
Gains (%)	3.0	8.1	2.1	1.9	12.0	11.4	5.9	8.7	1.8	7.4	4.1	4.4

Table 6: Experimental results under $ptb_rate = 0.5$.

Dataset	FB15K				FB15K-237				WN18			
	Metrics	MRR	H@10	H@3	H@1	MRR	H@10	H@3	H@1	MRR	H@10	H@3
<i>Knowledge graph embedding methods</i>												
TransE	0.279	0.540	0.363	0.136	0.151	0.342	0.190	0.053	0.158	0.337	0.285	0.018
TransR	0.256	0.500	0.323	0.127	0.141	0.294	0.160	0.064	0.150	0.327	0.267	0.020
DistMult	0.328	0.536	0.372	0.224	0.210	0.366	0.228	0.133	0.279	0.370	0.319	0.224
ComplEx	0.322	0.525	0.365	0.220	0.201	0.358	0.220	0.123	0.314	0.373	0.335	0.280
Rotate	0.350	0.547	0.398	0.249	0.227	0.387	0.246	0.149	0.307	0.377	0.333	0.266
<i>Graph neural network methods on KG</i>												
RGCN	0.134	0.271	0.142	0.065	0.116	0.234	0.117	0.058	0.253	0.323	0.287	0.209
CompGCN	0.378	0.600	0.429	0.266	0.223	0.377	0.240	0.149	0.345	0.315	0.336	0.279
<i>Uncertain knowledge graph embedding method</i>												
UKGE	0.257	0.299	0.213	0.088	0.143	0.284	0.172	0.053	0.228	0.297	0.210	0.115
<i>Negative sampling method</i>												
NSCaching	0.272	0.454	0.310	0.179	0.176	0.297	0.190	0.115	0.123	0.182	0.133	0.093
SANS	0.335	0.545	0.387	0.224	0.225	0.382	0.244	0.147	0.313	0.379	0.337	0.275
<i>Positive-Unlabeled learning method</i>												
PUDA	0.329	0.525	0.369	0.229	0.202	0.347	0.217	0.131	0.231	0.329	0.264	0.178
nPUGraph	0.417*	0.663*	0.470*	0.291*	0.258*	0.433*	0.285*	0.171*	0.373*	0.443*	0.379*	0.327*
Gains (%)	10.4	10.5	9.6	9.6	13.9	12.0	16.1	14.8	8.2	16.9	12.6	16.9

Table 7: Experimental results under $ptb_rate = 0.7$.

Dataset	FB15K				FB15K-237				WN18			
	Metrics	MRR	H@10	H@3	H@1	MRR	H@10	H@3	H@1	MRR	H@10	H@3
<i>Knowledge graph embedding methods</i>												
TransE	0.219	0.457	0.294	0.087	0.104	0.273	0.128	0.020	0.114	0.229	0.199	0.020
TransR	0.195	0.396	0.248	0.087	0.096	0.217	0.108	0.036	0.096	0.207	0.167	0.016
DistMult	0.250	0.425	0.282	0.163	0.173	0.308	0.186	0.105	0.181	0.240	0.211	0.143
ComplEx	0.241	0.408	0.271	0.157	0.158	0.290	0.170	0.092	0.202	0.243	0.219	0.178
Rotate	0.271	0.439	0.309	0.185	0.198	0.337	0.212	0.129	0.192	0.250	0.212	0.159
<i>Graph neural network methods on KG</i>												
RGCN	0.129	0.229	0.134	0.075	0.086	0.178	0.086	0.040	0.166	0.215	0.191	0.135
CompGCN	0.354	0.578	0.402	0.243	0.193	0.325	0.204	0.129	0.212	0.262	0.229	0.183
<i>Uncertain knowledge graph embedding method</i>												
UKGE	0.186	0.358	0.199	0.076	0.133	0.201	0.115	0.075	0.099	0.176	0.153	0.116
<i>Negative sampling method</i>												
NSCaching	0.174	0.309	0.194	0.105	0.157	0.276	0.169	0.099	0.057	0.085	0.062	0.041
SANS	0.292	0.478	0.332	0.196	0.201	0.342	0.216	0.131	0.202	0.254	0.222	0.171
<i>Positive-Unlabeled learning method</i>												
PUDA	0.254	0.421	0.283	0.170	0.165	0.291	0.176	0.104	0.109	0.171	0.127	0.077
nPUGraph	0.365*	0.600*	0.427*	0.277*	0.243*	0.390*	0.266*	0.155*	0.247*	0.303*	0.257*	0.209*
Gains (%)	3.0	3.8	6.3	13.9	20.9	14.1	23.0	18.5	16.8	15.5	12.3	14.4