

Manuscript Number:

Title: Probabilistic Belief Embedding for Large-scale Knowledge
Population

Article Type: Full Length Article

Keywords: knowledge population; belief embedding; entity inference;
relation prediction; triplet.

Corresponding Author: Mr. Miao Fan,

Corresponding Author's Institution:

First Author: Miao Fan

Order of Authors: Miao Fan; Qiang Zhou, Ph.D.; Andrew Abel, Ph.D.; Fang
Zheng, Ph.D.; Ralph Grishman, Ph.D.

Abstract: This paper contributes a novel embedding model which estimates the probability of each candidate belief $\langle h, r, t, m \rangle$ in a large-scale knowledge repository via simultaneously learning distributed representations for entities (h and t), relations (r), and the words in relation mentions (m). It facilitates knowledge population by means of simple vector operations to discover new beliefs. Given an imperfect belief, we can not only infer the missing entities, predict the unknown relations, but also tell the plausibility of the belief, just leveraging the learnt embeddings of remaining evidences. To demonstrate the scalability and the effectiveness of our model, we conduct experiments on several large-scale repositories which contain millions of beliefs from WordNet, Freebase and NELL, and compare it with other cutting-edge approaches via competing the performances assessed by the tasks of entity inference, relation prediction and triplet classification with their respective metrics. Extensive experimental results show that the proposed model outperforms the state-of-the-arts with significant improvements.

Suggested Reviewers:

Opposed Reviewers:

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Probabilistic Belief Embedding for Large-scale Knowledge Population

Miao Fan^a, Qiang Zhou^a, Andrew Abel^a, Fang Zheng^a, Ralph Grishman^c

^aCenter for Speech and Language Technologies, Department of Computer Science and Technology, Tsinghua University, 100084, Beijing, China.

^bComputing Science and Mathematics, School of Natural Sciences, Room 4B59, Cottrell Building, University of Stirling, Stirling FK9 4LA, U.K.

^cProteus Group, Computer Science Department, Courant Institute of Mathematical Sciences, New York University, 715 Broadway, 10003, NY, U.S.A.

Abstract

This paper contributes a novel embedding model which estimates the probability of each candidate belief $\langle h, r, t, m \rangle$ in a large-scale knowledge repository via simultaneously learning distributed representations for entities (h and t), relations (r), and the words in relation mentions (m). It facilitates knowledge population by means of simple vector operations to discover new beliefs. Given an imperfect belief, we can not only infer the missing entities, predict the unknown relations, but also tell the plausibility of the belief, just leveraging the learnt embeddings of remaining evidences. To demonstrate the scalability and the effectiveness of our model, we conduct experiments on several large-scale repositories which contain millions of beliefs from WordNet, Freebase and NELL, and compare it with other cutting-edge approaches via competing the performances assessed by the tasks of *entity inference*, *relation prediction* and *triplet classification* with their respective metrics. Extensive experimental results show that the proposed model outperforms the state-of-the-arts with significant improvements.

Keywords: Knowledge population, belief embedding, entity inference, relation prediction, triplet classification

URL: fanmiao.cs1t.thu@gmail.com (Miao Fan), zq-lxd@mail.tsinghua.edu.cn (Qiang Zhou), aka@cs.stir.ac.uk (Andrew Abel), fzhengetsinghua.edu.cn (Fang Zheng), grishman@cs.nyu.edu (Ralph Grishman)

1. Introduction

Information extraction [1, 2] has drawn much attention in recent years because of the explosive growth in the number of web pages. It is the study of extracting structured beliefs from unstructured online texts to populate knowledge bases. Thanks to the long-term efforts made by experts, crowdsourcing and even machine learning techniques, several web-scale knowledge repositories, such as Wordnet¹, Freebase² and NELL³, have been built. Among these knowledge repositories, WordNet [3] and Freebase [4, 5] follow the RDF format [6] that represents each belief as a triplet, i.e. $\langle head\ entity, relation, tail\ entity \rangle$, but NELL [7] goes a step further to extend each triplet with a *relation mention* which is a snatch of extracted free text to indicate the corresponding *relation*. Here we take a belief recorded in NELL as an example:

$\langle city : caroline, citylocatedinstate, stateorprovince : maryland, county\ and\ state\ of \rangle$, in which *county and state of* is the mention between the head entity *city : caroline*, and the tail entity *stateorprovince : maryland*, to indicate the relation *citylocatedinstate*. In some cases, NELL also provides the *confidence* of each belief automatically learnt by machines.

Although we have gathered colossal quantities of beliefs, state-of-the-art work [8] reports that our knowledge bases are far from complete. For instance, nearly 97% persons in Freebase have no records about their parents, whereas we human beings can still find the clue of their immediate family for most of the Freebase persons via searching on the web and looking up their Wiki. To populate the incomplete knowledge repositories assisted by computers, scientists either compete the performance of relation extraction between two named entities on manually annotated text datasets, such as ACE⁴ and MUC⁵, or look for effective approaches on improving the accuracy of link prediction within the

¹<http://wordnet.princeton.edu/>

²<https://www.freebase.com/>

³<http://rtw.ml.cmu.edu/rtw/>

⁴<http://www.itl.nist.gov/iad/mig/tests/ace/>

⁵<http://www.itl.nist.gov/iaui/894.02/relatedprojects/muc/>

1
2
3
4
5
6
7
8
9 knowledge graphs constructed by the repositories without using extra free texts.

10 Recently, studies on text-based knowledge population have benefited a lot
11 from a grateful paradigm called distantly supervised relation extraction (DSRE
12 [9]) which bridges the gap between structured knowledge bases and unstructured
13 free texts. It alleviates the labor of manual annotation by means of automat-
14 30 [9]) which bridges the gap between structured knowledge bases and unstructured
15 free texts. It alleviates the labor of manual annotation by means of automat-
16 ically aligning each triplet $\langle h, r, t \rangle$ from knowledge bases to the corresponding
17 relation mention m in free texts. However the latest research [10] points out
18 that DSRE still suffers from the problem of sparse and noisy features. Although
19 Fan et al. fix the issue to some extent via leveraging the low-dimensional matrix
20 factorization, the approach could not handle large-scale datasets as discussed in
21 their academic article [10].
22 35

23 Fortunately, the knowledge embedding techniques [11, 12] enlighten us to en-
24 code the high-dimensional sparse features into low-dimensional distributed rep-
25 resentations. A simple but effective model is TransE [13] which trains a vector
26 40 representation for each entity and relation in large-scale knowledge bases with-
27 out considering any text information. Even though Weston et al. [14], Wang et
28 al. [15] and Fan et al. [16] broaden this field by adding word embeddings, there
29 is still no comprehensive and elegant model that can integrate such large-scale
30 heterogeneous resources to satisfy multiple subtasks of knowledge population
31 including *entity inference*, *relation prediction* and *triplet classification*.
32
33
34
35
36
37
38 45

39 Therefore, we contribute a novel embedding model in this article, which es-
40 timates the probability of each candidate belief $\langle h, r, t, m \rangle$ in large-scale reposi-
41 tories. It breaks through the limitation of heterogeneous data, and establishes
42 the connection between the structured knowledge graph and unstructured free
43 texts. The distributed representations for entities (h and t), relations (r), as
44 well as the words in relation mentions (m) are simultaneously learnt within the
45 uniform framework of the probabilistic belief embedding (PBE) we propose.
46 Then knowledge population can be facilitated by means of simple vector oper-
47 ations to discover new beliefs. Given an imperfect belief, we can not only infer
48 the missing entities, predict the unknown relations, but tell the plausibility of
49 the belief as well, just by means of the learnt vector representations of remain-
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9 ing evidences. To prove the effectiveness and the scalability of PBE, we set up
10 extensive experiments on multiple tasks, including *entity inference*, *relation pre-*
11 *60 diction* and *triplet classification*, for knowledge population, and evaluate both
12 our model and the cutting-edge approaches with appropriate metrics on sever-
13 al well-known large-scale repositories, such as WordNet, Freebase and NELL,
14 which contain millions of beliefs. Elaborate results of comparison demonstrate
15 that the proposed model outperforms the state-of-the-arts with significant im-
16 provements.
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

2. Related Work

We generally group the studies of knowledge population into three categories according to the diverse resources they use: text-based knowledge extraction, repository-based knowledge inference and hybrid-based knowledge population. As their individual names imply, the first research community extracts the re-
70 lations between two recognized entities from text corpora, the second takes advantage of the link patterns within a knowledge graph to infer new triplets, and the third party suggests leveraging both the structure and unstructured information from both the text corpora and the knowledge graph. This paper
75 contributes a novel embedding model for hybrid-based knowledge population, which stands on the boundary between the second and the third research communities, and we thus conduct experiments that mainly compare our approach with several state-of-the-arts mentioned in Section 2.2 and 2.3.

2.1. Text-based knowledge extraction

80 There exists a huge amount of unstructured electronic texts on the Web. To better understand these online data, we would like to create an intelligent system that can annotate all the data with the structure of our concerns. Normally, we concern more about the knowledge on relations between named entities. So far, off-the-shelf softwares have been available to help recognize entities in texts,
85 and what we need to further study is to identify the semantic relations between

1
2
3
4
5
6
7
8
9 a pair of the annotated entities. But before we learn to extract the relations
10 with supervised learning, we should annotate a portion of the data first, and
11 the paradigms of annotation have two branches as follows.
12
13

14 2.1.1. Corpus-based extraction

15
16 90 Traditional approaches compete the performance of relation extraction on
17 the public corpora, including ACE and MUC, which have been annotated by
18 experts already. They choose different features extracted from the texts, like
19 syntactic [17], kernel [18] or semantic parser features [19], and adopt discrimina-
20 tive classifiers, such as Perceptron and Support Vector Machine (SVM) to help
21 predict the relations. There is a comprehensive survey [2] which shows more
22 details about this branch of research.
23
24 95
25
26
27

28 2.1.2. Distantly supervised extraction

29
30 Mintz et al. [9] firstly adopt Freebase to *distantly supervise* Wikipedia to
31 automatically generate annotated corpora. The basic alignment assumption is
32 that if a pair of entities participate in a relation, *all sentences* that mention these
33 entities in Wikipedia are labeled by the relation name from Freebase. Then we
34 100 can extract a variety of textual features and learn a multi-class logistic regression
35 classifier. Inspired by multi-instance learning, Riedel et al. [20] relax the strong
36 assumption and replace *all sentences* with *at least one sentence*. Hoffmann et
37 al. [21] point out that many entity pairs have more than one relation. They
38 extend the multi-instance learning framework to the multi-label circumstance.
39 Surdeanu et al. [22] propose a novel approach to multi-instance multi-label
40 learning for relation extraction, which jointly models all the sentences in texts
41 and all labels in knowledge bases for a given entity pair. The latest research
42 105 [10] points out that the distant supervision paradigm still suffers from sparse
43 and noisy features. Whereas Fan et al. [10] fix the issue by means of the low-
44 dimensional matrix factorization, as discussed in their scholar, the approach
45 could not handle large-scale datasets as well.
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

2.2. Repository-based knowledge inference

This research community aims at self-inferring new beliefs based on knowledge repositories without extra texts. It has two categories, namely graph-based inference models and embedding-based inference models. The principal differences between them are:

- *Symbolic representation v.s. Distributed representation:* Graph-based models regard the entities and relations as atomic elements, and represent them in a symbolic framework. In contrast, embedding-based models explore distributed representations via learning a low-dimensional continuous vector representation for each entity and relation.
- *Relation-specific v.s. Open-relation:* Graph-based models aim to induce rules or paths for a specific relation first, and then infer corresponding new beliefs. On the other hand, embedding-based models encode all relations into the same embedding space and conduct inference without any restriction on some specific relation.

2.2.1. Graph-based Inference

Graph-based inference models generally learn the representation for specific relations from the knowledge graph.

N-FOIL [23] learns first order Horn clause rules to infer new beliefs from the known ones. So far, it has helped to learn approximately 600 such rules. However, its ability to perform inference over large-scale knowledge repositories is currently still very limited.

PRA [24, 25, 26] is a data-driven random walk model which follows the paths from the head entity to the tail entity on the local graph structure to generate non-linear feature combinations representing the labeled relation, and uses logistic regression to select the significant features which contribute to classifying other entity pairs belonging to the given relation.

1
2
3
4
5
6
7
8
9
2.2.2. *Embedding-based Inference*

10 Embedding-based inference models usually design various scoring functions
11 $f_r(h, t)$ to measure the plausibility of a triplet $\langle h, r, t \rangle$. The lower the dissimilar-
12 ity of the scoring function $f_r(h, t)$ is, the higher the compatibility of the triplet
13 will be.
14
15
16

17 *Unstructured* [13] is a naive model which exploits the occurrence information
18 of the head and the tail entities without considering the relation between them.
19 It defines a scoring function $\|\mathbf{h} - \mathbf{t}\|$, and this model obviously can not discrim-
20 inate a pair of entities involving different relations. Therefore, *Unstructured* is
21 commonly regarded as the baseline approach.
22
23
24

25 *Distance Model (SE)* [11] uses a pair of matrices (W_{rh}, W_{rt}) , to characterize
26 a relation r . The dissimilarity of a triplet is calculated by $\|W_{rh}\mathbf{h} - W_{rt}\mathbf{t}\|_1$. As
27 pointed out by Socher et al. [27], the separating matrices W_{rh} and W_{rt} weak-
28 en the capability of capturing correlations between entities and corresponding
29 relations, even though the model takes the relations into consideration.
30
31
32

33 *Single Layer Model*, proposed by Socher et al. [27] thus aims to alleviate
34 the shortcomings of the *Distance Model* by means of the nonlinearity of a single
35 layer neural network $g(W_{rh}\mathbf{h} + W_{rt}\mathbf{t} + \mathbf{b}_r)$, in which $g = \tanh$. The linear output
36 layer then gives the scoring function: $\mathbf{u}_r^T g(W_{rh}\mathbf{h} + W_{rt}\mathbf{t} + \mathbf{b}_r)$.
37
38

39 *Bilinear Model* [28, 29] is another model that tries to fix the issue of weak
40 interaction between the head and tail entities caused by *Distance Model* with a
41 relation-specific bilinear form: $f_r(h, t) = \mathbf{h}^T W_r \mathbf{t}$.
42
43

44 *Neural Tensor Network (NTN)* [27] designs a general scoring function: $f_r(h, t) =$
45 $\mathbf{u}_r^T g(\mathbf{h}^T W_r \mathbf{t} + W_{rh}\mathbf{h} + W_{rt}\mathbf{t} + \mathbf{b}_r)$, which combines the *Single Layer Model* and
46 the *Bilinear Model*. This model is more expressive as the second-order corre-
47 lations are also considered into the nonlinear transformation function, but the
48 computational complexity is rather high.
49
50

51 *TransE* [13] is a canonical model different from all the other prior arts, which
52 embeds relations into the same vector space of entities by regarding the relation
53 r as a translation from h to t , i.e. $\mathbf{h} + \mathbf{r} = \mathbf{t}$. It works well on the beliefs
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9 with ONE-TO-ONE mapping property but performs badly on multi-mapping
10 beliefs. Given a series of facts associated with a ONE-TO-MANY relation r ,
11 e.g. $\langle h, r, t_1 \rangle, \langle h, r, t_2 \rangle, \dots, \langle h, r, t_m \rangle$, *TransE* tends to represent the embeddings
12 of entities on the MANY-side extreme close to each other which are hardly
13 discriminated.
14
15

175 *TransM* [30] leverages the structure of the whole knowledge graph, and ad-
16 justs the learning rate which is specific to each relation based on the multiple
17 mapping property of the relation.
18
19

21 *TransH* [31] is the state of the art approach as far as we know. It improves
22
23 180 *TransE* by modeling a relation as a hyperplane, which makes it more flexible
24 with regard to modeling beliefs with multi-mapping properties.
25
26

27 *2.3. Hybrid-based knowledge population*

29 Due to the diverse feature spaces between unstructured texts and structured
30 beliefs, the challenge of connecting natural language and knowledge turns out
31
32 185 to project the features into the same space and to merge them together for
33 knowledge population. Fan et al. [16] have recently proposed that they can
34 jointly learn the embedding representations for both relations and mentions to
35 predict unknown relations between entities in NELL. However, the functionality
36 of their latest method limits to the relation prediction task, as the correlations
37 between entities and relations are ignored. Therefore, we look forward to a
38
39 190 comprehensive model that can simultaneously consider entities, relations and
40 even the relation mentions, and can integrate the heterogeneous resources to
41 support multiple subtasks of knowledge population, such as *entity inference*,
42 *relation prediction* and *triplet classification*.
43
44
45
46
47
48
49

50 195 **3. Theory**

51
52 The intuition of the subsequent theory is that: Not each belief we have learnt,
53 i.e. $\langle head\ entity, relation, tail\ entity, mention \rangle$ abbreviated as $\langle h, r, t, m \rangle$, is
54 perfect and complete enough [32]. We thus explore modeling the probability of
55
56
57
58
59
60
61
62
63
64
65

each belief, i.e. $Pr(h, r, t, m)$. It is assumed that $Pr(h, r, t, m)$ is collaboratively influenced by $Pr(h|r, t)$, $Pr(t|h, r)$ and $Pr(r|h, t, m)$, where $Pr(h|r, t)$ stands for the conditional probability of inferring the head entity h given the relation r and the tail entity t , $Pr(t|h, r)$ represents the conditional probability of inferring the tail entity t given the head entity h and the relation r , and $Pr(r|h, t, m)$ denotes the conditional probability of predicting the relation r between the head entity h and the tail entity t with the relation mention m extracted from free texts. Therefore, we define that the probability of a belief equals to the geometric mean of $Pr(h|r, t)Pr(r|h, t, m)Pr(t|h, r)$ as shown in the subsequent equation,

$$Pr(h, r, t, m) = \sqrt[3]{Pr(h|r, t)Pr(r|h, t, m)Pr(t|h, r)}. \quad (1)$$

Suppose that we have a certain repository Δ , such as WordNet, which contains thousands of beliefs validated by experts. The learning object is intuitively set to maximize \mathcal{L}_{max} , where

$$\mathcal{L}_{max} = \prod_{\langle h, r, t, m \rangle \in \Delta} Pr(h, r, t, m). \quad (2)$$

In most cases, we can automatically build much larger but imperfect knowledge bases as well via crowdsourcing (Freebase) and machine learning techniques (NELL). However, each belief of NELL has a confidence-weighted score c to indicate its plausibility to some extent. Therefore, we propose an alternative goal which aims at minimizing \mathcal{L}_{min} , in which

$$\mathcal{L}_{min} = \prod_{\langle h, r, t, m, c \rangle \in \Delta} \frac{1}{2} [Pr(h, r, t, m) - c]^2. \quad (3)$$

To facilitate the optimization progress, we prefer using the loglikelihood of \mathcal{L}_{max} and \mathcal{L}_{min} , and the learning targets can be further processed as follows,

$$\begin{aligned} & \arg \max_{h, r, t, m} \log \mathcal{L}_{max} \\ = & \arg \max_{h, r, t, m} \sum_{\langle h, r, t, m \rangle \in \Delta} \log Pr(h, r, t, m) \\ = & \arg \max_{h, r, t, m} \sum_{\langle h, r, t, m \rangle \in \Delta} \frac{1}{3} [\log Pr(h|r, t) + \log Pr(r|h, t, m) + \log Pr(t|h, r)]; \end{aligned} \quad (4)$$

$$\begin{aligned}
& \arg \min_{h,r,t,m} \log \mathcal{L}_{min} \\
= & \arg \min_{h,r,t,m} \sum_{\langle h,r,t,m,c \rangle \in \Delta} \frac{1}{2} [\log Pr(h, r, t, m) - \log c]^2 \\
= & \arg \min_{h,r,t,m} \sum_{\langle h,r,t,m,c \rangle \in \Delta} \frac{1}{2} \left\{ \frac{1}{3} [\log Pr(h|r, t) + \log Pr(r|h, t, m) + \log Pr(t|h, r)] - \log c \right\}^2.
\end{aligned} \tag{5}$$

The advantage of the conversions above is that we can separate the factors out, compared with Equation (1), and what left for us is to figure out the approaches on modeling $Pr(h|r, t)$, $Pr(r|h, t, m)$ and $Pr(t|h, r)$.

$Pr(r|h, t, m)$ leverages the evidences from two different resources to predict the relation. If the concurrence of the two entities (h and t) in knowledge bases is independent of the appearance of the relation mention m from free texts, we can factorize $Pr(r|h, t, m)$ as shown by Equation (6):

$$Pr(r|h, t, m) = Pr(r|h, t)Pr(r|m). \tag{6}$$

Then we need to consider formulating $Pr(h|r, t)$, $Pr(r|h, t)$, $Pr(t|h, r)$ and $Pr(r|m)$, respectively.

Figure 1(a) illustrates the traditional way of recording knowledge as triplets. The triplets $\langle h, r, t \rangle$ can construct a knowledge graph in which entities (h and t) are nodes and the relation (r) between them is a directed edge from the head entity (h) to the tail entity (t). This kind of symbolic representation, whilst being very efficient for storing, is not flexible enough to statistical learning approaches [11]. But once we project each elements, including entities and relations in the knowledge repository, into the same embedding space, we can use

$$\mathcal{D}(h, r, t) = -\|\mathbf{h} + \mathbf{r} - \mathbf{t}\| + \alpha, \tag{7}$$

a simple vector operation to measure the distance between $\mathbf{h} + \mathbf{r}$ and \mathbf{t} , in which h, r and t are encoded in d dimensional vectors, and α is the bias parameter. To estimate the conditional probability of appearing t given h and r , i.e. $Pr(t|h, r)$,

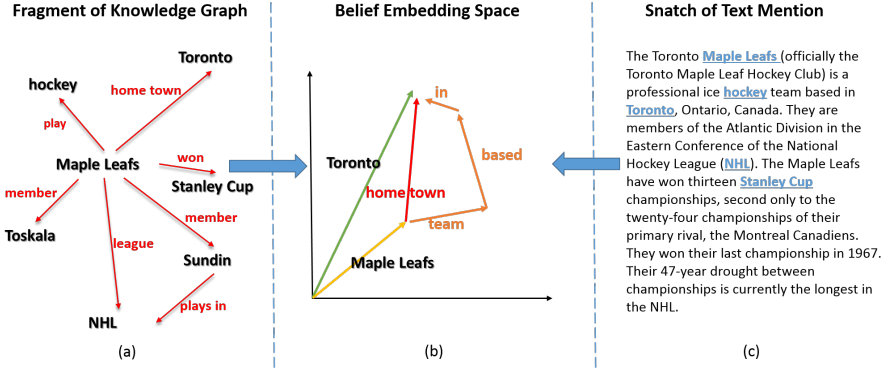


Figure 1: The whole framework of belief embedding. (a) shows a fragment of knowledge graph; (c) is a snatch of Wiki which describes the knowledge graph of (a); (b) illustrates how the belief $\langle Maple Leafs, home town, Toronto, team based in \rangle$ is projected into the same embedding space.

however, we need to adopt the softmax function⁶ as follows,

$$Pr(t|h, r) = \frac{\exp^{\mathcal{D}(h,r,t)}}{\sum_{t' \in E_t} \exp^{\mathcal{D}(h,r,t')}} \quad (8)$$

where E_t is the set of tail entities which contains all possible entities t' appearing in the tail position. Similarly, we can regard $Pr(h|r, t)$ and $Pr(r|h, t)$ as

$$Pr(h|r, t) = \frac{\exp^{\mathcal{D}(h,r,t)}}{\sum_{h' \in E_h} \exp^{\mathcal{D}(h',r,t)}} \quad (9)$$

and

$$Pr(r|h, t) = \frac{\exp^{\mathcal{D}(h,r,t)}}{\sum_{r' \in R} \exp^{\mathcal{D}(h,r',t)}} \quad (10)$$

in which E_h is the set of head entities which contains all possible entities h' appearing in the head position, and R is the set of all candidate relations r' .

One the other hand, Figure 1(c) shows that free texts can provide fruitful contexts between two recognized entities, but the one-hot⁷ feature space is rather high and sparse. Therefore, we can also project each words in relation

⁶http://en.wikipedia.org/wiki/Softmax_function

⁷<http://en.wikipedia.org/wiki/One-hot>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

mentions into the same embedding space of entities and relations. To measure the similarity between the mention m and the corresponding relation r , we adopt inner product of their embeddings as shown by Equation (11),

$$\mathcal{F}(r, m) = \mathbf{W}^T \phi(m) \mathbf{r} + \beta, \quad (11)$$

where \mathbf{W} is the matrix of $\mathbb{R}^{n_v \times d}$ containing n_v vocabularies with d dimensional embeddings, $\phi(m)$ is the sparse one-hot representation of the mention indicating absence or presence of words, $r \in \mathbb{R}^d$ is the embedding of relation r , and β is the bias parameter. Similar to Equation (8), (9) and (10), the conditional probability of predicting relation r given mention m , i.e. $Pr(r|m)$ can be defined as,

$$Pr(r|m) = \frac{\exp^{\mathcal{F}(r,m)}}{\sum_{r' \in R} \exp^{\mathcal{F}(r',m)}}. \quad (12)$$

Above all, we can finally model the probability of a belief via jointly embedding the entities, relations and even the words in mentions as demonstrated by

205 Figure 1(b).

4. Algorithm

To search for the optimal solutions of Equation (4) and (5), we can use *Stochastic Gradient Descent*⁸ (SGD) to update the embeddings of entities, relations and words of mentions in iterative fashion. However, it costs a lot to compute the normalization terms in $Pr(h|r, t)$, $Pr(r|h, t)$, $Pr(t|h, r)$ and $Pr(r|m)$ according to the their definitions made by Equation (8), (9), (10) and (12) respectively. For instance, if we directly calculate the value of $Pr(h|r, t)$ for just one belief, tens of thousands $\exp^{\mathcal{D}(h', r, t)}$ need to be re-valued, as there are tens of thousands candidate entities h' in E_h .

Enlightened by the work of Mikolov et al. [33], we have found an efficient approach that adopts negative sampling technique to approximate the conditional probability functions, i.e. Equation (8), (9), (10) and (12), by being transformed

⁸http://en.wikipedia.org/wiki/Stochastic_gradient_descent

to binary classification problems shown as the subsequent equations respectively,

$$\log Pr(h|r, t) \approx \log Pr(1|h, r, t) + \sum_{i=1}^k \mathbb{E}_{h'_i Pr(h' \in E_h)} \log Pr(0|h'_i, r, t), \quad (13)$$

$$\log Pr(t|h, r) \approx \log Pr(1|h, r, t) + \sum_{i=1}^k \mathbb{E}_{t'_i Pr(t' \in E_t)} \log Pr(0|h, r, t'_i), \quad (14)$$

$$\log Pr(r|h, t) \approx \log Pr(1|h, r, t) + \sum_{i=1}^k \mathbb{E}_{r'_i Pr(r' \in R)} \log Pr(0|h, r'_i, t), \quad (15)$$

$$\log Pr(r|m) \approx \log Pr(1|r, m) + \sum_{i=1}^k \mathbb{E}_{r'_i Pr(r' \in R)} \log Pr(0|r'_i, m), \quad (16)$$

where we sample k negative beliefs and discriminate them from the positive case. For the simple binary classification problems mentioned above, we choose the logistic function with the offset ϵ shown in Equation (17) to estimate the probability that the given triplet $\langle h, r, t \rangle$ is correct:

$$Pr(1|h, r, t) = \frac{1}{1 + \exp^{-\mathcal{D}(h, r, t)}} + \epsilon, \quad (17)$$

and with the offset η shown in Equation (18) to tell the probability of the occurrence of r and m :

$$Pr(1|r, m) = \frac{1}{1 + \exp^{-\mathcal{F}(r, m)}} + \eta. \quad (18)$$

We also display the framework of *PBE* learning algorithm written in pseudocode as shown by Algorithm 1.

5. Experiment

Besides its access to the efficient SGD algorithm, the learnt embeddings by *PBE* can contribute more effectiveness on multiple subtasks of knowledge population, such as entity inference, relation prediction, and triplet classification.

- *Entity inference*: Given a wrecked triplet, like $\langle h, r, ? \rangle$ or $\langle ?, r, t \rangle$, the subtask works on inferring the missing entities to complete the triplet.

ALGORITHM 1 : The Learning Algorithm of PBE

Input:

Training set $\Delta = \{(h, r, t, m, c)\}$, entity set E , relation set R , vocabulary set V of relation mentions; dimension of embeddings d , number of negative samples k , learning rate γ , maximum epochs n ; the bias α and β , the offset ϵ and η .

```
1: foreach  $e \in E$  do
2:    $e := \text{Uniform}(\frac{-6}{\sqrt{d}}, \frac{6}{\sqrt{d}})$ 
3: end foreach
4: foreach  $r \in R$  do
5:    $r := \text{Uniform}(\frac{-6}{\sqrt{d}}, \frac{6}{\sqrt{d}})$ 
6: end foreach
7: foreach  $v \in V$  do
8:    $v := \text{Uniform}(\frac{-6}{\sqrt{d}}, \frac{6}{\sqrt{d}})$ 
9: end foreach
10:  $i := 0$ 
11: while  $i < n$  do
12:   foreach  $\langle h, r, t, m, c \rangle \in \Delta$  do
13:     foreach  $j \in \text{range}(k)$  do
14:       Negative sampling:  $\langle h'_j, r, t, m \rangle \in \Delta'_h$ 
15:       /* $\Delta'_h$  is the set of  $k$  negative beliefs replacing  $h^*$ */
16:       Negative sampling:  $\langle h, r'_j, t, m \rangle \in \Delta'_r$ 
17:       Negative sampling:  $\langle h, r, t'_j, m \rangle \in \Delta'_t$ 
18:     end foreach
19:     Gradient ascent:  $\sum_{h,r,t,h',r',t',v \in m} \nabla \log Pr(h, r, t, m)$  according to Equation (4)
20:     OR
21:     Gradient descent:  $\sum_{h,r,t,h',r',t',v \in m} \nabla [\log Pr(h, r, t, m) - \log c]^2$  according to Equation (5)
22:     /*Updating embeddings of  $\langle h, r, t, m \rangle \in \Delta$ ;  $\langle h', r, t, m \rangle \in \Delta'_h$ ;  $\langle h, r', t, m \rangle \in \Delta'_r$ ;  $\langle h, r, t', m \rangle \in \Delta'_t$  with  $\gamma$  and the batch gradients derived from Equation (13), (14), (15) and (16).*/
23:   end foreach
24:    $i++$ 
25: end while
```

Output:

All the embeddings of h, t, r and v , where $h, t \in E$, $r \in R$ and $v \in V$.

- *Relation prediction*: Given a pair of entities and the text mentions indicating the semantic relations between them, i.e. $\langle h, ?, t, m \rangle$, this subtask predicts the best relations of the two entities.
- *Triplet classification*: It tells whether a completed triplet is correct or not $(\langle h, r, t \rangle? 1 : 0)$.

5.1. Entity inference

One of the benefits of knowledge embedding is that simple vector operations can apply to entity inference which contributes to knowledge graph completion. For example, if we would like to know which entity $h \in E_h$ is the exact head entity given the relation r and the entity t , we just need to compute the $\arg \max_{h \in E_h} Pr(h|r, t)$, with the help of the entity and relation embeddings. In the meanwhile, $\arg \max_{t \in E_t} Pr(t|h, r)$ will help us to find the best tail entity given the head entity h and the relation r .

5.1.1. Dataset

To demonstrate the wide adaptability of our approach, we prepare four datasets, i.e. **NELL-50K**, **WN-100K**, **FB-500K** and **NELL-1M** from the repositories of NELL [34], WordNet [3] and Freebase [4, 5], with varies scales as shown by Table 1. The NELL [35] designed and maintained by Carnegie Mellon University is an outstanding system which runs 24 hours/day and never stops learning the beliefs on the Web. Since the starting date of January 2010, it has acquired a knowledge repository with over 80 million confidence-weighted beliefs so far. The dataset **NELL-50K** we adopt, contains about fifty thousand training beliefs from NELL, and each belief has been validated to be true. We also extract a much larger one (**NELL-1M**) with one million training examples from NELL, where each belief is automatically learnt by machine and weighted ranging (0.5, 1.0). **WN-100K** is made by experts from WordNet, and owns only 11 kinds of relations but much more entities. Therefore, it is a sparse repository

in which fewer entities have connections. The last dataset (**FB-500K**⁹) we use was released by Bordes et al. [13]. It is a large but dense, crowdsourcing dataset extracted from Freebase, in which almost every two entities have connections, and each belief is a triplet without a confidence score.

DATASET	NELL-50K	WN-100K	FB-500K	NELL-1M
#(ENTITIES)	29,904	38,696	14,951	82,691
#(RELATIONS)	233	11	1,345	218
#(TRAINING EX.)	57,356	112,581	483,142	1,000,000
#(VALIDATING EX.)	10,710	5,218	50,000	24,864
#(TESTING EX.)	10,711	21,088	59,071	24,863

Table 1: Statistics of the datasets used for the entity inference task.

Table 1 shows the statistics of these four datasets. The statistical characteristic of these datasets are different, which may lead to the variety of tuning parameters.

5.1.2. Metric

For each testing belief, all the other entities that appear in the training set take turns to replace the head entity. Then we get a bunch of candidate triplets. The plausibility of each candidate triplet is firstly computed by various scoring functions, such as $Pr(h|r, t)$ in *PBE*, and then sorted in ascending order. Finally, we locate the ground-truth triplet and record its rank. This whole procedure runs in the same way when replacing the tail entity, so that we can gain the mean results. We use two metrics, i.e. *Mean Rank* and *Mean Hit@10* (the proportion of ground truth triplets that rank in Top 10), to measure the performance. However, the results measured by those metrics are relatively *raw*, as the procedure above tends to generate false negative triplets. In other

⁹We change the original name of the dataset (**FB15K**), so as to follow the naming conventions in our paper. Related studies on this dataset can be looked up from the website <https://www.hds.utc.fr/everest/doku.php?id=en:transe>

words, some of the candidate triplets rank rather higher than the ground truth triplet just because they also appear in the training set. We thus filter out those triplets to report more reasonable results.

DATASET	NELL-50K			
METRIC	MEAN RANK		MEAN HIT@10	
	Raw	Filter	Raw	Filter
TransE [13]	2,436 / 29,904	2,426 / 29,904	18.9%	19.6%
TransM [30]	2,296 / 29,904	2,285 / 29,904	20.5%	21.3%
TransH [31]	2,185 / 29,904	2,072 / 29,904	21.6%	28.8%
PBE	2,078 / 29,904	1,996 / 29,904	22.5%	26.4%

Table 2: Entity inference results on the **NELL-50K** dataset. We compared our proposed PBE with the state-of-the-art method TransH and other prior arts mentioned in Section 2.2.

DATASET	WN-100K			
METRIC	MEAN RANK		MEAN HIT@10	
	Raw	Filter	Raw	Filter
TransE [13]	10,623 / 38,696	10,575 / 38,696	3.8%	4.1%
TransM [30]	14,586 / 38,696	13,276 / 38,696	1.8%	2.0%
TransH [31]	12,542 / 38,696	12,463 / 38,696	2.3%	2.6%
PBE	8,462 / 38,696	8,409 / 38,696	9.0%	10.1%

Table 3: Entity inference results on the **WN-100K** dataset. We compared our proposed PBE with the state-of-the-art method TransH and other prior arts mentioned in Section 2.2.

5.1.3. Performance

We compare *PBE* with the state-of-the-art *TransH*, *TransM*, *TransE* mentioned in Section 2.2 via evaluating their performances on **NELL-50K**, **WN-100K**, **FB-500K**, and **NELL-1M** datasets. We tune the parameters of each previous model based on the validation set, and select the combination of parameters which leads to the best performance. To make responsible compar-

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

DATASET	FB-500K			
METRIC	MEAN RANK		MEAN HIT@10	
	Raw	Filter	Raw	Filter
TransE [13]	243 / 14,951	125 / 14,951	34.9%	47.1%
TransM [30]	196 / 14,951	93 / 14,951	44.6%	55.2%
TransH [31]	211 / 14,951	84 / 14,951	42.5%	58.5%
PBE	165 / 14,951	61 / 14,951	50.5%	64.6%

Table 4: Entity inference results on the **FB-500K** dataset. We compared our proposed PBE with the state-of-the-art method TransH and other prior arts mentioned in Section 2.2.

DATASET	NELL-1M			
METRIC	MEAN RANK		MEAN HIT@10	
	Raw	Filter	Raw	Filter
TransE [13]	29,059 / 82,691	29,052 / 82,691	6.5%	6.6%
TransM [30]	28,435 / 82,691	28,129 / 82,691	5.4%	5.5%
TransH [31]	27,455 / 82,691	26,980 / 82,691	7.8%	8.7%
PBE	7,528 / 82,691	7,485 / 82,691	8.7%	9.0%

Table 5: Entity inference results on the **NELL-1M** dataset. We compared our proposed PBE with the state-of-the-art method TransH and other prior arts mentioned in Section 2.2.

isons between *PBE* and the state-of-the-art approach *TransH*, we request its authors [31] to re-evaluate their system with all the four datasets and to report the best results. For *PBE*, we tried several combinations of parameters: $d = \{20, 50, 100\}$, $\gamma = \{0.1, 0.05, 0.01, 0.005\}$, and $norm = \{L_1, L_2\}$, and finally chose $d = 50$, $\gamma = 0.01$, $norm = L_2$ for **NELL-50K** and **WN-100K** datasets, and $d = 100$, $\gamma = 0.01$, $norm = L_2$ for **FB-500K** and **NELL-1M** datasets to conduct further experiments.

Table 2, 3, 4 and 5 demonstrate that *PBE* outperforms all the state-of-the-arts, including *TransE* [13], *TransM* [30] and *TransH* [31], and achieves significant improvements on all datasets. Overall, The *relative increments* performed by *PBE* compared with the best results of prior arts under all metrics are as subsequence,

- **NELL-50K**: {*Mean Rank Raw*: 4.9% \uparrow , *Hit@10 Raw*: 4.2% \uparrow , *Mean Rank Filter*: 3.7% \uparrow , *Hit@10 Filter*: 8.3% \downarrow };
- **WN-100K**: {*Mean Rank Raw*: 20.3% \uparrow , *Hit@10 Raw*: 136.8% \uparrow , *Mean Rank Filter*: 20.5% \uparrow , *Hit@10 Filter*: 146.3% \uparrow };
- **FB-500K**: {*Mean Rank Raw*: 15.8% \uparrow , *Hit@10 Raw*: 27.3% \uparrow , *Mean Rank Filter*: 13.3% \uparrow , *Hit@10 Filter*: 10.4% \uparrow };
- **NELL-1M**: {*Mean Rank Raw*: 72.5% \uparrow , *Hit@10 Raw*: 11.5% \uparrow , *Mean Rank Filter*: 72.2% \uparrow , *Hit@10 Filter*: 3.4% \uparrow }

5.2. Relation prediction

The scenario of this subtask is that: given a pair of entities and a short text/mention indicating the correct relations, we compute the $\arg \max_{r \in R} Pr(r|h, t)Pr(r|m)$ to predict the best relations.

5.2.1. Dataset

We continue using the datasets mentioned in Section 5.1 to compare the performances among all the competing methods. But, as the words in relation

1
 2
 3
 4
 5
 6
 7
 8
 9
 10
 11
 12
 13
 14
 15
 16
 17
 18
 19
 20
 21
 22
 23
 24
 25
 26
 27
 28
 29
 30
 31
 32
 33
 34
 35
 36
 37
 38
 39
 40
 41
 42
 43
 44
 45
 46
 47
 48
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

mentions are additionally concerned the in this subtask, we also show the vo-
 cabulary size of relation mentions in each dataset in Table 6 as follows, except
 for **WN-100K** and **FB-500K** which only contain triplets as beliefs, and the
 sizes of their vocabulary are null.

DATASET	NELL-50K	WN-100K	FB-500K	NELL-1M
#(ENTITIES)	29,904	38,696	14,951	82,691
#(RELATIONS)	233	11	1,345	218
#(VOCABULARY)	8,948	-	-	12,354
#(TRAINING EX.)	57,356	112,581	483,142	1,000,000
#(VALIDATING EX.)	10,710	5,218	50,000	24,864
#(TESTING EX.)	10,711	21,088	59,071	24,863

Table 6: Statistics of the datasets used for the relation prediction task.

5.2.2. Metric

We compare the performances between our models and other state-of-the-
 art approaches mentioned in Section 2.2 and 2.3, including *TransE* [13], *TransM*
 [30], *TransH* [31] and *JRME* [16], with the metrics as follows,

- *Average Rank*: Each candidate relation will gain a score calculated by
 Equation (7). We sort them in ascent order and compare with the corre-
 sponding ground-truth belief. For each belief in the testing set, we get the
 rank of the correct relation. The average rank is an aggregative indicator,
 to some extent, to judge the overall performance on relation extraction of
 an approach.
- *Hit@10*: Besides the average rank, scientists from the industrials concern
 more about the accuracy of extraction when selecting Top10 relations.
 This metric shows the proportion of beliefs that we predict the correct
 relation ranked in Top10.

- *Hit@1*: It is a more strict metric that can be referred by automatic system, since it demonstrates the accuracy when just picking the first predicted relation in the sorted list.

DATASET	NELL-50K		
METRIC	AVG. R.	HIT@10	HIT@1
TransE [13]	131.8 / 233	16.3%	3.0%
TransM [30]	70.2 / 233	18.9%	4.3%
TransH [31]	46.3 / 233	20.0%	5.1%
JRME [16]	6.2 / 233	87.8%	60.2%
PBE	2.5 / 233	96.6%	78.3%

Table 7: Performance of relation prediction on TransE, TransM, TransH, JRME and PBE evaluated by the metrics of Average Rank, Hit@10 and Hit@1 with **NELL-50K** dataset.

DATASET	WN-100K		
METRIC	AVG. R.	HIT@10	HIT@1
TransE [13]	3.8 / 11	98.3%	15.1%
TransM [30]	4.6 / 11	97.5%	14.8%
TransH [31]	3.4 / 11	99.0%	19.3%
JRME [16]	3.9 / 11	99.0%	15.9%
PBE	2.0 / 11	99.1%	72.6%

Table 8: Performance of relation prediction on TransE, TransM, TransH, JRME and PBE evaluated by the metrics of Average Rank, Hit@10 and Hit@1 with **WN-100K** dataset.

325 5.2.3. Performance

Table 7, 8, 9 and 10 illustrate the results of experiments on relation prediction with **NELL-50K**, **WN-100K**, **FB-500K** and **NELL-1M** datasets, respectively. All of them show that *PBE* performs best compared with all the latest approaches including the state-of-the-art *JRME* [16]. The relative increments are

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

DATASET	FB-500K		
METRIC	AVG. R.	HIT@10	HIT@1
TransE [13]	762.7 / 1,345	7.3%	1.9%
TransM [30]	402.3 / 1,345	13.4%	3.2%
TransH [31]	79.5 / 1,345	39.2%	15.6%
JRME [16]	60.9 / 1,345	27.4%	7.2%
PBE	2.6 / 1,345	97.3%	66.7%

Table 9: Performance of relation prediction on TransE, TransM, TransH, JRME and PBE evaluated by the metrics of Average Rank, Hit@10 and Hit@1 with **FB-500K** dataset.

DATASET	NELL-1M		
METRIC	AVG. R.	HIT@10	HIT@1
TransE [13]	70.4 / 218	5.4%	0.4%
TransM [30]	65.5 / 218	18.7%	3.4%
TransH [31]	62.9 / 218	26.8%	5.8%
JRME [16]	7.0 / 218	89.0%	54.5%
PBE	5.8 / 218	92.1%	65.0%

Table 10: Performance of relation prediction on TransE, TransM, TransH, JRME and PBE evaluated by the metrics of Average Rank, Hit@10 and Hit@1 with **NELL-1M** dataset.

- **NELL-50K**: { *Mean Rank*: 59.7% \uparrow , *Hit@10*: 10.0% \uparrow , *Hit@1*: 30.0% \uparrow };
- **WN-100K**: { *Mean Rank*: 41.1% \uparrow , *Hit@10*: 0.1% \uparrow , *Hit@1*: 276.2% \uparrow };
- **FB-500K**: { *Mean Rank*: 95.7% \uparrow , *Hit@10*: 148.2% \uparrow , *Hit@1*: 327.6% \uparrow };
- **NELL-1M**: { *Mean Rank*: 20.6% \uparrow , *Hit@10*: 3.5% \uparrow , *Hit@1*: 19.3% \uparrow }.

Moreover, the leading results of *PBE* and *JRME* on **NELL** datasets also inspire us that text mentions can contribute a lot on predicting the correct relations.

5.3. Triplet classification

Triplet classification is another inference related task proposed by Socher et al. [27] which focuses on searching a relation-specific threshold σ_r to identify whether a triplet $\langle h, r, t \rangle$ is plausible. If the probability of a testing triplet (h, r, t) computed by $Pr(h|r, t)Pr(r|h, t)Pr(t|h, r)$ is below the relation-specific threshold σ_r , it is predicted as positive, otherwise negative.

5.3.1. Dataset

It is emphasized that the head or the tail entity can be randomly replaced with another one to produce a negative training example, but in order to build much tough validation and testing datasets, we constrain that the picked entity should once appear at the same position. For example, $(Pablo\ Picaso, nationality, U.S.)$ is a potential negative example rather than the obvious nonsense $(Pablo\ Picaso, nationality, Van\ Gogh)$, given a positive triplet $(Pablo\ Picaso, nationality, Spain)$. Table 11 shows the statistics of the standard datasets that we used for evaluating models on the triplet classification subtask.

5.3.2. Metric

We use three metrics, i.e. *Classification Accuracy*, *Precision-recall Curve* and *Area Under Curve (AUC)*, to measure the performances among the competing methods.

DATASET	NELL-50K	WN-100K	FB-500K	NELL-1M
#(ENTITIES)	29,904	38,696	14,951	82,691
#(RELATIONS)	233	11	1,345	218
#(TRAINING EX.)	57,356	112,581	483,142	1,000,000
#(TC VALIDATING EX.)	21,420	10,436	100,000	49,728
#(TC TESTING EX.)	21,412	42,176	118,142	49,714

Table 11: Statistics of the datasets used for the triplet classification task.

- 360
 • *Classification Accuracy*: We sum up the correctness of each triplet $\langle h, r, t \rangle$ via comparing the probability of the triplet and the relation-specific threshold σ_r , which can be searched via maximizing the classification accuracy on the validation triplets which belong to the relation r .
- 365
 • *Precision-recall Curve*: It measures the global performance of classification by sorting all the triplets based on their estimated probability. We consider the positive testing triplets and draw the precision-recall curve for each approach.
- 370
 • *Area Under Curve (AUC)*: The AUC is a commonly used evaluation metric for binary classification problems like predicting a Buy or Sell decision (binary decision). The interpretation here is that given a random positive triplet and a negative triplet, the AUC gives the proportion of the time we guess which is which correctly. It is less affected by sample balance than accuracy. A perfect model will score an AUC of 1.0, while random guessing will score an AUC of around 0.5, a meager 50% chance on each other.

375 5.3.3. Performance

We use the best combination of parameter settings in the entity inference task: $d = 100$, $\gamma = 0.01$, $norm = L_2$ to generate the entity and relation embeddings, and learn the best classification threshold σ_r for each relation r .

DATASET	NELL-50K	WN-100K	FB-500K	NELL-1M
METRIC	<i>ACC.</i>	<i>ACC.</i>	<i>ACC.</i>	<i>ACC.</i>
TransE [13]	80.5%	64.2%	79.9%	64.0%
TransM [30]	82.0%	57.2%	85.8%	64.8%
TransH [31]	83.6%	59.5%	87.7%	67.0%
PBE	90.2%	67.8%	92.6%	86.2%

Table 12: The accuracy of triplet classification compared among several latest approaches: *PBE*, *TransH*, *TransM* and *TransE*.

DATASET	NELL-50K	WN-100K	FB-500K	NELL-1M
METRIC	<i>AUC</i>	<i>AUC</i>	<i>AUC</i>	<i>AUC</i>
TransE [13]	0.623	0.674	0.645	0.547
TransM [30]	0.683	0.610	0.772	0.558
TransH [31]	0.681	0.613	0.744	0.596
PBE	0.942	0.786	0.936	0.786

Table 13: The AUC of triplet classification compared among several latest approaches: *PBE*, *TransH*, *TransM* and *TransE*.

Compared with several of the latest approaches, i.e. *TransH* [31], *TransM* [30] and *TransE* [13], the proposed *PBE* approach still outperforms them within the metrics of *Classification Accuracy (ACC.)* and *Area Under Curve (AUC)*, as shown in Table 12 and 13. We also draw the precision-recall curves which indicate the capability of global discrimination by ranking the distance of all the testing triplets, and Figure 2, 3, 4 and 5 can intuitively show that *PBE* performs much better than the other approaches.

Compared with several of the latest approaches, i.e. *TransH* [31], *TransM* [30] and *TransE* [13], the proposed *PBE* approach outperforms with the relative improvements that

- **NELL-50K**: {*Accuracy*: 7.9% \uparrow , *AUC*: 37.9% \uparrow };
- **WN-100K**: {*Accuracy*: 5.6% \uparrow , *AUC*: 16.6% \uparrow };
- **FB-500K**: {*Accuracy*: 5.6% \uparrow , *AUC*: 21.2% \uparrow };
- **NELL-1M**: {*Accuracy*: 28.6% \uparrow , *AUC*: 31.8% \uparrow }.

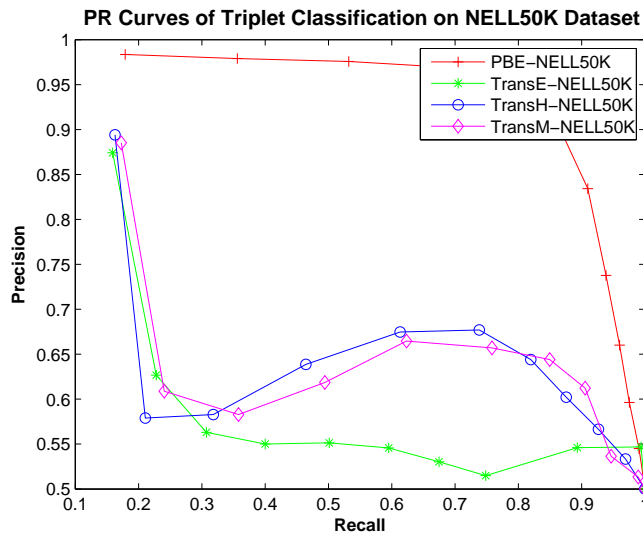


Figure 2: The precision-recall curves for triplet classification of *PBE*, *TransH*, *TransM* and *TransE* on **NELL-50K** dataset.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

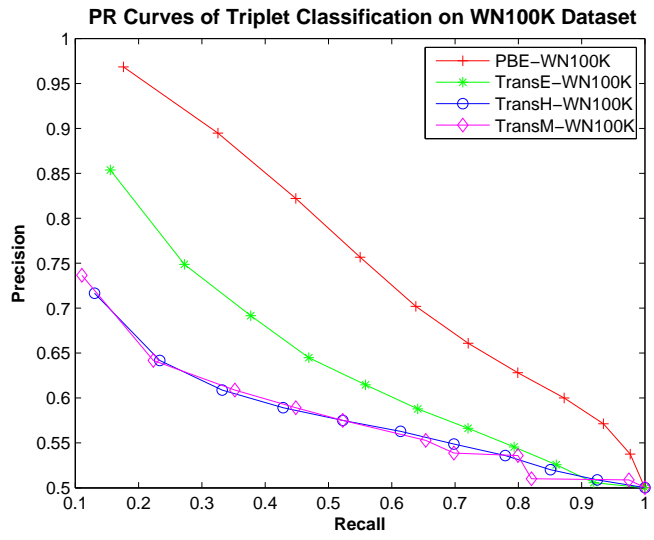


Figure 3: The precision-recall curves for triplet classification of *PBE*, *TransH*, *TransM* and *TransE* on **WN-100K** dataset.

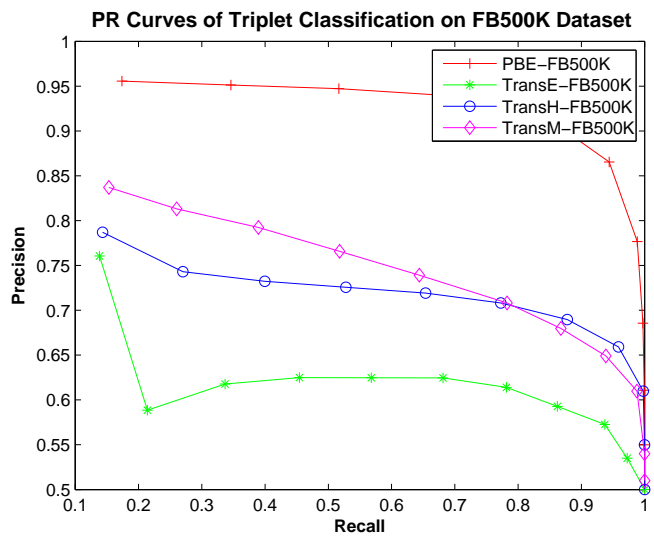


Figure 4: The precision-recall curves for triplet classification of *PBE*, *TransH*, *TransM* and *TransE* on **FB-500K** dataset.

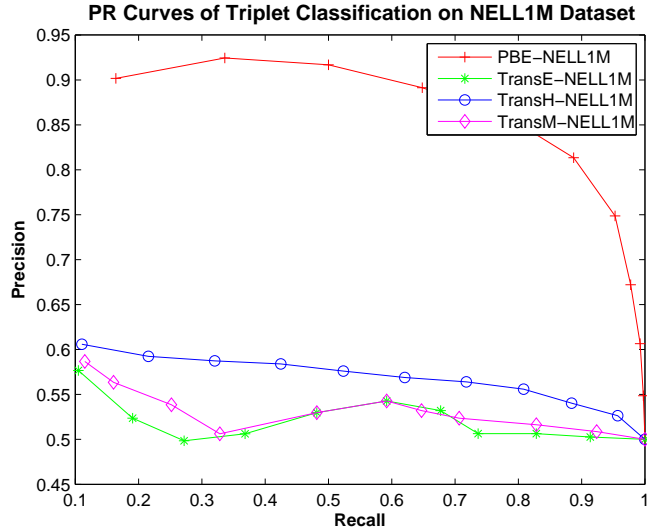


Figure 5: The precision-recall curves for triplet classification of *PBE*, *TransH*, *TransM* and *TransE* on **NELL-1M** dataset.

6. Conclusion

We challenge the problem of embedding beliefs which both contain structured knowledge and unstructured free texts and propose an elegant probabilistic model to tackle this issue at the first attempt by measuring the probability of a given belief $\langle h, r, t, m \rangle$. To efficiently learn the embeddings for each entity, relation, and word in mentions, we also adopt the negative sampling technique to transform the original model and display the algorithm based on stochastic gradient descend (SGD) to search the optimal solution. Extensive experiments on knowledge population including *entity inference*, *relation prediction* and *triplet classification* show that our approach achieves significant improvement on three large-scale repositories, compared with state-of-the-art methods. You can access to all the datasets through the publish link of OneDrive: <http://1drv.ms/1IDwZAR>.

We are pleased to see further improvements of the proposed model, which leaves open promising directions for the future work, such as taking advantage of

1
2
3
4
5
6
7
8
9 the probabilistic belief embeddings to enhance the studies of text summarization
10 and open-domain question answering.
11
12

13
14 ⁴¹⁰ **Acknowledgement**

15
16 The paper is dedicated to all the members of CSLT¹⁰ and Proteus Group
17
18 ¹¹. It was supported by National Program on Key Basic Research Project (973
19 Program) under Grant 2013CB329304, National Science Foundation of China
20 (NSFC) under Grant No.61433018 and No.61373075, and Chinese Scholarship
21 Council, when the first author was a joint-supervision Ph.D. candidate of Ts-
22 ⁴¹⁵ inghua University and New York University.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

55

¹⁰<http://cslt.riit.tsinghua.edu.cn/>

56 ¹¹<http://nlp.cs.nyu.edu/index.shtml>

1
2
3
4
5
6
7
8
9 **References**

- 10
11 [1] R. Grishman, Information extraction: Techniques and challenges, in: Inter-
12 national Summer School on Information Extraction: A Multidisciplinary
13 Approach to an Emerging Information Technology, SCIE '97, Springer-
14 420 Verlag, London, UK, UK, 1997, pp. 10–27.
15 URL <http://dl.acm.org/citation.cfm?id=645856.669801>
16
17
18
19
20 [2] S. Sarawagi, Information extraction, Foundations and trends in databases
21 1 (3) (2008) 261–377.
22
23
24 425 [3] G. A. Miller, Wordnet: a lexical database for english, Communications of
25 the ACM 38 (11) (1995) 39–41.
26
27
28 [4] K. Bollacker, R. Cook, P. Tufts, Freebase: A shared database of structured
29 general human knowledge, in: AAAI, Vol. 7, 2007, pp. 1962–1963.
30 URL <http://www.aaai.org/Papers/AAAI/2007/AAAI07-355.pdf>
31
32
33 430 [5] K. Bollacker, C. Evans, P. Paritosh, T. Sturge, J. Taylor, Freebase: a
34 collaboratively created graph database for structuring human knowledge,
35 in: Proceedings of the 2008 ACM SIGMOD international conference on
36 Management of data, ACM, 2008, pp. 1247–1250.
37 URL <http://dl.acm.org/citation.cfm?id=1376746>
38
39
40
41 435 [6] G. Klyne, J. J. Carroll, Resource description framework (rdf): Concepts
42 and abstract syntax, W3C Recommendation.
43
44
45
46 [7] A. Carlson, J. Betteridge, B. Kisiel, B. Settles, E. R. H. Jr., T. M. Mitchell,
47 Toward an architecture for never-ending language learning, in: Proceedings
48 of the Twenty-Fourth Conference on Artificial Intelligence (AAAI 2010),
49 2010.
50 440
51
52 [8] R. West, E. Gabrilovich, K. Murphy, S. Sun, R. Gupta, D. Lin, Knowledge
53 base completion via search-based question answering, in: WWW, 2014.
54 URL <http://www.cs.ubc.ca/~murphyk/Papers/www14.pdf>
55
56
57
58

- 1
2
3
4
5
6
7
8
9 [9] M. Mintz, S. Bills, R. Snow, D. Jurafsky, Distant supervision for relation
10 extraction without labeled data, in: Proceedings of the Joint Conference
11 445 of the 47th Annual Meeting of the ACL and the 4th International Joint
12 Conference on Natural Language Processing of the AFNLP: Volume 2-
13 Volume 2, Association for Computational Linguistics, 2009, pp. 1003–1011.
14
15
16
17 [10] M. Fan, D. Zhao, Q. Zhou, Z. Liu, T. F. Zheng, E. Y. Chang, Distant su-
18 pervision for relation extraction with matrix completion, in: Proceedings
19 450 of the 52nd Annual Meeting of the Association for Computational Linguis-
20 tics (Volume 1: Long Papers), Association for Computational Linguistics,
21 Baltimore, Maryland, 2014, pp. 839–849.
22
23 URL <http://www.aclweb.org/anthology/P14-1079>
24
25
26
27 [11] A. Bordes, J. Weston, R. Collobert, Y. Bengio, et al., Learning structured
28 embeddings of knowledge bases., in: AAAI, 2011.
29
30 URL [http://www.aaai.org/ocs/index.php/AAAI/AAAI11/paper/
31 viewPDFInterstitial/3659/3898](http://www.aaai.org/ocs/index.php/AAAI/AAAI11/paper/viewPDFInterstitial/3659/3898)
32
33
34 [12] A. Bordes, X. Glorot, J. Weston, Y. Bengio, A semantic matching energy
35 function for learning with multi-relational data, Machine Learning 94 (2)
36 460 (2014) 233–259.
37
38 URL [http://link.springer.com/article/10.1007/
39 s10994-013-5363-6](http://link.springer.com/article/10.1007/s10994-013-5363-6)
40
41
42 [13] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, O. Yakhnenko, Trans-
43 lating embeddings for modeling multi-relational data, in: Advances in Neu-
44 465 ral Information Processing Systems, 2013, pp. 2787–2795.
45
46
47 [14] J. Weston, A. Bordes, O. Yakhnenko, N. Usunier, Connecting language and
48 knowledge bases with embedding models for relation extraction, in: Pro-
49 ceedings of the 2013 Conference on Empirical Methods in Natural Language
50 Processing, Association for Computational Linguistics, Seattle, Washing-
51 470 ton, USA, 2013, pp. 1366–1371.
52
53 URL <http://www.aclweb.org/anthology/D13-1136>
54
55
56
57
58

- 1
2
3
4
5
6
7
8
9 [15] Z. Wang, J. Zhang, J. Feng, Z. Chen, Knowledge graph and text jointly
10 embedding, in: Proceedings of the 2014 Conference on Empirical Methods
11 in Natural Language Processing (EMNLP), Association for Computational
12 475 Linguistics, 2014, pp. 1591–1601.
13 URL <http://aclweb.org/anthology/D14-1167>
14
15
16
17 [16] M. Fan, K. Cao, Y. He, R. Grishman, Jointly embedding relations and
18 mentions for knowledge population, arXiv preprint arXiv:1504.01683.
19
20
21 [17] N. Kambhatla, Combining lexical, syntactic, and semantic features with
22 480 maximum entropy models for extracting relations, in: Proceedings of the
23 ACL 2004 on Interactive poster and demonstration sessions, Association
24 for Computational Linguistics, 2004, p. 22.
25
26
27
28 [18] D. Zelenko, C. Aone, A. Richardella, Kernel methods for relation extrac-
29 tion, The Journal of Machine Learning Research 3 (2003) 1083–1106.
30 485
31
32 [19] Z. GuoDong, S. Jian, Z. Jie, Z. Min, Exploring various knowledge in re-
33 lation extraction, in: Proceedings of the 43rd Annual Meeting on As-
34 sociation for Computational Linguistics, ACL '05, Association for Com-
35 putational Linguistics, Stroudsburg, PA, USA, 2005, pp. 427–434. doi:
36 10.3115/1219840.1219893.
37 490
38 URL <http://dx.doi.org/10.3115/1219840.1219893>
39
40
41
42 [20] S. Riedel, L. Yao, A. McCallum, Modeling relations and their mentions
43 without labeled text, in: Machine Learning and Knowledge Discovery in
44 Databases, Springer, 2010, pp. 148–163.
45
46
47
48 [21] R. Hoffmann, C. Zhang, X. Ling, L. Zettlemoyer, D. S. Weld, Knowledge-
49 based weak supervision for information extraction of overlapping relations,
50 in: Proceedings of the 49th Annual Meeting of the Association for Compu-
51 tational Linguistics: Human Language Technologies-Volume 1, Association
52 for Computational Linguistics, 2011, pp. 541–550.
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
- 500 [22] M. Surdeanu, J. Tibshirani, R. Nallapati, C. D. Manning, Multi-instance multi-label learning for relation extraction, in: Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, Association for Computational Linguistics, 2012, pp. 455–465.
- 505 [23] J. R. Quinlan, R. M. Cameron-Jones, Foil: A midterm report, in: Proceedings of the European Conference on Machine Learning, ECML '93, Springer-Verlag, London, UK, UK, 1993, pp. 3–20.
URL <http://dl.acm.org/citation.cfm?id=645323.649599>
- [24] N. Lao, W. W. Cohen, Relational retrieval using a combination of path-constrained random walks, *Machine learning* 81 (1) (2010) 53–67.
- 510 [25] N. Lao, T. Mitchell, W. W. Cohen, Random walk inference and learning in a large scale knowledge base, in: Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Edinburgh, Scotland, UK., 2011, pp. 529–539.
URL <http://www.aclweb.org/anthology/D11-1049>
- 515 [26] M. Gardner, P. P. Talukdar, B. Kisiel, T. M. Mitchell, Improving learning and inference in a large knowledge-base using latent syntactic cues., in: EMNLP, ACL, 2013, pp. 833–838.
URL <http://dblp.uni-trier.de/db/conf/emnlp/emnlp2013.html#GardnerTKM13>
- 520 [27] R. Socher, D. Chen, C. D. Manning, A. Ng, Reasoning with neural tensor networks for knowledge base completion, in: Advances in Neural Information Processing Systems, 2013, pp. 926–934.
- [28] I. Sutskever, R. Salakhutdinov, J. B. Tenenbaum, Modelling relational data using bayesian clustered tensor factorization., in: NIPS, 2009, pp. 1821–1828.
- 525

- 1
2
3
4
5
6
7
8
9 [29] R. Jenatton, N. Le Roux, A. Bordes, G. Obozinski, et al., A latent factor
10 model for highly multi-relational data., in: NIPS, 2012, pp. 3176–3184.
11
- 12 [30] M. Fan, Q. Zhou, E. Chang, T. F. Zheng, Transition-based knowledge graph
13 embedding with relational mapping properties, in: Proceedings of the 28th
14 Pacific Asia Conference on Language, Information, and Computation, 2014,
530 pp. 328–337.
15
16
17
18
- 19 [31] Z. Wang, J. Zhang, J. Feng, Z. Chen, Knowledge graph embedding by trans-
20 lating on hyperplanes, in: Proceedings of the Twenty-Eighth AAAI Con-
21 ference on Artificial Intelligence, July 27 -31, 2014, Québec City, Québec,
22 Canada., 2014, pp. 1112–1119.
23 535
24 URL [http://www.aaai.org/ocs/index.php/AAAI/AAAI14/paper/view/](http://www.aaai.org/ocs/index.php/AAAI/AAAI14/paper/view/8531)
25 [8531](http://www.aaai.org/ocs/index.php/AAAI/AAAI14/paper/view/8531)
26
27
28
- 29 [32] M. Fan, Q. Zhou, T. F. Zheng, Learning embedding representations
30 for knowledge inference on imperfect and incomplete repositories, arXiv
31 preprint arXiv:1503.08155.
32 540
33
34
- 35 [33] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed rep-
36 resentations of words and phrases and their compositionality, in: C. Burges,
37 L. Bottou, M. Welling, Z. Ghahramani, K. Weinberger (Eds.), Advances in
38 Neural Information Processing Systems 26, 2013, pp. 3111–3119.
39 545
40
- 41 [34] A. Carlson, J. Betteridge, B. Kisiel, B. Settles, E. R. H. Jr., T. M. Mitchell,
42 Toward an architecture for never-ending language learning, in: Proceedings
43 of the Twenty-Fourth Conference on Artificial Intelligence (AAAI 2010),
44 2010.
45
46
47
- 48 [35] T. Mitchell, W. Cohen, E. Hruschka, P. Talukdar, J. Betteridge, A. Carlson,
49 B. Dalvi, M. Gardner, B. Kisiel, J. Krishnamurthy, N. Lao, K. Mazaitis,
50 T. Mohamed, N. Nakashole, E. Platanios, A. Ritter, M. Samadi, B. Set-
51 tles, R. Wang, D. Wijaya, A. Gupta, X. Chen, A. Saparov, M. Greaves,
52 J. Welling, Never-ending learning, in: Proceedings of the Twenty-Ninth
53 AAAI Conference on Artificial Intelligence (AAAI-15), 2015.
54 555
55
56
57
58