

# Knowledge Graph Embedding With Attentional Triple Context

Huan Gao<sup>1,2</sup>, Jun Shi<sup>1,2</sup>, Guilin Qi<sup>3</sup>

<sup>1</sup>School of Computer Science and Engineering, Southeast University, Nanjing, China

## Abstract

Because representation learning of knowledge graphs increasingly adds value to various applications that require machines to recognize and understand queries and their semantics, Knowledge graph embedding has increasingly gained attention for statistical modeling of knowledge graph. Existing methods treat each triple independent, but, cannot handle the graph structural information, such as connectivity pathes and neighbor entities, which contains rich information for inference. In this paper, we proposes an Attentional-Triple-Context-based knowledge Embedding model(ATCE). For each triple, two kinds of graph structural information are considered as its context, which we refer to as *triple context* : 1) Neighbor context is the outgoing relations and neighboring entities of an entity; 2) Path context is connective paths between a pair of entities, both of which contains rich useful and unrelated information for entities and relations. ATCE learns embedding for entities and relations with an attention mechanism which is designed to choose the useful inference information in triple context. The experimental results show that our methods outperforms the state-of-the-art methods for link prediction and triple classification.

## 1 Introduction

Recent advances in information extraction have led to huge Knowledge graphs(KGs), such as DBpedia [Auer *et al.*, 2007], YAGO [Suchanek *et al.*, 2007] and NELL [Mitchell *et al.*, 2015]. These KGs contain facts which represent relations between entities as triples  $\langle h, r, t \rangle$ . A triple indicate that entities  $h$  and  $t$  are connected by relation  $r$ . Even a KG contains a very large number of triples, it is still far from complete. The completeness of KGs damage their usefulness in downstream task. Knowledge graph completion or link predictions is thus important approaches for populating existing KGs.

Knowledge graph embedding models for KG completion have attracted much attention, due to their outstanding performance on various applications that require machines to recognize and understand queries and their semantics. These em-

bedding model is to represent entitles and relations in a KG into a low dimensional continuous vector space, such vectors contain rich semantic information, and can benefit many downstream tasks especially knowledge graph completion or linked predictions. Whether two entities have a previously unknown relationship can be predicted by simple functions of their corresponding vectors.

Despite the success of previous approaches in KG embedding [Bordes *et al.*, 2013] [Wang *et al.*, 2014b], most of them mainly model triples individually, ignore lots of information implicitly provided by the structure of the KG. In fact, triples are connected to each other and many triples around a triple could be regarded as a description of it. Recently, Several authors have addressed this issue by incorporating relation path information into model learning [Lin *et al.*, 2015a] [Toutanova K, 2016] learning and have shown that the relation paths between entities in KGs provide useful information and improve KG completion. These approaches only consider relation information while miss more structure information, such as K-degree neighbors of a given entity, a connected subgraph which n could be exploited for better KB completion. For instance, the whole neighborhood of entities and a connected subgraph between two entities could provide lots of useful information for predicting the relationship between two entities.

In this paper, we present a novel approach to embed a knowledge graph by utilizing the structure information called Attentional-Triple-Context-based knowledge Embedding model(ATCE) which utilizes and chooses the proper context of each triple in the knowledge graph. We define triple context consisting of neighbor context and path context, and define a new score function to evaluate the correlation between a triple and its contexts. Instead of using each triple independently, we incorporate triple context into the score function which is used to evaluate the confidence of a triple. In this way, we make use of a triple context while learning embeddings.

The advantages of our approach are three-fold: 1) We embed a triple by utilizing a local subgraph around a triple instead of a set of independent triples, and extract two kinds of context.

2) Based on the local structure information, we proposed a novel embedding learning approach which named ATCE and a new loss function which convert the score function in

TransE to a conditional probability.

3) In order to overcome the noisy data in the triple context of a triple, an attention mechanism in our approach are proposed to choose the proper information for embedding. In the meanwhile, the attention mechanism can learn the representation power of different neighbor entities and connective path in its context.

Finally, We have conducted preliminary experiment on two benchmark data sets and assessed our method on link prediction task and triple classification. In the experiments we shows chosen context through the attention mechanism to improve the effectiveness of this mechanism. The experimental results show impressive improvements on predictive accuracy compared to other baselines.

## 2 Triple Context

Firstly, we introduce some notations that are used in this paper. Let  $\mathcal{K}$  be a knowledge graph,  $\mathcal{E}$  and  $\mathcal{R}$  the set of all entities and relations respectively in  $\mathcal{K}$ . Each triple is denoted as  $(h, r, t)$ , in which  $h$  is the head entity,  $t$  is the tail entity and  $r$  is the relation between  $h$  and  $t$ . The embeddings of each entity and relation are denoted in bold, e.g.,  $\mathbf{h}$  is the embedding of  $h$ . All the embeddings are in  $d$ -dimensional space  $\mathbb{R}^d$ . Our goal is to learn embeddings of all entities and relations, which is denoted as  $\Theta$ . In the following subsections, we define neighbor context and path context, and then give the framework of our model.

### 2.1 Neighbor Context

Neighbor context of an entity is the surroundings of it in KG. It is the local structure that interacts most with the entity and can reflect various aspects of the entity. Specifically, given an entity  $e$ , the neighbor context of  $e$  is a set  $C_N(e) = \{(r, t) | \forall r, t, (e, r, t) \in \mathcal{K}\}$ , where  $r$  is an outgoing edge (relation) from  $e$  and  $t$  is the entity it reaches through  $r$ . In other words, the neighbor context of  $e$  is all the *relation-tail* pairs appearing in triples with  $e$  as the head. For example, as shown in Figure 1, the neighbor context of entity  $h$  is  $C_N(h) = \{(r_4, e_1), (r_3, e_2), (r_2, e_3), (r_1, e_8), (r_1, e_{10})\}$ . We predict the appearance of an entity based on its neighbor context in our model, as a measurement of the compatibility of the entity and its neighbor context.

### 2.2 Path Context

Path context of a pair of entities is the set of paths that starts from an entity to the other in a KG. It is helpful in modeling the relation and capturing interactions between the pair of entities. Given a pair of entities  $(h, t)$ , the path context of  $(h, t)$  is a set  $C_P(h, t) = \{p_i | \forall r_{m_1}, \dots, r_{m_i}, e_1, \dots, e_{m_i-1}, p_i = (r_{m_1}, \dots, r_{m_i}, (h, r_{m_1}, e_1) \in \mathcal{K}, \dots, (e_{m_i-1}, r_{m_i}, t) \in \mathcal{K})\}$ , where  $p_i$  is a list of relations (labeled edges) through which it can traverse from  $h$  to  $t$ ,  $m_i$  is the length of path  $p_i$ . In Figure 1, the path context between  $h$  and  $t$  is  $C_P(h, t) = \{(r_1, r_2), (r_2, r_1, r_2)\}$ . We use the path context to predict the tail entity of a triple given the head entity.

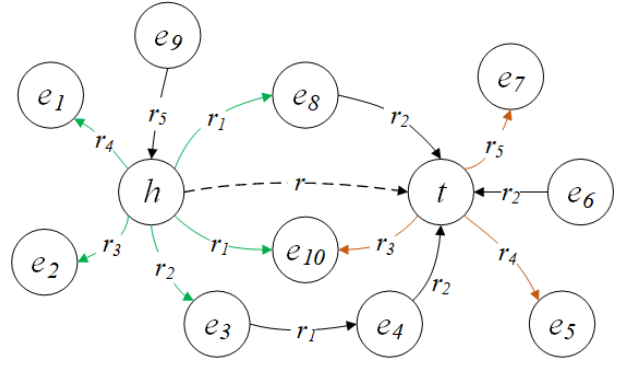


Figure 1: An illustration of the *triple context* of a triple  $(h, r, t)$  in a knowledge graph.

## 3 Knowledge Graph Embedding With Attentional Triple Context

So far, we have introduced neighbor context and path context, based on which we can define triple context. The triple context of triple  $(h, r, t)$  is composed of the neighbor context of the head entity  $h$ , the path context of the entity pair  $(h, t)$ , which can be formalized as:

$$C(h, r, t) = C_N(h) \cup C_P(h, t) \quad (1)$$

The triple context of a triple can be considered to embody the surrounding structures of it in the graph, which makes the model aware of the information contained in graph structures.

We then introduce our approach in detail. In general KG embedding models, the score function of a triple is only related to the embeddings of entities and relations. For example, TransE defines the score function as  $f_{TransE}(h, r, t) = \|\mathbf{h} + \mathbf{r} - \mathbf{t}\|_{L_1/L_2}$ . In our method, triple context is introduced in the score function. Given a candidate triple  $(h, r, t)$ , the score function is the conditional probability that the triple holds given the triple context and all the embeddings, as follows:

$$f(h, r, t) = P((h, r, t) | C(h, r, t); \Theta) \quad (2)$$

where  $C(h, r, t)$  is the triple context of  $(h, r, t)$ . A higher score of a triple indicates that it holds to a greater extent.

We define an objective function by maximizing the joint probability of all triples in knowledge graph  $\mathcal{K}$ , which can be formulated as:

$$P(\mathcal{K} | \Theta) = \prod_{(h, r, t) \in \mathcal{K}} f(h, r, t) \quad (3)$$

For the score function in Eq. (2), we use conditional probability formula to decompose the probability  $P((h, r, t) | C(h, r, t); \Theta)$  as:

$$f(h, r, t) = P(h | C(h, r, t); \Theta) \cdot P(t, r | C(h, r, t), h; \Theta) \quad (4)$$

where the evaluation of the whole triple is decomposed into two parts. The probabilities that  $h$ ,  $t$  and  $r$  appear given respective condition are determined in turn in these two parts.

The first part  $P(h | C(h, r, t); \Theta)$  in Eq. (4) represents the conditional probability that  $h$  is the head entity given the

triple context and all the embeddings. Since whether  $h$  appears is decided mostly by the neighboring structures of  $h$  in the KG, we can approximate  $P(h|C(h, r, t); \Theta)$  as  $P(h|C_N(h); \Theta)$ , where  $C_N(h)$  is the neighbor context of  $h$  in the KG. The approximated probability  $P(h|C_N(h); \Theta)$  can be considered as the compatibility between  $h$  and its neighbor context, it is formalized as a softmax-like representation, which is also used in [Wang *et al.*, 2014a] and has been validated, as follows:

$$P(h|C_N(h); \Theta) = \frac{\exp(g_N(h, C_N(h)))}{\sum_{h' \in \mathcal{E}} \exp(g_N(h', C_N(h)))} \quad (5)$$

where  $g_P(\cdot, \cdot)$  is the function that describes the correlation between an arbitrary entity  $h'$  and a neighbor context of the specific entity  $h$ . In reality, different neighbors may have different power of influence to represent  $h$ . For example, a target entity is *Terminate2 : JudgementDay* which is a famous movie in 1991. The director of this movie or the type of this movie may more important than others. In order to filter unrelated entities and obtain  $g_P(\cdot, \cdot)$ , we first obtain  $a_i$  which means the represent power between the  $i$ th neighbor entity in  $C_N(h)$  and an arbitrary triple  $\langle h', r, t \rangle$ . Inspired by score function of TransE, we substitute the neighbor entities in  $C_N(h)$  for the head  $h'$  in triple to compute  $a_i$  and then we define  $a_i$

$$a_i = \|t_{n_i} - r_{n_i} + r - t\| \quad (6)$$

Where  $t_{n_i}$  is the  $i$ th neighbor entity in  $C_N(h)$  and  $r_{n_i}$  is the relation between  $t_{n_i}$  and  $h'$ .

Then we use attention model  $\alpha_i$  to represent how  $h'$  selectively focuses on  $C_N(h)$ . If the value of  $\alpha_i$  is higher, the corresponding neighbor entity in  $C_N(h)$  is more important for  $h'$ .

$$\alpha_i = \frac{\exp(-a_i)}{\sum_j \exp(-a_j)} \quad (7)$$

In Eq. (8), given a neighbor context  $C_N(h)$ , the embedding vector of each neighbor has different weights by the attention mechanism in Eq. (7). Finally, based on the attention results we obtain the correlation between an arbitrary entity  $h'$  and neighbor context  $C_N(h)$ :

$$g_N(h, C_N(h)) = - \sum_i \alpha_i \|\mathbf{h}' + \mathbf{r}_{n_i} - \mathbf{t}_{n_i}\| \quad (8)$$

The second part  $P(r, t|C(h, r, t), h; \Theta)$  in Eq. (4) is the conditional probability that  $t$  is the tail entity and  $r$  is the relation given the head entity  $h$ , triple context and all the embeddings. In this part we introduce path context that means  $t$  could be related to  $h$  through a potential connective path in a knowledge graph. In the second part two kinds of relatedness should be considered that one is the relatedness between  $h$  and  $t$  in a potential connective path  $p_i$ . And the other is the relatedness between  $r$  and  $p_i$ . Then We introduce path context among  $h$ ,  $t$  and  $r$  to measure the relatedness of them and approximate  $P(r, t|C(h, r, t), h; \Theta)$  as  $P(r, t|C_P(h, t), h; \Theta)$ , where  $C_P(h, t)$  is the path context between  $h$  and  $t$ . The approximated probability  $P(r, t|C_P(h, t), h; \Theta)$  is formalized as follows:

$$P(r, t|C_P(h, t), h; \Theta) = \frac{\exp(g_P(r, t, C_P(h, t)))}{\sum_{r' \in \mathcal{R}, t' \in \mathcal{E}} \exp(g_P(r', t', C_P(h, t)))} \quad (9)$$

where  $g_P(\cdot, \cdot)$  is a function of correlation among an arbitrary entity  $t'$ , an arbitrary relation  $r'$  and path context of the specific entity pair  $(h, t)$ . Similar to the neighbor context, given a triple  $\langle h, r, t \rangle$  different paths in a path context  $C_P(h, t)$  have different power of influence the triple. For example when predicting the entity *Englis*, relations like *locate\_in\_Country* will have less attentions, and the relation *Speak\_Language* will have greater attention to represent *Englis*. In order to obtain  $g_P(\cdot, \cdot)$ , we firstly calculate the correlation  $b_i$  between path  $p_i$  in  $C_P(h, t)$  and a tail entity:

$$b_i = \|\mathbf{h} + \mathbf{p}_i - \mathbf{t}\| \quad (10)$$

Where  $\mathbf{p}_i$  composes all relations in  $p_i$  into a single vector by summing over all their embeddings and this approach is also used in PTransE [Lin *et al.*, 2015a]. For example, for path  $p_i = (r_{m_1}, \dots, r_{m_i})$ , the embedding of it is  $\mathbf{p}_i = \mathbf{r}_{m_1} + \dots + \mathbf{r}_{m_i}$ . Eq. (10) has a similar meaning with Eq. (6).

We then choose important paths from  $C_P(h, t)$  through the attention mechanisms. For each  $p_i$ , the weight  $\beta_i$  is then defined as

$$\beta_i = \frac{\exp(-b_i)}{\sum_j \exp(-b_j)} \quad (11)$$

If  $p_i$  is more closer to  $h$  and  $t$ ,  $\beta_i$  will indicate  $p_i$  have greater attentions on translating from  $h$  to  $t$ .

Finally given an arbitrary entity  $t'$  and an arbitrary relation  $r'$ , the correlation  $g_P(\cdot, \cdot)$  between  $t', r'$  and  $C_P(h, t)$  is:

$$g_P(r', t', C_P(h, t)) = - \sum_i \beta_i (\|\mathbf{h} + \mathbf{p}_i - \mathbf{t}'\| + \|\mathbf{p}_i - \mathbf{r}'\|) \quad (12)$$

In Eq. (12), given a path context  $C_P(h, t)$ , the embedding vector of each path  $p_i$  has different weight  $\beta_i$ . Eq. (12) has two parts, the first part is  $\|\mathbf{h} + \mathbf{p}_i - \mathbf{t}'\|$  that indicates the correlation between  $t'$  and  $p_i$ . Similarly, the second part is  $\|\mathbf{p}_i - \mathbf{r}'\|$  is the correlation between  $r'$  and  $p_i$ . If  $p_i$  is more related to  $r'$  and  $t'$ , the probability of Eq. (9) will be greater.

To utilize these two kinds of context, we combine them by jointly maximizing the probability in Eq. (4) of a triple which is exist in a knowledge graph.  $P(h|C(h, r, t); \Theta)$  and  $P(r, t|C(h, r, t), h; \Theta)$  can be approximated as  $P(h|C_N(h); \Theta)$  and  $P(r, t|C_P(h, t), h; \Theta)$ , respectively. Thus Eq. (4) Thus, Eq. (4) can be approximated as:

$$f(h, r, t) \approx P(h|C_N(h); \Theta) \cdot P(r, t|C_P(h, t), h; \Theta) \quad (13)$$

in which way the neighbor context and the path context of a triple are incorporated.

### 3.1 Model Learning

By feasible approximation, the score function is transformed to Eq. (13), each part is represented in softmax form as Eq. (5) and Eq. (9). However, it is impractical to compute

these softmax functions directly because of high computational overhead. Hence, we adopt negative sampling, which is proposed in [Mikolov *et al.*, 2013] to approximate full softmax function efficiently, to approximate softmax functions in our model. For the whole knowledge graph  $\mathcal{K}$ , it is approximated via negative sampling as follows:

$$\begin{aligned}
P(\mathcal{K}|\Theta) &= \prod_{(h,r,t) \in \mathcal{K}} f(h,r,t) \\
&= \prod_{(h,r,t) \in \mathcal{K}} P(h|C(h,r,t);\Theta) \cdot P(r,t|C(h,r,t),h;\Theta) \\
&\approx \prod_{(h,r,t) \in \mathcal{K}} [\sigma(g_N(h,C_N(h))) \cdot \prod_{h'} \sigma(-g_N(h',C_N(h)))] \\
&\quad \cdot [\sigma(g_P(r,t,C_P(h,t))) \cdot \prod_{r',t'} \sigma(-g_P(r',t',C_P(h,t)))]
\end{aligned} \tag{14}$$

where  $\{h',r,t\}$  is the corrupted triples by replacing the head entity with an arbitrary entity.  $\{h,r',t'\}$  also is the corrupted triples by replacing the tail entity with an arbitrary entity and a relation  $r$  with an arbitrary relation.  $\sigma(\cdot)$  is the logistic function. Simultaneity, for convenient application, we convert Eq. (14) to the negative logarithm form which can be optimized by stochastic gradient descent (SGD). The objective function of a knowledge graph is formulated as:

$$\begin{aligned}
-\log P(\mathcal{K}|\Theta) &= - \sum_{(h,r,t) \in \mathcal{K}} [\log \sigma(g_N(h,C_N(h)))] \\
&\quad + \sum_{h'} \log \sigma(-g_N(h',C_N(h))) \\
&\quad + \log \sigma(g_P(r,t,C_P(h,t))) \\
&\quad + \sum_{r',t'} \sigma(-g_P(r',t',C_P(h,t)))
\end{aligned} \tag{15}$$

In real data sets, the size of neighbor context and path context may be very large, which is computationally expensive for model learning. For this reason, we sample from neighbor context and path context to make triple context tractable. Specifically, we set a threshold  $n_N$  for neighbor context and  $n_P$  for path context; if the size of the original context exceeds the threshold, we sample a subset, size of which is the threshold, for model learning. Moreover, the length of relation path is constrained to 2 and 3 in our model.

## 4 Experiments

### 4.1 Experimental Setup

**Data Set.** We use two widely-used benchmark data sets FB15k [Bordes *et al.*, 2013] and FB15k-237 for evaluation, which are extracted from Freebase. FB15k has 592,213 triples with 14,951 entities and 1,345 relationships. Triples in FB15k-237 are a subset of the FB15K set. FB15k-237 excludes redundant relations and direct training links for held-out triples, with the goal of making the task more realistic [?]. The two datasets are further divided into three parts for model training, tuning and evaluation. Specifically, we use

Table 3: Link prediction results

Metric	Mean Rank		HITS@10(%)	
	Raw	Filter	Raw	Filter
TransE	243	125	34.9	47.1
TransH (unif)	211	84	42.5	58.5
TransH (bern)	212	87	45.7	64.4
TransR (unif)	226	78	43.8	65.5
TransR (bern)	198	77	48.2	68.7
CTransR (unif)	233	82	44.0	66.3
CTransR (bern)	199	75	48.4	70.2
PTransE	207	58	51.4	<b>84.6</b>
GAKE	228	119	44.5	64.8
ATCE	<b>110</b>	<b>25</b>	<b>55.3</b>	83.1

FB15k and FB15k-237 since their triples are rich and closer to the real popular knowledge graph.

**Evaluation protocol.** Following the same protocol used in [Bordes *et al.*, 2013], we use *Mean Rank* and *Hits@10* as evaluation protocols of our model. For each test triple  $(h,r,t)$ , we replace tail head  $t$  (or the head  $h$ ) with each entity  $e$  in  $\mathcal{E}$  to generate *corrupted triples* and calculate the scores of each triple using the score function. After ranking the scores in descending order, we then get the rank of the correct entity. *Mean Rank* is the mean of all the predicted ranks, and *Hits@10* denotes the proportion of correct entities ranked in the top 10. Note that, a corrupted triple ranking above a test triple could be valid, which should not be counted as an error. To eliminate the effects of such condition, corrupted triples that already exist in the KG are filtered before ranking. In this case, the setting of evaluation is called "Filter", while the original one is called "Raw". Since this effect, the "Filter" setting is more preferred. In both settings, a higher Hits@10 imply the better performance of a model.

**Baselines.** We use a few outstanding models in recent years as baselines and compare our model with them, including TransE [Bordes *et al.*, 2013], TransH [Wang *et al.*, 2014b], TransR [Lin *et al.*, 2015b], CTransR [Lin *et al.*, 2015b], PTransE [Lin *et al.*, 2015a] and GAKE [Feng *et al.*, 2016]. For each baseline model, we first learn representations of all entities and relations. After that the conditional probability of a triple is calculated by the score function. While in ATCE, we construct the triple context with each triple. At last, if the conditional probability of  $\langle h,r,t \rangle$  is larger than  $p_r$ ,  $\langle h,r,t \rangle$  is regard as a correct triple.

**Implementation.** We construct the knowledge graph using Apache TinkerPop<sup>1</sup>, an open source graph computing framework. In a few cases, the reverse relation, an edge labeled  $r^{-1}$  from  $t$  to  $h$  for the triple  $(h,r,t)$ , would be useful when representing some patterns in the graph. For instance, the relation path  $a \xrightarrow{motherOf} b \xleftarrow{fatherOf} c$ , i.e.,  $(a, motherOf, b)$  and  $(c, fatherOf, b)$ , indicates a potential relation *marriedTo* between  $a$  and  $c$ . Therefore, we add reverse relation of each relation into KG. Specifically, for each edge labeled  $r$  from  $h$  to  $t$  in the graph, we add another edge labeled  $r^{-1}$  from  $t$  to  $h$ .

<sup>1</sup><http://tinkerpop.apache.org/>

Table 1: Results on FB15k by relation category

Task	Predicting head(HITS@10(%))				Predicting tail(HITS@10(%))			
Relation Category	1-To-1	1-To-N	N-To-1	N-To-N	1-To-1	1-To-N	N-To-1	N-To-N
TransE	43.7	65.7	18.2	47.2	43.7	19.7	66.7	50.0
TransH (unif)	66.7	81.7	30.2	57.4	63.7	30.1	83.2	60.8
TransH (bern)	66.8	87.6	28.7	64.5	65.5	39.8	83.3	67.2
TransR (unif)	76.9	77.9	38.1	66.9	76.2	38.4	76.2	69.1
TransR (bern)	78.8	<b>89.2</b>	34.1	69.2	79.2	37.4	<b>90.4</b>	72.1
CTransR (unif)	78.6	77.8	36.4	68.0	77.4	37.8	78.0	70.3
CTransR (bern)	<b>81.5</b>	89.0	34.7	71.2	<b>80.8</b>	38.6	90.1	73.8
ATCE	71.0	60.3	<b>83.9</b>	<b>81.9</b>	70.3	<b>89.9</b>	76.0	<b>89.2</b>

Table 2: Results on FB15k-237 by relation category

Task	Predicting head(HITS@10(%))				Predicting tail(HITS@10(%))			
Relation Category	1-To-1	1-To-N	N-To-1	N-To-N	1-To-1	1-To-N	N-To-1	N-To-N
TransE	43.7	65.7	18.2	47.2	43.7	19.7	66.7	50.0
TransH (unif)	66.7	81.7	30.2	57.4	63.7	30.1	83.2	60.8
TransH (bern)	66.8	87.6	28.7	64.5	65.5	39.8	83.3	67.2
TransR (unif)	76.9	77.9	38.1	66.9	76.2	38.4	76.2	69.1
TransR (bern)	78.8	<b>89.2</b>	34.1	69.2	79.2	37.4	<b>90.4</b>	72.1
CTransR (unif)	78.6	77.8	36.4	68.0	77.4	37.8	78.0	70.3
CTransR (bern)	<b>81.5</b>	89.0	34.7	71.2	<b>80.8</b>	38.6	90.1	73.8
ATCE	71.0	60.3	<b>83.9</b>	<b>81.9</b>	70.3	<b>89.9</b>	76.0	<b>89.2</b>

For neighbor context generation, it’s expensive to consider all the neighbors of each entity in the graph for the reason that there are some entities connecting with a large amount of other entities, which would lead to a huge size of neighbor context. We use sampling to reduce the size of neighbor context. For those entity whose neighbor context size is larger than a fixed size  $n$ , we sample  $n$  neighbors randomly from it’s neighbor context. Similarly, a large number of paths between a pair of entities would result in high computational complexity. To solve the problem, firstly, we limit the length of paths by 2-step and 3-step, then, we use random walk to sample  $m$  paths between a pair of entities. In our experiment,  $n$  and  $m$  are all set as 10. Note that for some pairs of entities, there may be no 2 or 3 step relation paths. In such case, we suppose that the relatedness between those pairs of entities are relatively low and the values of  $g_P(\cdot, \cdot)$  in Eq. (12) are set as -100. We store all triple contexts in MongoDB<sup>2</sup>, a NoSQL database, that can help us for rapid seeking the context of a specified entity.

For all tasks, all parameters were initialized from a uniform distribution  $U[-\sqrt{6k}, \sqrt{6k}]$  as suggested by TransE. ATCE can also be initialized with pre-trained embeddings. We use mini-batch SGD to train our model. We choose the learning rate  $\alpha$  of SGD among  $\{0.1, 0.01, 0.001\}$ , the dimension of embeddings  $k$  among  $\{50, 75, 100\}$ , the batch size  $B$  among  $\{120, 480, 960, 1920, 4800\}$ . The best parameters are determined by the performance on valid set. The optimal parameters are  $\alpha = 0.001$ ,  $k = 50$  and  $B = 4800$ . ATCE is implemented in Python using Pytorch. The code and data are available at <https://github.com/>

<sup>2</sup>[www.mongodb.org/](http://www.mongodb.org/)

gaohuan2015/IJCAI-Model

## 4.2 Link Prediction

Link prediction [Bordes *et al.*, 2013] is a task to predict the missing head or tail entity in a given triple based on training triples. Metrics *Mean Rank* and *Hits@10* are used to measure the performance of our model.

We collected the result of link prediction in Table 3. From the results we can see that, our model outperforms other baselines on most of the metrics significantly and consistently, while slightly worse than PTransE on HITS@10. The result implies that triple contexts do improve the performance on link prediction. Although using similar types of contexts in the graph, GAKE’s performance is inferior to our model, which shows the superiority of our framework. Note that the experimental results of HOLE are absent here, for it uses a different metric, MRR (Mean Reciprocal rank), instead of *Mean rank* for evaluation. But according to *Hits@10* reported in [Nickel *et al.*, 2016], the results of our model are better than HOLE.

In Table 1 and Table 2, we show separate evaluation results by category of relationships on the two benchmark datasets. We can see that ATCE brings promising improvements on modeling complex relations, such as predicting tail of 1-To-N relations, predicting head of N-To-1 relations and N-To-N relations. Specifically, ATCE behaves well when predicting the “N” side of 1-To-N and N-To-1 relations, indicating that valid triples have higher scores than invalid triples in general. In some other simpler scenarios, such as 1-To-1 relations and predicting the “1” side of 1-To-N and N-To-1 relations, the performance of ATCE is still acceptable although not so good as some other baselines, such as TransH and TransR. The re-

Table 4: Triple Classification results

Start	FB15k	FB15k-237
TransE	0.81	0.82
TransH	0.81	0.81
TransR	0.82	0.82
TransD	0.87	0.87
GAKE	0.89	0.88
NTN	0.41	0.42
ATCE	<b>0.91</b>	<b>0.92</b>

sults suggest that the incorporation of triple context is helpful when handling complex relations, at the cost of precision in modeling simple relations, which seems complementary to some other baselines.

### 4.3 Triple Classification

We also test our model on triple classification. In this task, given a knowledge base and a triple  $\langle h, r, t \rangle$  we aim to determine whether it is correct or not. FB15K has only positive examples, thus we generated negative triples for FB15K by following strategy of [Socher *et al.*, 2013]. As the result, the classification accuracies on FB15K and FB15k-237 can be compare directly with previous studies. In this task, we choose TransE, TransH, TransR, TransD, NTN and GAKE as baseline models. The parameter values for training TransE, TransH, TransR, TransD, NTN and GAKE are borrowed from their reports.

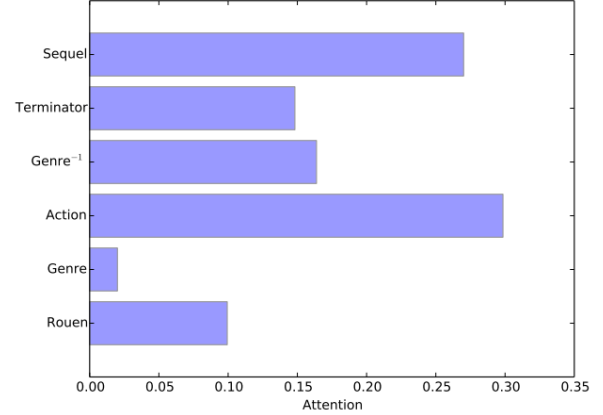
In Table 4, we shows the accuracies of triple classification on two datasets. The context preserving embeddings in general outperform their base model. As the table shows, ATCE always has higher accuracy than TransE, TransH, TransR, TransD, NTN and GAKE. One thing to note is that the improvements by the context preserving embeddings are always observed in FB15k, while those in FB15-237 are small or slightly negative. This can be explained with number of triples with a transitive or symmetric relation in the dataset. In FB15k-237, this kinds of triples are removed from FB15k. Thus the improvement in FB15K is remarkable, our approach outperforms others by 11.04% in terms of accuracy on average as the triple context brings more information especially relations between entities when learning the knowledge graph representations.

Furthermore, to better understand the attention mechanism in our approach, The examples of path are selected by attention mechanism in path context are shown in Table 5. The number in front each line is the weight score, which is computed by the attention mechanism for each relation. From the examples we can see that ATCE successfully combines structure learning and parameter learning. It not only choose multiple connective path between two entities to capture the complex structure in the knowledge base, but also learn weight score of the path for a specific relation.

For neighbor context, we also use attention mechanism to choose which neighbor is more related for the head. we demonstrate attentions of the 6 different neighbors when they are regards as the neighbor context of the entity *Terminate2 : JudgementDay*, which indicates a movie in

Table 5: Attention Results in Path Context

Weight	Path	Relation
0.45	<i>contains</i> $\rightarrow$ <i>contains</i>	<i>partially_contains</i>
0.35	<i>contains</i> $\rightarrow$ <i>nationality</i>	<i>marriage_location</i>
0.24	<i>nationality</i>	<i>marriage_location</i>
0.35	<i>contains</i> $\rightarrow$ <i>place_lived</i>	<i>marriage_location</i>
0.2	<i>nominated_for</i>	<i>film_edited_by</i>
0.3	<i>nominated</i> $\rightarrow$ <i>award_nominee</i>	<i>film_edited_by</i>

Figure 2: Attentions for neighbor context of *Terminate2:JudgementDay*

1991. Figure 2 shows the results, from the results we see two entities, *Action* and *Sequel*, have largest attention to represent the target entity *Terminate2 : JudgementDay*, as *Action* reflects the type of movie while only some of movies have sequels.

## 5 Related Work

In recent year, a new kind of approach named *knowledge graph embedding* (also known as *knowledge representation learning*) has been studied in [Nickel *et al.*, 2011][Socher *et al.*, 2013][Bordes *et al.*, 2013][Wang *et al.*, 2014b][Lin *et al.*, 2015b]. It makes KGs more manageable and performs well on several important tasks like link prediction [Bordes *et al.*, 2011] and triple classification [Socher *et al.*, 2013]. It aims to embed a KG into a low dimensional vector space, where each entity and relation are represented as a vector (also referred to as *embedding*) or a matrix. In general, a KG embedding model defines a score function for evaluating the confidence of a triple, and optimizes an objective function constructed from the score function to learn the embeddings. In this way, symbolic computation can be replaced by numerical computation, which makes it easier to manipulate KGs through linear algebra operations while the semantic information is retained meanwhile. A few typical models of knowledge graph embedding are as follows. TransE is one of the most widely used approach in those approaches, which views a head entity can be translation to a tail entity by a relationship on the same hyperplane. Based on the score function of the translation model which is designed as  $f_{TransE}(h, r, t) = \|\mathbf{h} + \mathbf{r} - \mathbf{t}\|_{L_1/L_2}$  variant models of

TransE [Bordes *et al.*, 2011] are introduced to improve the performance. In TransH [Wang *et al.*, 2014b], entities are projected into the relation specific hyperplane that is perpendicular to the relationship embedding. TransD [Guoliang Ji and Zhao, 2015] and TransR [Lin *et al.*, 2015b] still follow the principle of TransH, the entities and relations are embedded on different hyperplane by the different mapping matrices which are both related to the entity and relation.

Translation-based approaches treat knowledge graph as a set of triples. Unlike aforementioned models, TransA [Bordes *et al.*, 2011] utilizes a closed set of candidate entities to determine the relation. The similar works include RESCAL, SME and LFM. In these models, relations are represented as a bilinear operator that can interactions across the entities of a triples to explain its validity. In addition, some works discover incorporating additional information will improve the performance, such as text [Jiacheng Xu and Huang, 2017] and entity type [Denis KrompaB and Tresp, 2015].

There is also works that incorporates graph structures into KG embedding models. PTransE [Lin *et al.*, 2015a] is a path-based representation learning model, which utilizes relation paths to improve the performance of TransE. However, it merely uses relation paths in the KG, ignoring other information contained in graph structures. Besides, GAKE [Feng *et al.*, 2016], the most relevant model for our work, leverages three types of graph context for representation learning, in a way similar to learning language models. However, when dealing with contexts, it average the embeddings of all entities and relations in a context as the representation of it, in which way the structural information is still lost. Furthermore, its experimental results are inferior to TransH and TransR. It is a huge challenge in context selection and the feasible use of it.

## 6 Conclusion

In this paper, we proposed ATCE, a new KG embedding model which is able to take advantages of the triple context which includes neighbor context and path context in the KG. By defining two kinds of context of a triple and representing them in a unified framework, our model can learn embeddings that are aware of their context. ATCE not only can learn the embedding of a given KG, but also can capture the complex structure in KG. We evaluate our model on two benchmark datasets on link prediction and triple classification. In the meantime, we analysis the attention mechanism for choosing the neighbor enmities and connective pathes in triple context. The experimental results show significant improvements over the major baselines.

In the future, we will research on the following aspects: (1) it will be interesting to incorporate multilingual knowledge graph in to our model to further improve the performance. (2) Current results show complementarity to some other methods such as TransH, TransR. We would think about a combination of those methods and our model.

## References

- [Auer *et al.*, 2007] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. Dbpedia: A nucleus for a web of open data. In *Proceedings of ISWC 2007 + ASWC 2007*, pages 722–735, 2007.
- [Bordes *et al.*, 2011] Antoine Bordes, Jason Weston, Ronan Collobert, and Yoshua Bengio. Learning structured embeddings of knowledge bases. In *Proc. of the AAAI*, 2011.
- [Bordes *et al.*, 2013] Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In *NIPS*, pages 2787–2795, 2013.
- [Denis KrompaB and Tresp, 2015] Stephan Baier Denis KrompaB and Volker Tresp. Type-constrained representation learning in knowledge graphs. In *Proc. of the ISWC*, 2015.
- [Feng *et al.*, 2016] Jun Feng, Minlie Huang, Yang Yang, and Xiaoyan Zhu. Gake: Graph aware knowledge embedding. In *COLING*, pages 641–651, 2016.
- [Guoliang Ji and Zhao, 2015] Liheng Xu Kang Liu Guoliang Ji, Shizhu He and Jun Zhao. Knowledge graph embedding via dynamic mapping matrix. In *Proc. of ACL*, pages 687–696, 2015.
- [Jiacheng Xu and Huang, 2017] Kan Chen Jiacheng Xu, Xipeng Qiu and Xuanjing Huang. Knowledge graph representation with jointly structural and textual encoding. In *Proc. of IJCAI*, pages 1318–1324, 2017.
- [Lin *et al.*, 2015a] Yankai Lin, Zhiyuan Liu, Huan-Bo Luan, Maosong Sun, Siwei Rao, and Song Liu. Modeling relation paths for representation learning of knowledge bases. In *Proc. of EMNLP*, pages 705–714, 2015.
- [Lin *et al.*, 2015b] Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. Learning entity and relation embeddings for knowledge graph completion. In *Proc. of AAAI*, pages 2181–2187, 2015.
- [Mikolov *et al.*, 2013] Tomas Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, pages 3111–3119, 2013.
- [Mitchell *et al.*, 2015] Tom M. Mitchell, William W. Cohen, Estevam R. Hruschka Jr., and et al. Never-ending learning. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA.*, pages 2302–2310, 2015.
- [Nickel *et al.*, 2011] Maximilian Nickel, Volker Tresp, and Hans-Peter Kriegel. A three-way model for collective learning on multi-relational data. In *Proc. of ICML*, pages 809–816, 2011.
- [Nickel *et al.*, 2016] Maximilian Nickel, Lorenzo Rosasco, and Tomaso A. Poggio. Holographic embeddings of knowledge graphs. In *Proc. of AAAI*, pages 1955–1961, 2016.
- [Socher *et al.*, 2013] Richard Socher, Danqi Chen, Christopher D. Manning, and Andrew Y. Ng. Reasoning with neural tensor networks for knowledge base completion. In *NIPS*, pages 926–934, 2013.

- [Suchanek *et al.*, 2007] Fabian M. Suchanek, Gjergji Kasneci, and Gerhard Weikum. Yago: A core of semantic knowledge. In *Proceedings of WWW 2007*, pages 697–706, 2007.
- [Toutanova K, 2016] Yih W et al. Toutanova K, Lin V. Compositional learning of embeddings for relation paths in knowledge base and text. In *Proc. of ACL*, pages 1120–1136, 2016.
- [Wang *et al.*, 2014a] Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. Knowledge graph and text jointly embedding. In *Proc. of EMNLP*, pages 1591–1601, 2014.
- [Wang *et al.*, 2014b] Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. Knowledge graph embedding by translating on hyperplanes. In *Proc. of AAAI*, pages 1112–1119, 2014.