



开放科学  
(资源服务)  
标识码  
(OSID)

# 基于混合上下文的知识表示学习方法研究

张良 石璐

东南大学计算机科学与工程学院 南京 211189

**摘要:** [目的/意义] 知识图谱是一种包含了丰富实体和关系的数据结构, 然而当前绝大多数知识图谱都是不完备的。[方法/过程] 知识表示学习是目前对知识图谱进行补全的一种热门方法, 本文在对知识表示学习方法梳理的基础上探究了文本信息和知识图谱自身三元组信息互补的机制, 并利用远程监督与注意力机制将两类信息结合。[结果/结论] 提出一种基于混合上下文的知识表示学习方法, 实验结果表明, 结合文本与三元组这类混合上下文信息的模型能够明显提升知识图谱表示学习的效果, 并在一些指标上优于已有的一些模型。

**关键词:** 知识图谱; 表示学习; 文本; 远程监督; 注意力机制

**中图分类号:** G35; TP391

## The Research of Knowledge Representation Learning Based on Hybrid Context

ZHANG Liang SHI Jun

School of Computer Science and Engineering, Southeast University, Nanjing 211189, China

**Abstract:** [Objective/Significance] Knowledge graph is a kind of data structure which contains rich entities and relationships. However, most knowledge graph is incomplete. [Method / Process] Knowledge representation learning is a popular method to complete knowledge graph. Based on the analysis of knowledge representation learning methods, this paper explores the complementary mechanism of text information and triple information, and combines the two types of information with distance supervision and attention mechanism. [Results / Conclusions] Finally, a knowledge representation learning method based on mixed context is proposed. The experimental results show that the model combining text and triples can significantly improve the effect of knowledge representation learning, and is better than some existing models in some indicators.

**Keywords:** Knowledge graph; representation learning; text; distance supervision; attention mechanism

**基金项目** 国家自然科学基金重点项目“面向大规模多源数据的人物画像和定位分析关键技术”(U509000121)。

**作者简介** 张良(1989-), 博士研究生, 研究方向为知识图谱、语义 Web, E-mail: lzhang@seu.edu.cn; 石璐(1993-), 硕士, 研究方向为知识图谱。

## 引言

近年来,知识图谱越来越受到人工智能领域的广泛关注,各类知识图谱在反欺诈、智能问答、人物分析,旅游、医疗等垂直领域中发挥了重要作用<sup>[1-5]</sup>。知识图谱旨在刻画客观世界中的实体及其之间的关系,其中实体表示成知识图谱中的节点,实体之间的关系对应着知识图谱中各实体之间的连线,可以看成一种有向图形式的数据结构。知识图谱不仅能清晰地展现各实体间原有的关系,更能通过推理计算等方法来获得数据之间潜在的联系<sup>[6]</sup>。由于知识图谱特别是大规模知识图谱包含着非常复杂的结构化信息,以往在处理的时候需要专门设计相应的图算法来实现知识图谱的计算和存储,这样的方法复杂度高,可移植性差,效率上往往不能得到保证<sup>[7]</sup>。

随着深度学习技术的发展,知识图谱表示学习方法给这一问题的解决带来了契机。深度学习在处理大规模复杂数据时有着一定的优势,而现代社会的知识也随着网络信息的暴涨呈几何指数的增长,对应着知识图谱的规模越来越大,也越来越难以处理。为了高效地对大规模知识图谱进行计算处理,知识表示学习方法被提了出来<sup>[8]</sup>。知识表示学习方法是将实体与关系投影到统一的低维向量空间中,利用向量对实体和关系进行语义计算。知识表示学习在语义搜索,智能问答以及知识图谱补全等方面有着广泛的应用。然而在实际应用中,知识图谱都会存在不同程度的信息缺失,这样导致很多知识图谱都是不完备的。为了尝试解决这一问题,研究者们提出引入一些额外的上下文信息

(例如文本)来弥补知识图谱中的信息缺失,许多引入了额外信息的知识表示学习模型被提出。这些上下文信息都与知识图谱中的实体和关系信息是有关联的,这在一定程度上弥补了知识图谱信息的缺失,但依然存在许多问题。比如引入的文本信息包含了大量噪音,影响了模型的性能,另外这些模型大多只关注如何处理引入的外部信息,却忽视了知识图谱自身的结构化信息,没有考虑到上下文信息的多样性。因为知识图谱内部的结构化信息对于每一个实体和关系也是一种上下文信息,比如与某个实体直接相连接的若干实体,或者某几个实体组成的路径等,这些都可以看作知识图谱内部的结构化信息。由于这类信息不会包含噪音,所以会比较准确地刻画对应的实体和关系。因此有必要在引入了文本等其它额外信息的同时对知识图谱自身的结构化信息也加以考虑,即有必要对结合了多种类型上下文信息的知识表示学习模型进行探索。

目前大多数知识表示学习模型在建模时只考虑了某一类上下文信息,比如只考虑内部的结构化信息或者只考虑外部的文本等其它信息,这样会使得上下文信息的类别不够丰富,无法从多个角度对知识图谱进行补全。针对这一问题,本文提出了一种基于混合上下文的知识表示学习方法,即同时考虑文本信息以及知识图谱内部的结构化信息。该方法在充分利用了知识图谱中固有的结构化信息的基础上引入了外部的文本信息,利用包含实体的文本对知识图谱进行补充。本文余下的内容分为五个小节,其中第一小节介绍了目前知识表示学习的一些代表工作,第二小节对引入外部文本的问题进行了定义,第三小节重点介绍处理文本的方法,

第四小节介绍实验的评估标准，并与其它模型进行对比，展示实验结果，最后进行总结与展望，下面分别对这几个小节进行介绍。

## 1 相关工作

目前的知识表示学习模型可以分为两类：一类注重利用知识图谱本身的结构化信息，例如知识图谱里的三元组，实体和关系路径，子图等等，这些都可以看作实体和关系的上下文信息，这类上下文信息可以称作知识图谱的内部上下文信息。另一类模型主要关注怎样把文本、图像或规则等信息加入到知识图谱中，这类信息可以看成知识图谱的外部上下文信息。

基于内部上下文信息的模型又可以大致分为三类：一类是基于三元组的模型，不考虑三元组以外的任何额外信息。这些模型主要包含了以 TransE<sup>[8]</sup>、TransH<sup>[9]</sup>、TransR<sup>[10]</sup>、TransD<sup>[11]</sup>以及 TransSparse<sup>[12]</sup> 等为代表的一系列模型，它们通过对头实体和尾实体之间的距离进行测量，从而评估三元组。第二类是基于路径的模型，这类模型充分利用了知识图谱中的路径信息，例如 PTransE<sup>[13]</sup> 将实体间的多个关系进行语义组合得到路径的向量表示，然后将这种路径信息加入到模型的构建中。第三类是将三元组周围的上下文信息与模型的构建结合起来。例如 GAKE<sup>[14]</sup> 模型，它将实体和关系作为主语定义了三种上下文，并利用这些上下文信息来预测当前的主语信息，但它并没有对实体和关系进行有效地区分。TCE<sup>[15]</sup> 模型在 GAKE<sup>[14]</sup> 模型的基础上进行了改进，它也提出了两类上下文信息：一类是邻居上下文；另一类是路径上下文，

这两类上下文信息都能在一定程度上反映了目标实体的语义信息。

这些只考虑了内部上下文信息的模型只关注知识图谱自身的结构信息，对于图谱中缺失的信息却无能为力。实际上大多数知识图谱都会存在不同程度的信息缺失，为了弥补知识图谱中缺失的信息，研究者们开始将文本等外部上下文信息加入到模型中。第一个引入外部信息的模型<sup>[15]</sup> 利用词向量将文本向量化，然后利用维基百科锚等操作将文本和知识图谱联系起来。清华大学的刘知远等提出了 DKRL<sup>[17]</sup> 模型，他们分别将实体和文本进行向量表示，其中对文本的处理利用了循环神经网络。类似地，TEKE<sup>[18]</sup> 模型通过对文本中的实体进行标注，构造实体与词语的共现网络，从而得到实体与关系在文本中的上下文。其中实体在文本中的上下文为与实体共现的词，关系在文本中的上下文为同时与头尾实体共现的词。另外还有一些引入其它类型信息的模型，如引入逻辑规则的模型<sup>[19,20]</sup>，引入了实体类别的模型<sup>[21,22]</sup> 等。这些模型都从不同的角度对知识图谱进行了信息的补全。

无论是基于内部上下文还是外部上下文信息构建的模型，都只考虑了某一种类型的上下文信息，为了充分利用这两类上下文信息，我们提出了一种将文本信息与三元组上下文信息结合的模型，为了使模型更有效，我们选取了已充分利用了内部上下文信息的 TCE<sup>[15]</sup> 模型作为基础，在此基础上加入文本信息。

## 2 文本的引入与问题定义

由于我们的模型是在 TCE<sup>[15]</sup> 模型的基础上

构建的,所以重点介绍文本的引入与处理方法。许多知识图谱中的三元组很不完整,这样会对许多应用带来影响,所以需要通过一些方法对缺失的信息进行补充。近年来从文本中抽取实体和关系的方法越来越受到关注,特别是基于关系抽取的远程监督模型<sup>[23]</sup>越来越受到关注。这个模型包含了一个著名的假设,即如果两个

实体  $h$  和  $t$  在知识图谱中存在某种关系  $r$ , 那么包含这两个实体指称的句子通常也表达了其间存在的这种关系(图1),利用这种方法可以将知识图谱中的实体与关系和文本进行关联。本文基于远程监督的假设,期望通过包含实体指称的句子对知识图谱表示学习提供语义信息的补充,从而提升效果。

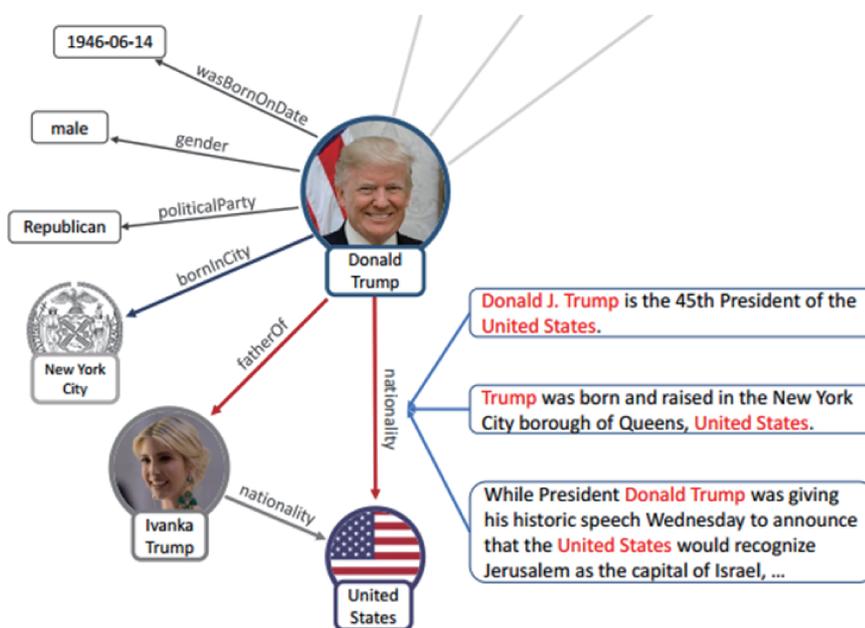


图1 引入文本的示例

首先引入一些符号的定义。知识图谱中关系的集合表示为  $R$ , 实体的向量表示为  $\theta_E \in R^{m \times k}$ , 关系的向量表示为  $\theta_R \in R^{n \times k}$ , 其中  $m$  和  $n$  分别表示为实体和关系的总个数,  $k$  表示实体和关系的维度。文本中词语的向量用参数  $\theta_V \in R^{|\mathcal{V}| \times k}$  表示。模型包含的参数为  $\theta = \{\theta_E, \theta_R, \theta_V\}$ 。文本语料记为  $D$ , 它包含若干个句子  $s$ , 即  $D = \{s_1, \dots, s_{|D|}\}$ , 而每个句子又都是若干个词语  $w$  的序列, 记为  $s = \{w_1, \dots, w_{|s|}\}$ 。文本中的每个句子均包含两个实体的指称 (mention), 句子  $s$

中两个实体之间存在的潜在关系记为  $r_s \in R$ 。由于加入了文本, 所以我们的目标是对知识图谱与文本进行联合训练, 从而获取知识图谱中实体与关系以及文本中词语的向量表示。为此我们需要得到在引入文本的条件下最大化知识图谱的联合概率  $P$ , 对于一个知识图谱  $G$ , 定义优化目标:  $\hat{\theta} = \arg_{\theta} \max P = (G|D; \theta)$ , 然后将目标函数定义为:

$$P = (G|D; \theta) = \prod f(h, r, t) \quad (1)$$

其中  $h$  代表知识图谱中三元组的头实体,  $r$  表示

关系， $t$ 表示尾实体。将文本加入到TCE<sup>[15]</sup>模型中，TCE<sup>[15]</sup>模型包含两类上下文信息，分别是邻居上下文和路径上下文。邻居上下文是指与某个实体有直接关系的上下文信息。路径上下文是指从某个实体 $h$ 出发到另一个实体 $t$ 之间所包含的所有关系信息的总和。

### 3 相关方法

为了实现知识图谱与文本的联合建模，需要将知识图谱中的实体与关系和文本进行关联。对于知识图谱中的实体及文本中的实体指称，我们通过共享其向量的方法来实现语义信息的关联。例如，实体“Apple Inc.”在文本中可能以指称“Apple”的形式出现，但其对应于同一个向量，使得其既能够包含知识图谱中的结构信息，又能够包含文本中的语义信息。如果知识图谱中出现同名实体，则利用实体消歧方法进行预处理。为了对文本进行处理，我们采用深度学习的方法来进行文本语料的训练。

#### 3.1 关系与文本的表示学习

根据远程监督的假设，文本中的句子 $s$ 包含了实体 $h$ 与 $t$ 的指称，且整个句子表达了两个实体之间的关系 $r$ 。由于无法直接对句子进行计算，为了计算关系 $r$ 的概率分布 $P(r|h, t, D; \Theta)$ ，我们利用卷积神经网络模型对句子进行编码。整个文本处理框架主要分为词表示层、特征提取层和输出层。

在词表示层里，句子中的每个词都被表示成一个向量。为了更精确地描述每一个词的向

量信息，我们用两部分向量来共同描述每个词的向量 $x$ 。除了每个词本身的向量 $w$ 以外，我们还给每个词添加了一个位置特征向量 $p$ ，用于描述句子中每个单词与两个实体指称之间的相对位置。例如在句子“[Donald John Trump] is the 45th and current President of the [United States]”中，单词“is”相对于两个实体指称的位置表示为 $[1, -8]$ ，说明“is”在实体指称Donald John Trump右侧，与其距离为1，在实体指称United States左侧，与其距离为8。句子中每个单词 $w$ 相对于两个实体指称的位置表示为 $[d_1, d_2]$ ，其中 $d_1$ 和 $d_2$ 分别表示 $w$ 相对于两个实体指称的方向和距离，将两个距离向量进行首尾拼接可以得到位置特征向量 $p \in R^{k \times 2}$ ，其中 $k$ 表示位置向量的维度。将单词的词向量和位置特征向量进行拼接，可以得到每个单词的向量表示，即 $x = [w; p]$ 。因此，每个句子 $s$ 都可以表示为一个向量序列：

$$s = \{x_1, \dots, x_{|s|}\} = \{[w_1; p_1], \dots, [w_{|s|}; p_{|s|}]\} \quad (2)$$

在特征提取层我们使用卷积神经网络从词表示层输出的向量中提取出句子级别的特征。卷积层使用大小为 $d$ 的窗口在输入句子的向量序列上进行滑动，对窗口里的向量进行组合，将同一窗口里的向量首尾拼接，得到向量 $z_i \in R^k$ 。然后，对组合向量 $z_i$ 进行卷积操作：

$$h_i = \tanh(Wz_i + b) \quad (3)$$

其中，参数矩阵 $W \in R$ 是卷积核， $b \in R$ 是一个偏置向量， $\tanh(\cdot)$ 表示双曲正切函数。为了进一步在隐藏层向量的每一个维度中确定最有用的特征，在进行卷积操作以后，对隐藏层向量进行最大池化操作。假设隐藏层中向量的个数为 $n_h$ ，那么输出向量 $q$ 的每一维可以通过以

下方式确定:

$$q^i = \max(h_1^i, \dots, h_{nh}^i) \quad (4)$$

输出层将特征提取层输出的向量转化为每个关系的概率分布。首先将特征提取层的输出向量通过一个矩阵变换, 转换为一个  $n$  维向量  $y = Mq$ , 矩阵  $M$  将  $q$  转换为  $n$  维向量, 其中的每一维表示每个关系的评分, 然后将向量  $y$  输入到 Softmax 分类器中, 这样就得到概率了分布。

### 3.2 基于注意力机制的文本编码

在实际场景中, 一对实体可能出现在多个句子中, 不同的句子对这对实体之间关系预测的作用也不同, 我们需要将多个句子中的信息进行融合, 从而获取关系的概率分布。对于一个知识图谱中的三元组  $(h, r, t)$ , 并不是所有包含头尾实体指称的句子都蕴含关系  $r$ 。因此, 需要模型在预测关系  $r$  的同时, 根据每个句子与该关系的关联程度有所侧重地选择句子中的信息, 以不同的权重对不同的句子进行组合。因此本文使用注意力机制<sup>[24]</sup>来对不同句子的语义特征进行组合, 图 2 是引入文本的整体流程。

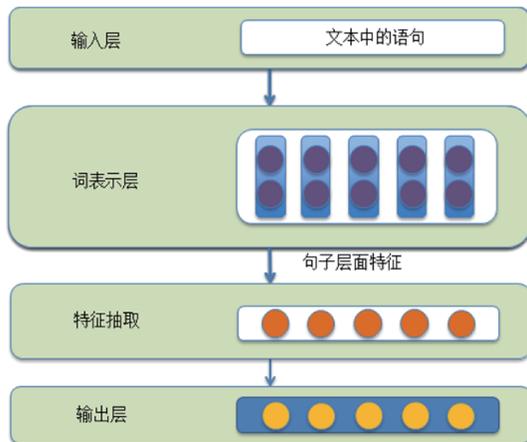


图 2 文本处理流程

将包含实体  $h$  和  $t$  指称的句子的集合记作  $S_{(h,t)}$ ,  $S_{(h,t)} = s_1, \dots, s_c$ , 其中包含了  $c$  个句子。对于其中的每一个句子, 利用卷积神经网络将其编码成一个向量  $q$ , 并将其输入到一个全连接层, 得到每个句子的中间向量表示  $e$ , 然后, 根据实体与句子的向量表示, 计算得出每个句子的权重:

$$\alpha_j = \frac{\exp((t-h) \cdot e_j)}{\sum_{i=1}^c \exp((t-h) \cdot e_i)} \quad (5)$$

其中,  $t-h$  是基于 TransE 的假设  $h+r \approx t$ , 用于表示  $h$  与  $t$  之间存在的潜在关系。通过  $t-h$  与句子中间表示  $e_j$  的内积, 可以得到每个句子与实体对之间潜在关系的关联程度。该式通过一个 Softmax 的形式, 使得所有句子的权重符合概率分布条件。随后, 可以将每个句子的向量加权求和得到所有句子的向量表示:

$$s = \sum_{j=1}^c \alpha_j q_j \quad (6)$$

将输出向量通过转换矩阵  $M$  变换为  $n$  维向量得到每个关系的评分  $y = Ms$ , 最后将  $y$  输入到 Softmax 函数中得到最终概率分布:

$$P(r|h,t,D;\Theta) = \frac{\exp(y_r)}{\sum_{r \in R} \exp(y_r)} \quad (7)$$

整个模型的学习采用如下算法:

#### 算法 4.1 基于混合上下文的模型学习算法

输入: 知识图谱  $G$ , 文本语料  $D = \{s\}$ , 实体集合  $E$ , 关系集合  $R$ , 向量维度  $k$

输出: 更新后的参数  $\Theta$

1: **initialize:**  $e \leftarrow \text{uniform}\left(\frac{-6}{\sqrt{d}}, \frac{6}{\sqrt{d}}\right)$  for each entity  $e \in E$

2:  $r \leftarrow \text{uniform}\left(\frac{-6}{\sqrt{d}}, \frac{6}{\sqrt{d}}\right)$  for each relation  $r \in R$

3:  $w \leftarrow \text{word2vec}(w)$  for each relation  $r \in R$

```

4: loop:
5:    $\mathbf{e} \leftarrow \mathbf{e}/\|\mathbf{e}\|$  for each entity  $e \in E$ 
6:    $\mathbf{r} \leftarrow \mathbf{r}/\|\mathbf{r}\|$  for each relation  $r \in R$ 
7:    $T_{batch} \leftarrow \emptyset$ 
8:   for each batch  $S_{batch}$  in  $G$ :
9:     for each triple  $(h, r, t)$  in  $S_{batch}$ :
10:       $H^- \leftarrow \text{sample}(\{h' | (h', r, t) \notin G\})$ 
11:       $R^- \leftarrow \text{sample}(\{r' | (h, r', t) \notin G\})$ 
12:       $T^- \leftarrow \text{sample}(\{t' | (h, r, t') \notin G\})$ 
13:       $D \leftarrow \{s | m_h, m_t \in s\}$ 
14:       $T_{batch} \leftarrow T_{batch} \cup ((h, r, t), (H^-, R^-, T^-), D)$ 
15:     end for
16:   Update the gradient of  $-\sum_{((h,r,t),(H^-,R^-,T^-),D) \in T_{batch}} \nabla P(G|D; \Theta)$ 
17:   end for
18: end loop

```

在参数初始化阶段，首先初始化实体与关系的向量。引入的文本语料通过 word2vec 向量进行初始化。然后对每个实体与关系的向量进行标准化，每次从知识图谱训练集中取出一批数据  $S_{batch}$ ，对于其中的每个三元组分别对头实

体、关系和尾实体进行负采样，得到三个负例的集合。对于  $S_{batch}$  中的每个实体对，从文本语料中找出包含这两个实体指称的句子，随后，对这批数据利用目标函数计算损失，并根据学习率对参数进行更新，当算法达到最大迭代上限则停止迭代。

## 4 实验与评估

### 4.1 数据集

本文实验的所用数据集包含两部分，第一部分是知识图谱，第二部分是文本语料。我们选取的知识图谱数据分别是 FB15K 和 FB15K-237。FB15k<sup>[8]</sup> 是从 Freebase 中抽取的一个子集，它总共有 592,213 个三元组，其中包含了 14,951 个实体和 1,345 个关系。FB15k-237<sup>[25]</sup> 是 FB15k 的一个子集，它包含了 310,116 个三元组，其中包含了 14,541 个实体和 237 个关系。这两个数据集被分为三个部分：训练集、验证集和测试集。训练集用于模型的训练，验证集用于选择超参数，测试集用于测试模型最终的效果，详见表 1。

表 1 知识图谱数据集详细数据

数据集	实体个数	关系个数	训练集规模	验证集规模	测试集规模
FB15k	14,951	1,345	483,142	50,000	59,071
FB15k-237	14,541	237	272,115	17,535	20,466

对于文本语料的选择，本文使用纽约时报数据 NYT10<sup>[26]</sup>，该语料对其中出现的 Freebase 中的实体进行了标注。本文根据远程监督的假设，抽取出包含以上两个数据集中实体对的句子作为训练文本。对于

FB15k，本文从 NYT10 中抽取了 194,385 个句子，对于 FB15k-237，本文从 NYT10 中抽取了 78,978 个句子，然后根据这些实体对在知识图谱中存在的关系对抽取的这些句子进行标注。

## 4.2 评估标准

本文使用 Mean Rank 和 Hits@10 作为评估标准。给定一个测试三元组  $(h, r, t)$ ，首先将其中的头实体  $h$ 、尾实体  $t$  或关系  $r$  去掉，形成一个不完整的三元组。以预测头实体  $h$  为例，将头实体  $h$  从三元组中去掉以后得到  $(?, r, t)$ ，然后用任意实体  $e$  进行替换，组成新的三元组  $(e, r, t)$ ，通过评分函数计算其评分  $f(e, r, t)$ 。由于  $e$  可以取任意实体，因此对于每个实体  $e$  都能够得到一个评分  $f(e, r, t)$ ，将这些实体按得分的降序排列，得到一个实体的序列  $L = (e_1, e_2, \dots, e_m)$ ，在这个序列的基础上，可以分别得出以下两种评估标准：

Mean Rank 表示对于所有测试三元组，原始头实体  $h$  在序列  $L$  中排名的平均值。Mean Rank 越小，表示模型的性能越好。

Hits@10 也是一个常用的评估标准，它表示在所有测试三元组中，原始头实体  $h$  在序列  $L$  中的排名小于等于  $n$  的三元组个数占测试三元组总个数的比例。Hits@10 越大，表示模型的性能越好。在实体预测的场景下，一般来说正确的实体在序列  $L$  中的排名尽可能靠前，而错误的实体排名尽可能靠后。但在序列  $L$  中，排名在原始头实体  $h$  前面的某个实体  $e$  与  $(?, r, t)$  组成的三元组  $(e, r, t)$  可能是正确的，即存在于原始知识图谱中。对于这些实际正确的实体，由于模型已经让其获得了一个靠前的排名，因此不能将其视作错误的实体，需要将这些实体对当前测试实体排名的影响消除。具体做法是从序列  $L$  中去掉除原始头实体  $h$  之外实际正确的实体  $e$ ，即如果三元组  $(e, r, t)$  存在于原始知

识图谱中，则将实体  $e$  从序列中去除。按照惯例，将这种设置记为 Filtered（或 Filt.）。

## 4.3 模型评估

基于 TCE<sup>[15]</sup> 模型，我们引入了一些文本参数，分别如下：卷积神经网络滑动窗口大小  $d=3$ ，隐藏层维度  $k_h=230$ ，实体、关系和单词的维度  $k_w=50$ ，位置向量  $k_p$  的维度为  $k_p=5$ ，模型的训练和 TCE<sup>[15]</sup> 模型一样用负采样的方法。我们使用 DKRL<sup>[17]</sup>、TEKE<sup>[18]</sup>、DESP<sup>[27]</sup>、E + DISTMULT<sup>[28]</sup> 和 Conv-E + DISTMULT<sup>[28]</sup> 作为对比实验，因为这几个模型都将文本信息引入了知识图谱表示学习模型中。

表 2 基于混合上下文的模型在 FB15k 和 FB15k-237 上进行实体预测的结果

数据集	模型	Hits@10(Filtered)
FB15k	DKRL	67.4
	TEKE	73.0
	DESP	77.3
	TCE + Text	77.5
FB15k-237	E + DistMult	60.2
	Conv-E + DistMult	61.1
	TCE + Text	60.8

基于混合下文与文本的表示学习模型实验结果如表 2 所示，其中，“TCE + Text”项表示本文提出的基于混合上下文的表示学习模型。从表中可以看出，基于混合上下文的模型在 FB15k 数据集上超越了其它同类模型。另外与传统的表示学习模型相比，本文提出的模型在 Hits@10（Filtered）指标上取得的效果也能够超越大部分模型，说明这种基于混合上下文信息的模型确实可以提升知识图谱表示学习的

效果。

## 5 总结与展望

本文提出了一种基于混合上下文的知识表示学习方法,将文本信息与三元组上下文信息结合。通过大量丰富文本信息的引入,进一步对知识图谱的内容进行了补充。实验表明这种加入了混合信息的模型在实体预测方面优于现有的一些模型,但由于该方法在建模的过程中已经用到了关系路径的信息,所以无法对关系进行预测,这也是后续工作需要完善的方向。无论是外部的文本信息还是知识图谱内部的结构化信息,都只是众多信息类别中的一小部分。目前已有将多种类别的信息结合起来进行建模的多模态知识表示方法<sup>[29]</sup>,但都往往顾此失彼,并不能真正地对多种类型的信息进行有效地整合。随着信息的不断积累,知识图谱的规模也会越来越大,知识种类也会越来越多,如何从多种维度进行知识表示学习的研究是我们应该考虑的问题。因为当前的人工智能应用往往达不到人们的预期,或知识表达能力不足,或缺少逻辑推理等等,其中原因之一可能就是缺乏对知识多样性的表征。所以将时间、空间、文本、规则、图像,甚至视频等多种知识源进行有效地整合,可能会使知识图谱成为人们迈向真正人工智能时代的坚实基础。

## 参考文献

- [1] 漆桂林,高桓,吴天星.知识图谱研究进展[J].情报工程,2017,3(1):4-25.
- [2] 曹明宇,李青青,杨志豪,等.基于知识图谱的原发性肝癌知识问答系统[J].中文信息学报,2019,33(6):88-93.
- [3] 于娟,黄恒琪,席运江,等.基于图数据库的人物关系知识图谱推理方法研究[J].情报科学,2019,37(10):8-12.
- [4] 侯梦薇,卫荣,陆亮,等.知识图谱研究综述及其在医疗领域的应用[J].计算机研究与发展,2018,55(12):2587-2599.
- [5] 冯小兰,赵小兵.汉藏双语旅游领域知识图谱系统构建[J].中文信息学报,2019,33(11):64-72.
- [6] 陈果,许天祥.小规模知识库指导下的细分领域实体关系发现研究[J].情报学报,2019,38(11):1200-1211.
- [7] 刘知远,孙茂松,林衍凯,等.知识表示学习研究进展[J].计算机研究与发展,2016,53(2):247-261.
- [8] Bordes A, Usunier N, García-Durán A, et al. Translating Embeddings for Modeling Multi-relational Data[C]. Advances in Neural Information Processing Systems 26: 27<sup>th</sup> Annual Conference on Neural Information Processing Systems, 2013: 2787-2795.
- [9] Wang Z, Zhang J, Feng J, et al. Knowledge Graph Embedding by Translating on Hyper-planes[C]. Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence, 2014: 1112-1119.
- [10] Lin Y, Liu Z, Sun M, et al. Learning Entity and Relation Embeddings for Knowledge Graph Completion[C]. Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, 2015, 2181-2187.
- [11] Ji G, He S, Xu L, et al. Knowledge Graph Embedding via Dynamic mapping Matrix[C]. In: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, 2015: 687-696.
- [12] Ji G, Liu K, He S, et al. Knowledge Graph Completion with Adaptive Sparse Transfer Matrix[C]. In: Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, 2016: 985-991.
- [13] Lin Y, Liu Z, Luan H, et al. Modeling Relation Paths

- for Representation Learning of Knowledge Bases[C]. In: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, 2015: 705-714.
- [14] Feng J, Huang M, Yang Y, et al. GAKE: Graph Aware Knowledge Embedding[C]. In: Proceedings of the 26th International Conference on Computational Linguistics, 2016: 641-651.
- [15] Shi J, Gao H, Qi G, et al. Knowledge graph embedding with triple context[C]. Proceedings of the 2017 ACM on Conference on Information and Knowledge Management. ACM, 2017: 2299-2302.
- [16] Wang Z, Zhang J, Feng J, et al. Knowledge Graph and Text Jointly Embedding[C]. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, 2014: 1591-1601.
- [17] Xie R, Liu Z, Jia J, et al. Representation Learning of Knowledge Graphs with Entity Descriptions[C]. In: Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, 2016: 2659-2665.
- [18] Wang Z, Li J. Text-Enhanced Representation Learning for Knowledge Graph[C]. In: Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, 2016: 1293-1299.
- [19] Rocktäschel T, Singh S, Riedel S. Injecting Logical Background Knowledge into Embeddings for Relation Extraction[C]. The 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2015: 1119-1129.
- [20] Demeester T, Rocktäschel T, Riedel S. Lifted Rule Injection for Relation Embeddings[C]. Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, 2016: 1389-1399.
- [21] Krompaß D, Baier S, Tresp V. Type-Constrained Representation Learning in Knowledge Graphs[C]. In: The Semantic Web - ISWC 2015 14<sup>th</sup> International Semantic Web Conference, 2015: 640-655.
- [22] Hu Z, Huang P, Deng Y, et al. Entity Hierarchy Embedding[C]. Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7<sup>th</sup> International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, 2015: 1292-1300.
- [23] Mintz M, Bills S, Snow R, et al. Distant supervision for relation extraction without labeled data[C]. Proceedings of the 47<sup>th</sup> Annual Meeting of the Association for Computational Linguistics and the 4th International Joint Conference on Natural Language Processing of the AFNLP, 2009: 1003-1011.
- [24] Bahdanau D, Cho K, Bengio Y. Neural Machine Translation by Jointly Learning to Align and Translate[J]. arXiv preprint arXiv:1409.0473, 2014.
- [25] Toutanova K, Chen D, Pantel P, et al. Representing Text for Joint Embedding of Text and Knowledge Bases[C]. Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, 2015: 1499-1509.
- [26] Riedel S, Yao L, McCallum A. Modeling Relations and Their Mentions without Labeled Text[C]. Proceedings of Machine Learning and Knowledge Discovery in Databases, European Conference, 2010. 148-163.
- [27] Zhong H, Zhang J, Wang Z, et al. Aligning Knowledge and Text Embeddings by Entity Descriptions[C]. Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, 2015: 267-272.
- [28] Toutanova K, Chen D, Pantel P, et al. Representing Text for Joint Embedding of Text and Knowledge Bases[C]. Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, 2015: 1499-1509.
- [29] Pezeshkpour P, Chen L, Singh S. Embedding multimodal relational data for knowledge base completion[J]. arXiv preprint, arXiv:1809.01341, 2018.