

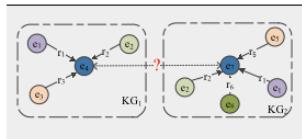
基于上下文的跨语言知识图谱实体对齐方法

Northeastern University

马新月, 聂铁铮, 申德荣, 寇月
东华大学 计算机科学与工程学院, 沈阳, 中国

背景

- 知识图谱之间存在知识丢失、知识重复、知识关联不明确等问题；
- 传统的实体对齐方法存在训练缓慢、精确度低等问题，有局限性；
- 现有的实体对齐方法未充分考虑实体的上下文信息，当中心实体没有传出任何信息时，很难为其匹配到相似实体。



当知识图谱中的实体不存在传出三元组时，无法对其进行有效建模，**难以考虑传入关系来提升实体嵌入效果。**

引言

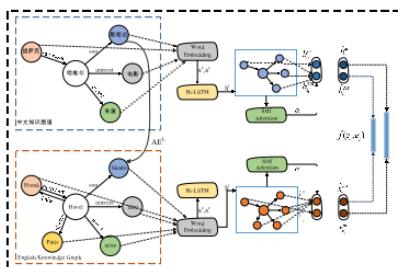
面向知识图谱的实体对齐技术可以实现知识的连接，将两类知识图谱融合为规模更大、质量更权威的领域知识图谱。**实体对齐的目标**是识别出不同知识图谱中指代对象为现实世界中同一事物的实体对。实体对齐方法有传统的实体对齐方法，**或**基于知识表示学习的实体对齐方法。

- 基于翻译的表示学习方法能够在低维稠密空间中学习实体的向量表示。
- 基于图神经网络的嵌入方法在三元组的维度上结合节点的属性信息进行复杂信息建模。

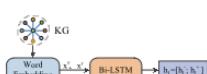
为解决多语言知识图谱实体对齐的精确度问题，在实体嵌入的过程中充分考虑实体的**上下文信息**，提出基于上下文的知识图谱实体对齐方法。

方法

本文的基于上下文的跨语言实体对齐方法，通过中心实体的上下文信息来解决前期工作的缺陷，总体框架如下图所示：



上图为知识图谱实体对齐方法总体框架，首先根据词嵌入来分别对两个知识图谱中实体的关系进行语义编码，利用Bi-LSTM得到节点间关系的初始特征表示，接着利用图注意力机制进一步处理实体的嵌入表示，捕获每个三元组中的内部依赖性对邻居节点对中心实体嵌入表示的**重要程度**，最后利用基于距离的度量函数进行实体对齐。



特征提取层

通过基于词的Bi-LSTM来学习三元组中的实体间关系的嵌入表示，将深层次的语义特征向量表示作为实体的初始特征向量表示。

实体嵌入层

结合图注意力机制对知识图谱中的中心实体的上下文建模，最终的嵌入表示结合了所有中心实体相关的邻居节点信息及关系信息，添加超参数 β 来平衡邻居类别(传入/传出)的重要程度。

$$\begin{aligned} \hat{\mathbf{h}}_e &= \text{SEAL}(O_{\text{out}}, \alpha_{\text{out}} \mathbf{h}_e^{\text{in}} + O_{\text{in}}, \alpha_{\text{in}} \mathbf{h}_e^{\text{out}}) \\ &+ \frac{1}{|\mathcal{N}_e|} \left(\sum_{n \in \mathcal{N}_e} \sum_{r \in \mathcal{R}_{en}} \mathbf{h}_r^{\text{in}} \right)^T \left(\sum_{n \in \mathcal{N}_e} \mathbf{h}_n^{\text{in}} \right) \\ &+ \frac{1}{|\mathcal{N}_e|} \left(\sum_{n \in \mathcal{N}_e} \sum_{r \in \mathcal{R}_{en}} \mathbf{h}_r^{\text{out}} \right)^T \left(\sum_{n \in \mathcal{N}_e} \mathbf{h}_n^{\text{out}} \right) \\ &+ \alpha_{\text{out}} \nabla^T \left[\sum_{n \in \mathcal{N}_e} \mathbf{h}_n^{\text{in}} \right] \left(\frac{1}{|\mathcal{N}_e|} \sum_{n \in \mathcal{N}_e} \mathbf{h}_n^{\text{in}} \right) \\ &+ \alpha_{\text{in}} \nabla^T \left[\sum_{n \in \mathcal{N}_e} \mathbf{h}_n^{\text{out}} \right] \left(\frac{1}{|\mathcal{N}_e|} \sum_{n \in \mathcal{N}_e} \mathbf{h}_n^{\text{out}} \right) \\ &+ \alpha_{\text{out}} \text{softmax}(\mathbf{z}_{\text{out}}) = \frac{\exp(\alpha_{\text{out}} z_{\text{out}})}{\sum_{i \in \mathcal{N}_e} \exp(\alpha_{\text{out}} z_i)} \\ &+ \alpha_{\text{in}} \text{softmax}(\mathbf{z}_{\text{in}}) = \frac{\exp(\alpha_{\text{in}} z_{\text{in}})}{\sum_{i \in \mathcal{N}_e} \exp(\alpha_{\text{in}} z_i)} \end{aligned}$$

实体对齐

度量实体对之间的距离，从距离排序列表中选择候选配对实体，距离越近的实体越相似，在候选配对列表中的排名越靠前。

$$D(c_i, c_j) = \sqrt{h_{c_i}^{(k)} - h_{c_j}^{(k)}} + \sqrt{h_{c_i}^{(l)} - h_{c_j}^{(l)}}$$

模型训练

利用基于距离排序损失函数，使两个知识图谱中相似的实体对在向量空间中的距离尽可能地接近。

$$\begin{aligned} L &= \sum_{i,j \in \mathcal{C}} \sum_{k,l \in \mathcal{C}} D(c_i, c_k) + D(c_j, c_l) \\ &+ \sum_{i,j \in \mathcal{C}} \sum_{k \in \mathcal{C}, k \neq i, j} D(c_i, c_k) + D(c_j, c_k) \end{aligned}$$

结果

实验对比结果如下表所示，粗体数字代表最佳性能。

	DST100K	DST210K	DST310K
HTransE	30.8	61.4	93.4
MTransE	37.3	74.3	94.9
JAP	41.3	72.1	94.4
GCN-Align	53.0	73.3	96.3
KGEA	56.4	72.4	94.6

可以看出本文的实体对齐方法KGEA对比之前的方法有一定的提升，这表明将Bi-LSTM与注意力机制相结合可以挖掘出更多元组间的潜在关联信息。

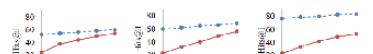
本文的方法KGEA与对比方法MTransE, JAP和GCN-Align在DST100K数据集上的实体对齐结果

	DST100K	DST210K	MTransE
HTransE	28.1	52.6	0.163
JAP	31.8	58.6	0.111
GCN-Align	53.0	73.3	0.603
KGEA	56.1	77.5	0.695

本文的方法KGEA与对比方法MTransE, JAP和GCN-Align在DST210K数据集上的实体对齐结果

使用不同比例的先验种子集作为训练数据来验证种子集数目对实体对齐结果的影响，对比实验结果如下图所示。

从图中可以看出，在所有的数据集中，KGEA的性能均为最优，KGEA利用了详细的实体关联信息进行嵌入表示，即使种子集比例较小，也能获得良好的效果。



不同的种子集比例对本文方法KGEA对齐效果的影响，其x轴是用于对齐的种子集的比例，y轴是Hits@1分数

结论

- 本文提出了一种新的基于上下文的知识图谱实体对齐方法，通过探索中心实体周围的复杂关联信息来学习实体嵌入。
- 利用Bi-LSTM模型与图注意力机制构造邻居的关系信息结构信息，更好地学习实体的向量表示来实现实体对齐。
- 实验结果表明本文的方法在三个多语言知识图谱数据集两个大规模数据集中均取得了最好的实体对齐效果。

致谢

本文得到国家重点研发计划项目(No. 2018YFB1003404)、国家自然科学基金(No. U1811261, 61672142, 61472070)、国防基础科研计划(JCKY2018205C012)资助。