

《智能信息处理》课程作业

知识表示学习方法研究

贾康

考核	到课 [10]	作业 [20]	考试 [70]	课程成绩 [100]
得分				

2021 年 12 月 21 日

知识表示学习方法研究

贾康

(大连海事大学 信息科学与技术学院 大连 116026)

摘要 近年来,知识表示学习已经成为知识图谱领域研究的热点。为了及时掌握当前知识表示学习方法的研究现状,通过归纳与整理,将具有代表性的知识表示方法进行了介绍和归类,主要分为传统的知识表示模型、改进的知识表示模型、其他的知识表示模型。对每一种方法解决的问题、算法思想、应用场景、评价指标、优缺点进行了详细归纳与分析。通过研究发现,当前知识表示学习主要面临关系路径建模、准确率、复杂关系处理的挑战。针对这些挑战,展望了采用关系的语义组成来表示路径、采用实体对齐评测指标、在实体空间和关系空间建模,以及利用文本上下文信息以扩展 KG 的语义结构的解决方案。

关键词 知识图谱; 知识表示学习; 实体对齐; 链接预测; 三元组分类

Survey of knowledge representation learning methods

Kang Jia

(Department of Information Science Technology, Dalian Maritime University, Dalian 116026)

Abstract: In recent years, knowledge representation learning has become a hot topic in the field of knowledge graph. In order to grasp the current research status of knowledge representation learning methods in time, this paper introduces and classifies the representative knowledge representation methods through induction and sorting, which are mainly divided into traditional knowledge representation model, improved knowledge representation model and other knowledge representation models. This paper summarizes and analyzes the problems, algorithm ideas, application scenarios, evaluation indicators, advantages and disadvantages of each method in detail. Through research, this paper finds that the current knowledge representation learning mainly faces the challenges of relationship path modeling, accuracy and complex relationship processing. Aiming at these challenges, this paper looks forward to using semantic composition of relationship to represent path, using entity alignment evaluation index, modeling in entity space and relationship space, and using text context information to expand the solution of kg's semantic structure Resolution.

Keywords: knowledge graph; knowledge representation learning; entity alignment; link prediction; triple classification

1 介绍

知识图谱(knowledge graph, KG)是一种语义网络。这种语义网络由实体节点和关系组成,通过三元组(h, r, t)的方式进行表示。自从 KG 被提出以后,已经有不少的大型 KG 被构建出来,比较典型的有 WordNet、知立方、知心等^[1,2]。这些大型 KG 极大方便了人们对知识的查找。然而,由于大型 KG 采用了三元组(h, r, t)的方式表示知识,导致计算起来较复杂,给知识的推理带来了一定的麻烦^[3-5]。因此,近年来学者们将注意

力转向了知识表示学习,希望通过知识表示学习来解决 KG 面临的问题。知识表示学习指的是通过不断学习 KG 中的实体和关系,从而进行知识的表示。受到词向量的启发,知识表示学习最近取得了显著的突破,采用了将实体和关系嵌入到向量空间的方法,从而很好解决了 KG 存在的数据稀疏、不便于计算的问题。同时,这种方法让知识获取、知识融合、知识推理的效率得到了进一步的提升^[6-8]。由于以上的优点,知识表示

学习越来越受到学者们的关注。但是,该方向仍然面临很多的挑战,如关系路径问题、复杂关系处理问题等。本文将对现有的知识表示学习方法进行分析与研究,找出该方向面临的主要挑战,总结出可行的解决办法,并对未来的研究趋势进行展望。

2 知识表示学习基本概念与特点

2.1 知识表示学习基本概念

知识必须经过合理的表示才能被计算机处理。知识表示学习是对现实世界的一种抽象表达^[9-11]。评价知识表示的两个重要因素是表达能力与计算效率。一个知识表示应该具有足够强的表达能力,才能充分、完整地表达特定领域或者问题所需的知识。同时,基于这一知识表示的计算求解过程也应有足够高的执行效率。知识的表示方法主要分为符号表示和数值表示。比如,人们通常用“柏拉图”这三个字符指代哲学家柏拉图,用关联图(点、边等符号)表示关系,用 \Rightarrow 表示两个命题之间的逻辑蕴涵关系。而知识的数值表示用标量、向量、概率分布等数值表达事实与知识。

2.2 知识表示学习的特点

知识表示学习通过将实体和关系映射到向量空间,实现了知识的有效表示,它主要有以下的特点:

a)在词向量的启发下,KG中的实体和关系被映射到连续的向量空间,并包含一些语义层面的信息,便于在下游任务中更加方便地操作KG,例如问答任务、关系抽取等;

b)简化了计算,提高了效率。KG中的数据是采用三元组的方式进行表示,这种表示方式使得数据比较稀疏,在进行实体关系的推理时,需要设计专门的算法,并进行复杂的计算效率比较低。但是,在进行知识表示学习后,通过损失函数就可以计算语义相似度,计算效率显著提升。

c)提高了计算的精确性。知识表示学习能够通过实体向量与关系向量距离的计算,处理低频对象的语义表示问题,从而提高了精确性。

d)使知识融合变得更容易。知识可能来源于不同的知识库(knowledge base, KB),但是,表达的对象可能是同一个对象,只有将知识进行了融合,才能够得到准确的信息,而知识表示学习能使知识融合变得更容易。

3 知识表示学习方法研究进展

当前,知识表示学习方法的研究主要集中在三大类:传统的知识表示模型、改进的知识表示模型、其他的知识表示模型。

3.1 传统的知识表示模型

传统的知识表示模型指的是翻译模型里的经典模型,包括TransE模型、TransR模型、TransD模型、TransG模型、TransH模型。这些模型的出现,使得知识图谱嵌入(knowledge graph embedding, KGE)得到了快速的发展。

3.1.1 TransE 模型

考虑在低维向量空间中嵌入实体和多关系数据的关系的问题,TransE模型被提出。该模型的算法思想是:通过将关系解释为对实体的低维嵌入进行操作的转换来对关系进行建模(图1),采用基于翻译的思想,用三元组 (h, r, t) 表示头实体 h 到尾实体 t 利用关系 r 所进行的翻译。TransE的损失函数为

$$f_r(h, t) = \|h + r - t\|_{l_1/l_2} \quad (1)$$

在实际应用中,为了使习得的表示更有区分度,TransE模型使用了Hinge Loss目标函数,使得正负例尽可能分开:

$$L = \sum_{(h,r,t) \in S} \sum_{(h',r,t') \in S'} [\gamma + f_r(h, t) - f_r(h', t')]_+ \quad (2)$$

3.1.2 TransR 模型

由于实体、关系是不同的对象,不同的关系所关注的实体的属性也不尽相同,将它们映射到同一个语义空间,在一定程度上就限制了模型的表达能力。因此TransR模型被提了出来。该模型的算法思想是:首先将KB中的每个三元组 (h, r, t) 的头实体与尾实体向关系空间中投影,然后希望满足 $l_{hr} + l_r \approx l_{tr}$ 的关系(见图2),最后计算损失函数。TransR模型提出为每个关系构造相应的向量空间,将实体与关系在不同的向量空间分开表示。这种模型要求头尾实体在关

系 r 相对应的向量空间中的投影彼此接近即可。TransR 设定所有的计算都发生在关系空间中，并在计算三元组得分之间首先将实体向量通过关系矩阵投影向关系表示空间，即：

$$h_r = hM_r \quad (3)$$

$$t_r = tM_r \quad (4)$$

然后，利用投影到关系表示空间的头实体向量和尾实体向量进行三元组得分的计算。TransR 的损失函数为

$$f_r(h, t) = \|h + r - t\|_{l_2} \quad (5)$$

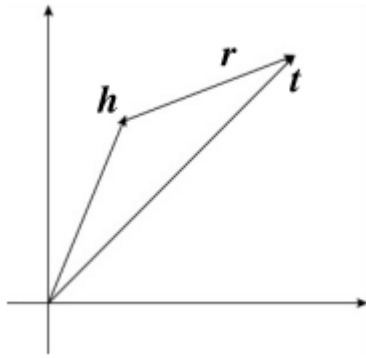


图 1 TransE 向量空间假设

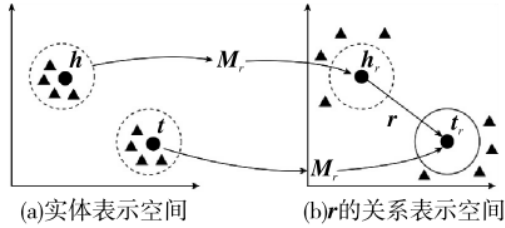


图 2 TransR 实体表示空间和关系表示空间

在图 2 中，圆圈表示特定的关系投影，三角形表示不具有 $h_r + t_r \approx r$ 关系的实体。圆点指的是满足 $h_r + t_r \approx r$ 关系的实体。

3.1.3 TransD 模型

考虑到在 KB 的三元组中，头实体和尾实体表示的含义、类型以及属性可能有较大差异，之前的 TransR 模型使它们被同一个投影矩阵进行映射，在一定程度上就限制了模型的表达能力。除此之外，将实体映射到关系空间体现的是从实体到关系的语义联系，而 TransR 模型中提出的投影矩阵仅考虑了不同的关系类型，而忽视了实体与关系之间的交互。因此，TransD 模型被提了出

来。该模型的算法思想是：将映射函数与实体、关系同时关联起来，分别定义头实体与尾实体在关系空间上的投影矩阵(见图 3 所示)。TransD 模型修改了映射函数，对于头尾实体向量 h 与 t ，它们的映射函数分别为

$$M_{rh} = r_p h_p^T + I^{m \times n} \quad (6)$$

$$M_{rt} = r_p t_p^T + I^{m \times n} \quad (7)$$

在式(6)(7)中， $I^{m \times n}$ 是单位矩阵。在这里，头实体和尾实体分别用两个不同的映射矩阵 M_{rh} 和 M_{rt} 进行投影。头实体的映射矩阵由关系向量 r_p 与头实体映射向量 h_p^T 共同决定；尾实体的映射矩阵由关系向量 r_p 与尾实体映射向量 t_p^T 共同决定。那么映射后得到的头实体和尾实体向量分别为

$$h_{\perp} = M_{rh}h \quad (8)$$

$$t_{\perp} = M_{rt}t \quad (9)$$

因此，TransD 模型的损失函数为

$$f_r(h, t) = \|h_{\perp} + r - t_{\perp}\|_{l_1/l_2} \quad (10)$$

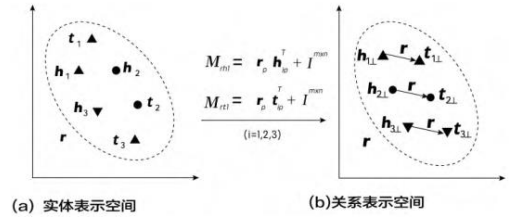


图 3 TransD 实体表示空间和关系表示空间

3.1.4 TransG 模型

在 KG 中，关系可能与对应的三元组相关联的实体对存在潜在的多个含义，即多重关系语义问题。因此，TransG 模型被提了出来。该模型的算法思想是：考虑关系 r 的不同语义，使用高斯混合模型来描述 KB 中每个三元组 (h, r, t) 的头实体与尾实体之间的关系(见图 4)，边训练边聚类，并利用特定于关系的组件矢量的组合来嵌入事实三元组。TransG 模型的生成过程如下：

a) 对于一个实体 $e \in E$ ，从标准正态分布中提取每个实体嵌入平均向量作为先验： $u_e \sim N(0, 1)$ 。

b) 对于三元组 $(h, r, t) \in \Delta$ ，首先，从关系中绘制语义成分： $\Pi_{rm} \sim \text{CRP}(\beta)$ ；其次，从正态分布中绘制头部实体嵌入向量：接着，

从正态分布中提取尾实体嵌入向量；最后，为此语义绘制关系嵌入向量：

$$u_{r,m} = t - h \sim N(u_t - u_h, (\sigma_h^2 + \sigma_t^2)E) \quad (11)$$

其中 u_h 和 u_t 分别表示头尾的平均嵌入向量， σ_h 和 σ_t 分别表示相应实体分布的方差， $u_{r,m}$ 是第 m 个关系 r 的分量平移向量。CRP 是一个狄利克雷过程，它可以自动检测语义成分。在这种情况下，得到的得分函数如下：

$$\begin{aligned} P\{(h, r, t)\} &\propto \sum_{m=1}^{M_r} \Pi_{r,m} P(u_{r,m} | h, t) \\ &= \sum_{m=1}^{M_r} \Pi_{r,m} e^{-\frac{\|u_h + u_{r,m} - u_t\|_2^2}{\sigma_h^2 + \sigma_t^2}} \end{aligned} \quad (12)$$

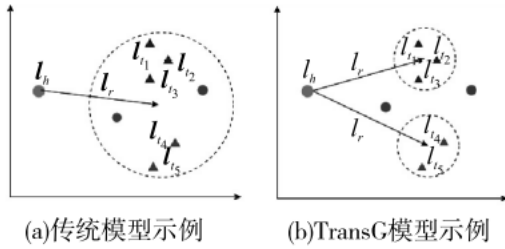


图 4 传统模型示例与 TransG 模型示例
在图 4 中，三角形指的是正确的尾实体，圆点指的是错误的尾实体，圆圈指的是投影矩阵。

3.1.5 TransH 模型

TransE 模型不能很好处理自反，一对多，多对一和多对多复杂关系映射。因此，TransH 模型被提出。此模型的算法思想是：尝试通过不同的形式表示不同关系中的实体结构，对于同一个实体而言，它在不同的关系下也扮演着不同的角色；模型首先通过关系向量与其正交的法向量选取某一个超平面，然后将头实体向量和尾实体向量沿法向量的方向投影到超平面(见图 5)，最后计算损失函数。TransH 使不同的实体在不同的关系下拥有了不同的表示形式，但由于实体向量被投影到了关系的语义空间中，故它们具有相同的维度。

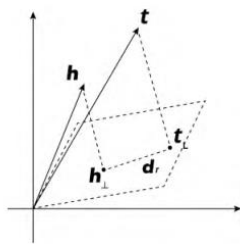


图 5 TransH 向量空间假设

4 知识表示学习面临的挑战

随着知识表示学习应用领域的不断拓展，知识表示学习方法面临着全新的挑战。知识表示学习方法正面临着关系路径建模、准确率、复杂关系处理 3 类挑战。

4.1 关系路径建模问题

关系路径建模问题指的是从关系路径的角度去推测实体之间的语义关系，但是，不是所有的关系路径都是可靠的。PTransE 仅考虑直接关系和两个实体之间的关系路径间的推理模式，从而进行学习，没有考虑关系路径之间的复杂模式，以及这种关系路径是否可靠。在实际的应用中，随着应用领域的不断扩展，复杂模式的关系将会变得很普遍，可能在一些重要的行业变得越来越重要。因此，关系路径建模问题将是一个潜在的挑战。

4.2 准确率问题

在应用场景上，除了 MGTransE 外，其他的方法采用的是链接预测、三元组分类，链接预测主要通过 Mean Rank 和 Hist@10 等评价指标来进行评定；三元组分类主要是通过 ACC 评价指标进行评定。然而，仅仅做链接预测、三元组分类存在一定的局限性，导致最终预测的准确率偏低。因此，需要在此基础上再进行其他的一些测试任务。

4.3 复杂关系处理问题

TransE、PTransE、MGTransE、KG2 在处理复杂关系方面存在忽略了多重映射关系的影响、不能很好处理一对多、多对一、多对多关系问题。这些模型在处理复杂关系时，性能会明显下降，容易出现实体表示区分不明显，甚至出现错误。因此，通过模型学习后得到的实体存在明显偏差。

5 研究展望

上面总结了知识表示学习方法面临的主要挑战，这些挑战在给知识表示方法带来难题的同时，也引导着学者们研究的方向。因此，以下的讨论将会成为今后研究的热点：

a) 采用关系的语义组成来表示路径。提出了一种通过最小化 KGE 的基于路径特定边距的损失函数的嵌入方法。对于每个路

径,通过对任何给定实体对的关系和多步关系路径之间的相关性进行编码,来自适应地确定其基于边际的损失函数。学习过程采用随机梯度下降 SGD 方法,正三元组随机遍历多次。当正三元组被访问时,负三元组 (h, r, t) 通过将 (h, r, t) 的三个分量中的一个替换为 KG 中的其他实体或关系来随机构造。此外,给定头部实体和尾部实体之间任意长度的关系路径,需要对关系和连接它们的多步关系路径之间的相关关系进行建模。实验结果表明,这种方法仍有待进一步提高。它的目标函数定义为

$$\sum_{(h,t,t') \in A} [E_{h,r,t} + \frac{1}{Z} \sum_{p \in P_{ht}} R(p/h,t) H_{p,r}] \quad (13)$$

$$E_{h,r,t} = \sum_{(h',r',t') \in \Delta'} (||h + r - t|| + \gamma_{opt} - ||h' + r' - t'||) \quad (14)$$

b) 采用实体对齐评测指标。在前面的知识表示学习面临的问题中可以看出,前面的模型一部分做了链接预测、三元组分类任务,另一部分只做了链接预测的任务。为了提高预测的准确性,可以考虑做实体对齐的任务。在做实验的过程中,可以同时分析 Mean Rank、Hist@10、Hist@1、ACC 评价指标。此外,为了使实体对齐的实验更具有说服力,可以添加效率和鲁棒性评估测试。如果需要确定 h 、 t 两个实体指代同一个对象有多大可能,则可以使用知识表示学习的得分函数对三元组 (h, r, t) 打分,根据最后的得分情况判定模型的实体关系预测效果。

c) 在实体空间和关系空间建模。可以考虑将 SE 模型和 TransE 模型结合,并且采用单层神经网络的非线性操作来精确刻画实体与关系的语义联系。在这种模型中应用面向关系的超平面的投影思想,将头尾实体映射至给定关系的超平面加以区分。在基于 WordNet 和 FreeBase 的大规模真实数据集上进行链接预测和三元组分类这两项任务。实验结果表明,此模型仍然需要进一步完善。

d) 利用文本上下文信息以扩展 KG 的语义结构。可以考虑基于共现网络和学习到的文本表示,构建实体(关系)表示模型,并将文本上下文信息合并到 KG 上的表示学习中。

首先在语义上注释该语料库中的实体,并构建一个由实体和单词组成的共现网络,以将 KG 和文本信息桥接在一起。然后,为实体和关系定义文本上下文,并将这些上下文合并到 KG 结构中。最后,使用基于常规翻译的优化过程来学习实体和关系的嵌入。学习过程是使用随机梯度下降(SGD)进行的。为避免过度拟合,使用 TransE 的结果初始化实体(关系)嵌入。实验结果表明,这种方法还有很大的扩展空间。

参考文献

- [1] David Paulius, Yu Sun. A Survey of Knowledge Representation in Service Robotics [J]. Robotics and Autonomous System, 2019 (118): 13–30.
- [2] Debarpita Santra, Swapan Kumar Basu, Jyotsna Kumar Mandal, Subrata Goswami. Rough set based lattice structure for knowledge representation in medical expert systems: Low back pain management case study [J]. Expert Systems With Applications, 2020 (145): 113084.
- [3] Kaushalya Kumarasinghe, Nikola Kasabov, Denise Taylor. Deep learning and deep knowledge representation in Spiking Neural Networks for Brain-Computer Interfaces [J]. Neural Networks, 2020 (121): 169–185.
- [4] Gayathri R, Uma V. Ontology based knowledge representation technique, domain modeling languages and planners for robotic path planning: A survey [J]. ICT Express, 2018 (4): 69–74.
- [5] Wu Zhenyong, Liao Jihua, Song Wenyan, Mao Hanling, Huang Zhenfeng, Li Xinxin, Mao Hanying. Semantic hyper-graph-based knowledge representation architecture for complex product development [J]. Computers in Industry, 2018 (100): 43–56.
- [6] Qiu Jiangnan, Zuo Min, Yan Shunin, Shi Huayan. A qualitative knowledge representation model and application for crisis events [J]. Procedia Computer Science, 2018 (126): 1828–1836.
- [7] Pridi Siregar, Nathalie Julien, Peter Hufnagel, George Mutter. A general framework dedicated to computational morphogenesis Part II– Knowledge representation and architecture [J]. BioSystems, 2018 (173) 314–334.

- [8] Nathaniel Rabb, Philip M. Fernbach, Steven. Sloman. Individual Representation in a Community of Knowledge [J]. Trends in Cognitive Sciences, 2019 (23): 10.
- [9] Hu W, Li C. Cross-Lingual Entity Alignment via Joint Attribute Preserving Embedding [C]// International Semantic Web Conference. Springer, Cham, 2017: 628-644.
- [10] Dettmers T, Minervini P, Stenetorp P, et al. Convolutional 2d knowledge graph embeddings [C]// Thirty-Second AAAI Conference on Artificial Intelligence, 2018.
- [11] Amelec Vilorio, Omar Bonerge Pineda Lezama. An intelligent approach for the design and development of a personalized system of knowledge representation [J]. Procedia Computer Science, 2019 (151): 1225–1230.