

面向链接预测的知识图谱嵌入研究综述

王 瑞, 李智杰, 李昌华, 张 颖

(西安建筑科技大学 信息与控制工程学院, 西安 710055)

摘要: 知识图谱在人工智能领域有着广泛的应用, 如信息检索、自然语言处理、推荐系统等; 然而, 知识图谱的开放性往往意味着它们是不完备的, 具有自身的缺陷; 鉴于此, 需建立更完整的知识图谱, 以提高知识图谱的实际利用率; 利用链接预测通过已有关系来推测新的关系, 从而实现大规模知识库的补全; 通过比较基于翻译模型的知识图谱链接预测模型, 从常用数据集与评价指标、翻译模型、采样方法等方面分析了知识图谱链接预测模型的框架, 并对基于知识图谱的链接预测模型进行了综述, 可为大规模知识图谱嵌入提供简单高效的识别推理方法, 增加下游人工智能控制应用任务的多样性。

关键词: 开放知识图谱; 知识图谱嵌入; 知识图谱补全; 链接预测

A Survey of Knowledge Graph Embedding Study for Link Prediction

WANG Rui, LI Zhijie, LI Changhua, ZHANG Jie

(School of Information and Control Engineering, Xi'an University of Architecture and Technology, Xi'an 710055, China)

Abstract: Knowledge graphs have a wide range of applications in the field of artificial intelligence, such as information retrieval, natural language processing, recommender systems, etc. However, the openness of knowledge graphs often means that they are incomplete and have their own flaws. In view of this, it is necessary to establish a more complete knowledge graph to improve the actual utilization of the knowledge graph. The link prediction is used to infer new relations through existing relations, so the completion of large-scale knowledge base is realized. By compared with the knowledge graph link prediction models based on translation models, the framework of the knowledge graph link prediction models is analyzed from the aspects of common data sets, evaluation indicators, translation models, and sampling methods, and the link prediction models based on knowledge graphs are reviewed. It can provide a simple and efficient identification and reasoning method for large-scale knowledge graph embedding, and increase the diversity of the downstream artificial intelligence control application tasks.

Keywords: open knowledge graph; knowledge graph embedding; knowledge graph completion; link prediction

0 引言

伴随着 Web 技术的崛起与更新迭代, 人类先后经历了以文档互联的“Web 1.0”时代与数据互联“Web 2.0”时代, 正在迈向基于知识互联的“Web 3.0”时代^[1]。同时, 随之而来的海量网络数据资源推动着人类社会进入大数据时代。如何从内容多源异质、组织结构松散的网络数据资源中有效提取组织非结构化信息和存储结构化知识变得非常重要, 同时也给“Web 3.0”提出的“知识之网”带来了极大的挑战。强大的语义处理能力和开放互联能力使得知识图谱具有良好的知识表达能力和解释性, 同时也提供了一种更好组织、管理和理解互联网海量信息的能力^[2]。知识图谱的研究起源于语义 Web, 知识图谱的概念最早由 Google 公司提出以表达其升级的搜索引擎技术, 如今知识图谱概念已经被用来泛指各类包含实体与丰富关系的知识

库, 被广泛用于存储人工智能任务的结构化语义信息。过去几年中, 知识图谱在人工智能应用中具有巨大潜力, 受到了广泛的关注。知识图谱的实例通常以三元组的形式进行存储, 将实体表示为有向图中代表属性或概念信息的节点, 关系表示为两实体之间具有实际语义的边, 诸如(中国, 首都, 北京)的三元组形式。

尽管知识图谱已从现实世界中提取了包含数百万个实体和数十亿个关系事实, 但大型知识图谱中的数据仍然稀疏不完整^[3]。例如, 在开放知识图谱 Freebase^[4]中, 约有 71% 的人缺少出生地信息, 99% 的没有民族信息^[5]; DBpedia^[6]中有 58% 的科学家实体没有指出其相关的主要贡献。随着知识图谱中知识实例的高速增长, 知识的表示形式以及之间的关联也变得更加复杂化、异质化。因此, 研究人员需将缺失的实例添加到知识库中以扩大其覆盖范围, 操作耗时耗力且人工成本较高。此外, 传统三元组的符号表

收稿日期: 2022-05-15; 修回日期: 2022-06-01。

基金项目: 国家自然科学基金(61373112, 51878536); 陕西省自然科学基金(2020JQ-687); 陕西省住房城乡建设科技计划项目(2020-K09)。

作者简介: 王 瑞(1996-), 女, 河南驻马店人, 硕士研究生, 主要从事模式识别、知识图谱等方向的研究。

李智杰(1980-), 男, 河南荥阳人, 博士, 副教授, 硕士生导师, 主要从事模式识别、数字建筑等方向的研究。

李昌华(1963-), 男, 宁夏人, 博士, 教授, 博士生导师, 主要从事图形图像处理、模式识别、数字建筑等方向的研究。

引用格式: 王 瑞, 李智杰, 李昌华, 等. 面向链接预测的知识图谱嵌入研究综述[J]. 计算机测量与控制, 2022, 30(9): 8-16.

示还面临着计算效率低和数据稀疏等问题^[7],导致其在大规模知识图谱的使用具有局限性,限制了知识图谱的发展,为知识图谱的表示带来了挑战。

在本文中,通过对知识图谱链接预测相关知识介绍,同时对链接预测模型框架进行了分析,并且列出了当前典型的应用场景,从而系统全面的对面向链接预测的知识图谱嵌入模型做了综述。

1 知识图谱链接预测概述

1.1 知识图谱嵌入

受当前技术的制约以及网络数据的繁杂冗余,在大型知识图谱中,需不断向知识库中补充新的实体和关系,导致研究人员的工作量剧增。此外,知识图谱中信息的缺失限制了知识图谱的使用,影响了知识图谱在推理和检索应用时的准确率。由于不能直接对三元组进行操作,需要为知识图谱中的实体和关系找到更好的表示形式。早期时候,使用符号三元组数据进行统计关系学习。但是这些方法既不具有良好的泛化性能,也不适用于大规模的知识图谱。因此,引入了知识图谱嵌入技术。嵌入是根据代表真实世界的数据集中相应元素的发生方式和彼此之间的相互作用自动学习的。同时,嵌入可用于表示任何种类元素的数值向量,将实体与关系向量化可在向量空间中通过数值计算挖掘出潜在的三元组信息及语义知识。此外,当嵌入作为一种类型的先验知识辅助时,可对神经网络的训练过程加以约束和监督^[8]。知识图谱是由实体和关系组成的复杂图结构,知识图谱嵌入是有向图的矢量表示,利用知识图谱嵌入操作来高效计算实体与关系的语义联系,提高了模型推理的准确率,同时也保留了知识图谱的固有结构,体现了原始图的语义,可用于识别其中的新链接,从而解决了链接预测任务。

伴随着知识图谱日新月异的发展,一系列的知识图谱嵌入模型被学者们相继提出。通过从知识图谱包含的关系信息中学习低维连续空间中的嵌入操作,将实体和关系表示为低维度的带有结构信息与语义信息的实值特征向量^[9],捕获了实体和关系的连接属性,为知识图谱提供数值计算框架,同时使其固有结构得以保留。如图1所示,知识图谱嵌入实质上就是通过优化基于边距的损失函数,其中边距是一个非负数,用于将正负三元组分开。将实体表示为空间中的向量,并通过距离来量化实体对象之间的相似性,关系通常被视为向量空间中的运算,获得具有某些明确定义的目标函数的三元组,即 (h, r, t) 的矢量表示。此外,关系也可以表示矩阵、张量、高斯分布以及多元高斯分布。训练知识图谱嵌入模型是为了找到模型的最佳参数从而进行最佳的嵌入,通过优化算法来迭代更新实体和关系的表示。在迭代更新过程中,通过一定的负采样策略替换正三元组的头或尾实体,从而生成负例三元组。优化过程旨在最大化肯定事实的合理性以及最小化否定事实的合理性。

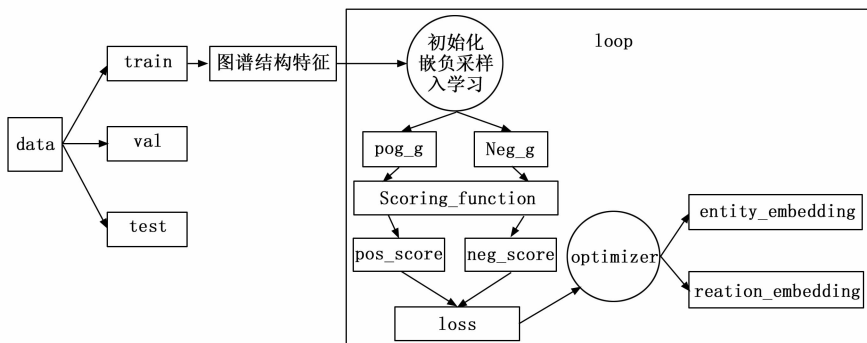


图1 知识图谱嵌入技术

知识图谱嵌入实现了对实体和关系的分布式表示,可高效地实现语义相似度计算等操作显著提升计算效率。同时,在低维实值向量空间中,可以度量任意对象之间的语义相似程度以及提高低频对象的语义表示的精确性^[10],实现异质知识对象之间的语义关联计算,有效缓解数据稀疏问题,实现异质信息融合。

1.2 链接预测

链接预测 (Link Prediction) 也称为知识图谱补全 (Knowledge Graph Completion), 利用评分函数计算并对候选实体或关系进行排序,旨在根据知识图谱中现有实体与关系推理出缺失的实体或关系。链接预测根据任务的不同,可分为头实体预测、尾实体预测和关系预测三种类型。例如,给定三元组实例 (h, r, t) ,首先利用嵌入模型学习实体与关系的向量特征;其次通过负采样策略破坏三元组中的任一实体或者关系生成知识图谱数据集中所没有的三元组 (h', r, t) 、 (h, r, t') 以及 (h, r', t) ;最后利用评分函数对其进行对应的评分 $f_r(h, t)$,并将所有实体进行由低到高的排序,输出最可能的实体或关系列表。这样可得到所有实体的排名,利用评估指标从而获得模型性能的评估。

链接预测是知识图谱嵌入的应用之一,是对存在于多对象总体中每个对象之间的相互作用及相互依赖关系推断的过程。链接预测旨在预测图谱中任意两个实体之间的关系以及实体间已存在关系的正确性,是对现有知识进行整合过滤以及筛选,进行更精准的知识发现,从而提高知识库中实例的质量,解决知识图谱中数据缺失不完整问题。既增加了下游应用的多样性,又可以作为预训练,利用实体与关系的表征向量支撑下游向量,为下游模型提供语义支持^[11]。即如图2所示,左侧图中的实线代表的是现有关系,虚线代表可能的关系,通过链接预测任务可计算出右侧图中不同颜色所代表的各种可能的关系。此外,在不同的链接预测任务中往往被赋予不同的功能,例如:在社交网络中链接预测被用于对用户或商品进行推荐;在生物学领域,被用于相互作用的发现;在知识图谱中被用于实体与关系的学习;在基础研究中,被用于图谱结构捕捉。链接预测任务是当前知识图谱嵌入模型研究的重点,面向链

接预测的知识图谱嵌入模型研究能够显著提升模型计算效率及性能,使知识获取、融合和推理的性能得到显著提升。对于基于知识图谱的人工智能应用等方面具有十分重要的意义,值得深入研究。

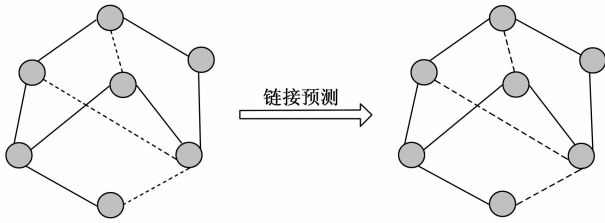


图 2 链接预测示例

1.3 知识图谱链接预测研究现状

为解决链接预测问题,已经提出了各种技术,包括基于翻译的方法、基于语义匹配的方法和基于神经网络的方法^[12]。其中,学习实体与关系的语义表示的知识图谱嵌入模型在当前研究中占有重要位置。基于此,本文从基于三元组结构信息和融合外部信息两个角度重点对面向链接预测的知识图谱嵌入模型进行了全面的综述。

1.3.1 基于三元组结构信息的知识图谱链接预测

目前绝大多数链接预测模型仅基于知识图谱中原始的实体与关系来推断新的事实。翻译模型是基于能量函数的平移模型,通过计算三元组的能量函数值来判断其是否为正例,一般情况下,负例三元组的能量计算数值较高。TransE^[13]在训练过程中引入负样本,通过学习正负例样本挖掘满足模型假设的实体和关系向量,促使语义相近的实体或者关系在向量空间中互相靠近,语义不相近的主动远离。TransE模型简单高效,但不能有效的对复杂关系建模。基于此,学者们提出了利用超平面让同一实体在不同关系下表示不同的 TransH^[14]模型、利用实体向关系空间投影并引用了投影映射的关系矩阵使不同关系拥有不同语义空间的 TransR^[15]模型、利用实体与关系之间的相互作用构建与实体与关系相关投影矩阵的 TransD^[16]模型。TransE、TransH、TransR 和 TransD 模型均是通过映射转换学习实体与关系的多样性来计算同一实体的三元组分数,有效避免了模型的收敛问题。自 2013 年首次提出 TransE 以来,基于这一框架提出了诸如通过关系映射属性转换嵌入的 TransM^[17]模型、通过更换损失函数中的度量函数为每一维的学习设置不同权重以实现自适应转换嵌入的 TransA^[18]模型等几十种基于不同架构的新模型。在最近的链接预测技术中,面向链接预测的知识图谱嵌入模型在一些基准测试中取得了很好的性能。

1.3.2 融合外部信息的知识图谱链接预测

基于三元组结构信息的知识图谱嵌入方法在一定程度上解决了当前主要问题,但是也仅仅考虑了知识图谱中的单个三元组同时假设三元组相互独立并对其单独建模。除了三元组本身的结构信息之外,知识图谱中往往还包括关系路径、实体描述、属性信息及实体类型等丰富的额外信

息,整合这些多源信息能够挖掘图谱底部更深层次语义信息,进一步提高模型的语义表示能力,从而实现更好的知识推理。

近年来,不少学者们还利用互联网语料库信息与三元组结构信息进行融合的知识表示学习,从而更好的实现开放式知识图谱的补全任务^[19]。Lin 等人^[20]提出了基于图谱自身结构信息的 PTransE 模型,在 TransE 模型的基础上加入路径信息,并使用路径约束资源算法来度量关系路径的置信度。其考虑了实体间多步间接路径的语义关系,将关系路径集成到学习过程中,在模型实验测试时取得很好的表现。在考虑实体描述信息方面,Xie 等人^[21]在模型训练时加入了实体描述信息,并将其与三元组结构信息进行联合建模,提出了基于实体描述的语义向量提出了 DKRL 模型;Xu 等人^[22]引入注意力机制并提出了联合学习模型,使实体在不同关系下表现出不同的语义向量;Gupta 等人^[23]提出了基于开放世界知识图谱的 CaRe 模型,通过学习实体邻域丰富的表示形式来捕获关系邻域的语义相似性;Shi 等人^[24]提出了使用依赖关系的内容屏蔽策略的 Con Mask,旨在从实体的文本信息中提取出与关系相关的语义信息;Wu 等人^[25]通过将数字属性预测损失添加到关系损失来扩展 TransE;An 等人^[26]提出了基于文本增强的知识表示学习模型,旨在处理三元组信息之间存在的歧义问题。此外,诸如 ConvE^[27]、ConvKB^[28]、HYPER^[29]、CompGCN^[30]、SACN^[31]和 CNN-BiLSTM^[32]等神经网络模型综合考虑了实体或关系的类型、时间信息、路径信息和子结构信息,同时卷积神经网络或注意力机制的使用也有助于产生更好的嵌入。

2 知识图谱链接预测框架分析

作为当前知识图谱方面研究热点的知识推理研究领域,受益于机器学习和深度学习技术的爆炸式增长,用于评价模型效果的链接预测更是成为衡量知识图谱表示模型效果最广泛使用的任务。链接预测是根据知识图谱中已存在的实体,通过对实体与关系的学习,并与知识库中对应实体或者关系进行链接从而实现知识库的补全^[33]。其本质思想是通过空间中已知的节点属性和不完全的链接来分析拓扑结构中存在的相似性,估计测试对象之间是否存在相应的链接^[34]。在过去几年中,作为学术界研究热点的知识图谱嵌入模型不断有新的研究成果产出,学者们也相继提出了基于不同方法的知识表示模型。本节先是按照时间线的前后简述了知识图谱嵌入模型分类,接着依据知识图谱建模过程是否有补充信息的加入,将翻译模型划分为仅基于三元组结构信息的知识图谱嵌入模型和融合外部信息的知识图谱嵌入模型,并对其进行详细介绍。

2.1 常用数据集与评价指标

知识图谱是基于大数据的,当前已经构建了许多开放的知识图谱,例如,Freebase、DBpedia、Yago^[35]和 NELL^[36-37]。它们通常包含大量使用数十亿实体和关系构建的事实,这些实体和关系分别表示为节点和链接这些节点的边。当前

在知识图谱链接预测领域主要使用如表 1 所示的数据集。

表 1 实验的数据集信息

数据集	关系数量	实体数量	训练集	验证集	测试集
WN18	18	40 943	141 442	5 000	5 000
WN11	11	38 696	112 581	2 609	10 544
WN18RR	11	40 943	86 835	3 034	3 134
FB13	13	75 043	316 232	5 908	23 733
FB15K	1 345	14 951	483 142	50 000	59 071
FB15K-237	237	14 541	272 115	17 535	20 466
YAGO10	37	123 000	1 000 000	5 000	5 000
NELL239	239	48 000	74 000	3 000	3 000

1) Freebase 是包含常见信息的世界知识, FB13、FB15K 和 FB15K-237 都是 Freebase 的子集。FB15K 中大约 70% 的三元组存在反向关系, 测试集中同样有 70% 左右的三元组, 在训练集中存在对应反向关系的三元组, 使得知识图谱表示模型可能倾向于学习反向关系^[38]; 其中, FB15K-237 是通过删除 FB15K 中训练集、测试以及验证集中的大量可逆关系数据创建得来的, 而且还过滤掉了所有琐碎的三元组, 确保训练集中连接的所有实体都没有直接连接到验证集或测试集中。其中, 15k 表示数据集中有 15k 个主题词, 237 表示共有 237 种关系。

2) WordNet 是覆盖范围比较广的英文语义知识库, 同时 WordNet 中的实体是具有不同概念的同义词, 关系表示同义实体之间的语义联系^[39]。WN11、WN18 和 WN18RR 都是 WordNet 的子集, 分别包含有 11 和 18 种关系。其由 WN18 删除可逆关系数据得到的子数据集, 消除了反向关系实例, 避免了表示任务中的信息泄露问题。

3) YAGO10: YAGO 数据集的子集, 主要包含关于人及其公民身份、性别和职业知识的信息。

4) NELL239: NELL 数据集的子集, 它包含有关人员、地点、团队、大学等实体类型的一般知识。

评价指标:

为了验证所提出的方法的性能, 通常在实验中设置“Raw”和“Filter”两种评价指标, 在“Raw”模式下生成的负样本不一定是实际意义上的错误三元组, 会扰乱排名, 降低 MR 指标, 故将其设置为“Filter”, 在排名之前用来过滤假的负例三元组。此外, 采用平均倒数排名 (Mean Reciprocal Rank, MRR)、平均排序 (Mean Rank, MR) 以及 Hits@k (k=1、3、10) 这三种通用的评价指标来衡量链接预测模型的性能。

1) MRR: 将测试集所有排名的倒数求均值, 即

$$MRR = \frac{1}{|T_{test}|} \sum_{(h,r,t) \in T_{test}} \frac{1}{2} \left(\frac{1}{rank_{r,t}(h)} + \frac{1}{rank_{h,r}(t)} \right) \quad (1)$$

其中: $rank_{r,t}(h)$ 表示头实体的排序, 同理, $rank_{h,r}(t)$ 表示尾实体的排序。MRR 主要用于衡量正三元组的最高排名, 第一个样本的贡献最大而且 MRR 具有平滑性, 受异常值的影响更小。MRR 的取值范围为 $MRR \in (0, 1)$,

计算值越大, 表示模型的链接预测性能越好。

2) MR: 指在得到的排序中对正确答案的实体排名求平均, 即

$$MR = \frac{1}{|T_{test}|} \sum_{(h,r,t) \in T_{test}} \frac{1}{2} (rank_{r,t}(h) + rank_{h,r}(t)) \quad (2)$$

MR 数值越小, 说明本模型在该任务上的模型性能越好。

3) Hits@k: 计算排名在前 k 位的正确实体所占的比例, 然后再对其求均值, 即

$$Hit@k = \frac{1}{|T_{test}|} \sum_{(h,r,t) \in T_{test}} \frac{1}{2} (|\{(h,r,t) | rank_{r,t}(h) \leq k\}| + |\{(h,r,t) | rank_{h,r}(t) \leq k\}|) \quad (3)$$

Hits@k 侧重于总体排名, 数值越大, 表示模型的链接预测性能越好。其中, K 的取值一般为 1、3 和 10。

2.2 知识图谱嵌入模型分类

伴随着知识图谱日新月异的发展, 一系列的知识图谱嵌入模型被学者们相继提出。一般情况下, 基于翻译模型的嵌入学习过程主要有三个步骤: 首先定义知识图谱中实体 $e \in E$ 和关系 $r \in R$ 在连续向量空间中的表示形式, 将实体表示为向量空间中带有结构信息与语义信息的特征向量, 关系表示为向量空间中实体间的翻译运算, 通常由随机初始化来获得实体和关系的嵌入向量; 其次定义三元组 (h, r, t) 的评分函数 $f_r(h, t)$, 根据嵌入向量 h 和 t 来评估任意事实三元组 (h, r, t) 在空间中成立的可能性, 得分越高表明事实成立的可能性越大; 最后通过优化算法来迭代更新实体和关系的表示。在迭代更新过程中, 通过一定的负采样策略替换正三元组的头或尾实体, 从而生成负例三元组。优化过程旨在最大限度提升真实事实的可能性, 同时降低无效事实的可能性。

由表 2 所示, 按照时间轴展示了知识图谱嵌入模型近几年的发展。同时, 在表 3 中总结了面向链接预测的知识图谱嵌入模型的优缺点。

表 2 知识图谱嵌入模型

	距离模型	翻译模型	语义匹配模型	神经网络模型	几何模型
2013	UM、SE	TransE	RESCAL LFM	NTN SLM	—
2014	—	TransH TransM	SME TATEC TRESICAL	MLP	—
2015	—	TransR TransD TransA	DisMult	—	—
2016	—	TranSparse TransF TransG	HolE ComplEx	NAM	—
2017	—	ITransF	ANALOGY	ProjE	Poincare
2018	—	TransAt	Simple HolEx	ConvE ConvKB	RotatE TorusE
2019	—	—	TuckER CrossE	CapsE ConvR	QuatE MuRP
2020	—	—	—	—	HAKE

表 3 链接预测模型优缺点总结

模型	优点	缺点
距离模型	假设简单	没有考虑数据的多关系问题,链接预测性能较差
翻译模型	参数量小,可应用到大型知识图谱的链接预测任务上;实验验证表明:TransE 配置适当时,能够超过多数公布的最先进结果;在嵌入维度设置较低的情况下,表现突出	TransE 模型存在不能建模一对多、多对一、多对多关系范式的缺陷,但可以通过适当的假设来弥补,如 TransH、TransR 等模型
语义匹配模型	参数量小;缩放方便;可处理对称性和自反性关系且大多数语义匹配模型都可完全表达。实验测试结果先进,并且在大量不同的配置中具有鲁棒性	完全表达只有在嵌入维度设置较大时才能凸显出来;不具有适用大规模知识图谱链接预测任务的能力
神经网络模型	能够充分利用知识图谱的物理结构,尤其围绕三元组的上下文信息,在复杂关系的学习中具有很强的表达能力;在嵌入维度设置较低时,具有最佳的性能。ConvE 模型在大量测试中的性能表现稳定且突出	参数多,存储空间和模型训练时间长;且不适用于建模实体的层次性、关系的多样性,可解释性差;存在不恰当的评估协议问题。
几何模型	可利用关系组件的混合对关系的潜在语义进行嵌入;RotatE 表现了目前公布的最先进链接预测结果	变相提升嵌入维度,模型假设复杂,可扩展性差;需要较长的训练时间和预测时间

2.3 翻译模型

翻译模型通常使用基于距离的评分函数,将三元组的合理性视为向量空间中两个实体节点间的距离。翻译模型本质上属于距离模型,同样是利用距离的评分函数来衡量事实成立的可能性。但相较于距离模型,翻译模型最大不同点是将关系建模为头实体到尾实体的翻译向量。

基于三元组的模型只关注实体与实体之间的一跳关系,依据知识图谱本身的结构化信息从三元组的视角对实体和实体之间的关系进行建模,认为不同事实三元组 (h, r, t) 之间相互独立。通常情况下很少考虑实体与关系的语义信息,即利用图谱的自身结构将每个关系解释为潜在空间中的平移,并将实体和关系表示为相同长度的一维向量。

TransE 模型是受 Word2Vec^[40] 启发所提出的第一个基于距离的模型,同时也是平移距离模型中最具代表性的模型。为有效捕获知识图谱的结构信息,将实体和关系表示为相同语义空间中的向量形式,使得嵌入的实体 h 和 t 可以通过 r 以低误差连接,即当三元组 (h, r, t) 成立时,有 $h+r \approx t$ 。TransE 参数简单训练效率高,但在处理 $N-1$ 、 $1-N$ 、 $N-N$ 等复杂关系上存在着一些缺陷,缺乏对各种关系的区分策略,可能会出现不同实体有着同样的含义。例如,(中国,首都,北京)和(英国,首都,伦敦)根据翻译原则在嵌入空间中会出现中国-首都=英国-首都这样

的情况,但很显然北京不等于伦敦。为了解决 TransE 不能很好的处理多关系实体的这一缺陷,学者们提出了一些基于 TransE 的变体,例如 TransH、TransR 等模型。表 4 中给出了 TransE、TransH、TransR 的得分函数以及参数空间类型,同时在图 3 中给出了具体的图示。

表 4 纯翻译模型相关信息

模型	得分函数 $f_r(h, t)$	参数
TransE	$-\ h+r-t\ _{1/2}$	$h, r, t \in R^k$
TransH	$\ (h-hw_r^T)+r-(t-tw_r^T)\ _{L1/L2}$	$h, t \in R^d, r, w_r \in R^d$
TransR	$-\ M_r h+r-M_r t\ _2$	$h, t \in R^d, r \in R^k, M_r \in R^{k \times d}$

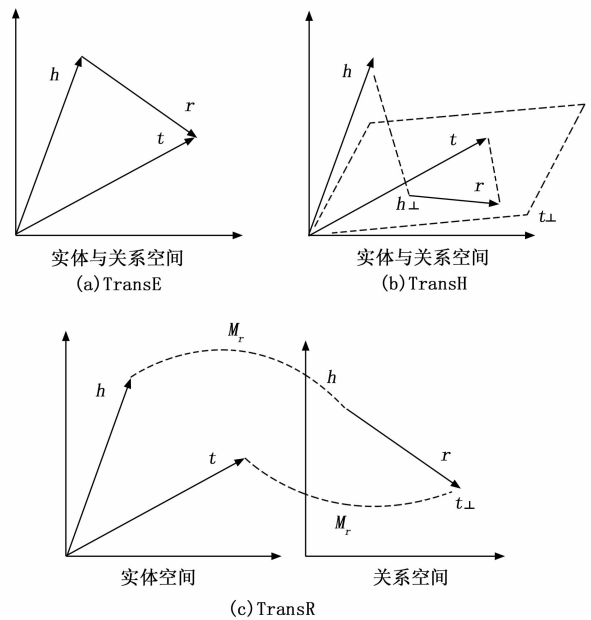


图 3 TransE、TransH、TransR 模型的嵌入

TransH 通过构建特定于关系的超平面,使得实体在连接不同语义的关系时拥有不同的表示形式。TransH 是在 TransE 模型的基础上加以改进的,TransH 假设将关系 r 都建模到以 w_r 为法向量的超平面上。也就是说在给定事实三元组 (h, r, t) 的情况下,先将给定实体对 (h, t) 的嵌入向量投影到关系 r 的超平面上,表示形式记作 h_{\perp} 和 t_{\perp} ,其中 $h_{\perp} = h - w_r^T h w_r$ 、 $t_{\perp} = t - w_r^T t w_r$ 。若三元组 (h, r, t) 成立,则有 $h_{\perp} + r \approx t_{\perp}$,实体的投影向量在超平面上通过 r 以低误差连接。TransH 可以使同一实体在不同关系中发挥不同的作用。

TransR 引入投影映射的特定关系矩阵来解决 TransH 不能区分不同实体类型的映射的问题,同时挖掘更多的语义信息。在 TransR 模型中,同一关系中实体的映射相同,同时不同的关系拥有不同的语义空间。将关系建模为实体空间到关系空间的投影矩阵,通过翻译操作建立实体向量与关系向量的联系。对于关系 r ,除本身结构中的关系语义,还使用映射矩阵 M_r 来描述关系所处的关系空间。假定

实体对 (h, r, t) , $h, t \in R^k$, $r \in R^d$, $k \neq d$ 。将头尾实体映射到对应的关系空间中形成向量 h_r, t_r , 然后在关系空间中对投影后的实体建模 $h_r + r \approx t_r$, 实现翻译操作, 其中 $h_r = hM_r, t_r = tM_r, M$ 为关系 r 对应的映射矩阵且 $M \in R^{k \times d}$, $\|h\|_2 \leq 1, \|r\|_2 \leq 1, \|t\|_2 \leq 1, \|hM_r\|_2 \leq 1, \|tM_r\|_2 \leq 1$ 。TransR 模型中引入的投影矩阵只与关系 r 有关, 忽略了实体与关系的交互。此外, 由于 TransR 为每个关系嵌入一个转换矩阵, 导致得分函数计算复杂度增高, 训练时间较长。

2.4 采样方法

负采样是在训练时从未观察到的三元组数据中抽取负例三元组, 也是知识图谱嵌入过程中的重要步骤。为了提高空间效率, 一般情况下知识图谱中只存储正样本而不存储负样本, 所以在模型训练期间, 向模型提供负样本是至关重要的。如果该模型只在真实样本上进行训练, 那么它可以通过简单地返回任何事实的大分数来将所有损失降至最低, 但这失去了模型训练的初衷。在知识图谱嵌入过程中, 否定事实的生成通常是通过负采样来完成的, 利用负采样来最小化边缘的排序损失, 同时也体现了知识图谱嵌入模型的性能在很大程度上取决于负采样的质量。直观地说, 利用负样本在嵌入空间中引入排斥力, 使事实三元组中不可互换的实体在嵌入时彼此远离。因此, 必须选择尽可能的训练生成高质量的负样本。随着训练的进行, 为模型提供越来越接近真实事实的负样本, 学习有效的表示方法, 以便更好地调整实体向量与关系向量的嵌入。

2.4.1 随机采样

随机采样是一种传统的负采样方法, 旨在从均匀分布中随机的选择实体替换事实三元组的头部或尾部实体生成负面事实。由于被采样的实体可能与被替换实体和目标关系完全无关, 所以生成的大多数负面事实很容易与正面事实区分开来, 未被充分训练的反例又很难被选择, 导致随机生成的负例三元组质量会很差, 有时也随之会出现“零损失”问题^[41]: 当生成的负例三元组质量较低时, 模型的评分函数会给予其较低的分值, 这将出现正、负三元组分值的差大于设置的边界值的情况, 随之的损失值也将为零。此时模型不会对实体向量与关系向量进行更新操作, 即模型在无效学习, 也就不能学习到更多的样本特征, 导致模型的训练程度评估出现偏差。如图 4 所示, 在训练初期时, 随机采样是非常有效的, 此时正、负例三元组在同一裕度内。随着随机采样训练的进行, 即对图中蓝色圆中的三元组进行采样, 此时这些三元组对于模型训练毫无意义。这是因为这些三元组超出了边界不在同一裕度内, 也就不会给模型带来任何的损失甚至减慢了模型收敛的速度。因此, 在边界内忽略一定数量的负三元组 (如黑色虚线圆圈所示) 可提高模型训练效率。

2.4.2 过滤采样

过滤采样是基于随机采样的一种采样方法, 只是在随机采样的过程中加入了过滤机制。通常情况下随机采样会出现假阴性负例三元组样本, 即有可能为正例三元组或者

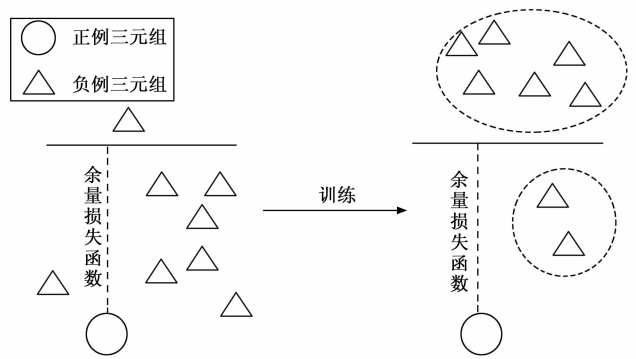


图 4 模型训练零损失状态

在数据集中曾出现过的三元组。当一些损坏的三元组最终成为有效的三元组时, 很明显这会影响到模型的表征能力与性能。在这种情况下, 当对所有三元组打分排名时, 会出现假阴性样本排在测试三元组之上的情况, 这并不是因为模型学习效果不好, 因为此时两个三元组都为真实实例。为了避免这种误导行为, 在排名之前, 过滤采样会从损坏的三元组列表中删除曾出现在实验数据集中的所有三元组, 保证所有损坏的三元组全部为真正的负样本。

2.4.3 伯努利采样

同样的, 伯努利采样也是基于随机采样的一种采样方法。对于给定三元组 (h, r, t) , 通过随机采样获得负三元组 (h', r, t') 。然而, 由于知识图谱不完善, 随机采样会将假阴性三元组引入训练。伯努利采样根据关系的映射属性, 设置不同的概率去替换头或尾部实体从而破坏原本数据集中的三元组。若是一对多关系时, 模型会给予更大的概率替换头部实体; 反之, 模型给予更大的概率替换尾部实体。通过这种方式, 减少了产生假阴性三元组的机会。具体地说, 在关系 r 的所有三元组中, 首先会计算每个头、尾实体的平均尾部实体数 tph 和 hph , 然后将伯努利分布用于负采样: 给定一个正三元组 (h, r, t) , 若概率为 $\frac{tph}{tph + hpt}$, 则采用替换头部实体的方式来破坏三元组; 若概率 $\frac{hph}{tph + hpt}$, 则采用替换尾部实体的方式来破坏三元组^[42]。

2.4.4 对抗生成采样

受生成对抗性深度模型^[43]的启发, 提出了对抗生成采样^[44]这一对抗学习框架, 其提供了对动态负样本分布进行建模的采样策略, 旨在提高模型训练时负例三元组的质量。将基于不同损失函数的嵌入模型作为生成器和鉴别器, 分别用来生成高质量的负例三元组和训练具有高表征能力的模型。如图 5 所示, 发生器用于训练原始模型, 随后通过基于概率的对数似然损失函数的生成器最大化鉴别器对其动作的响应, 动态地估计负样本分布, 通过高质量的负例三元组来改进知识图谱嵌入模型。对候选三元组上的概率分布进行计算采样, 并通过源于强化学习的策略梯度最小化生成的负例三元组的得分。基于距离的边缘损失函数的鉴别器将接收到正负样本三元组加以区分, 并采用优化函数来最小化边缘损

失。通过对分数较大的负例三元组进行采样，避免了梯度消失的问题。整个模型框架通过不断地训练模型，最终产生一个更好的鉴别器，从而获得更好的性能。

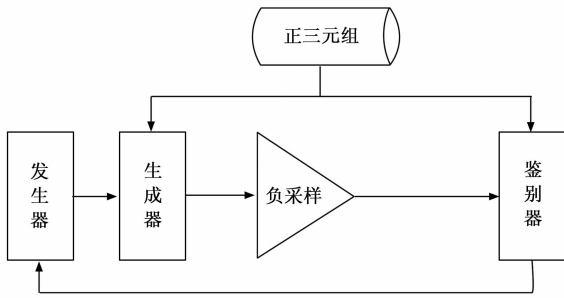


图 5 对抗生成采样框架

3 典型智能应用场景

知识图谱技术最早被 Google 公司提出并应用到其搜索引擎技术中，从而使搜索引擎具备了查询理解的能力。从字面匹配到概念理解，可更好的理解用户的真实想法为用户服务，让用户获得与搜索关键字最相关的词条链接以及获得与关键字更加智能化的信息，返回用户最希望的结果。如图 6 所示，当在搜索引擎中搜索《西游记》作者时，搜索引擎会将查询关键字理解现实世界中的概念和事物，然后搜索引擎根据“《西游记》”，“作者”两个实体来理解用户的意图，同时返回问题的答案和与搜索实体相关的其他实体。



图 6 百度搜索界面

人工智能的卓越发展使得知识图谱向量化表示得到了快速的发展。相较于传统 one-hot 编码的大维度、编码稀疏，无法体现实体间关系的远近程度，而嵌入技术可将实体和关系表示为向量的形式，更利于各种推理计算，同时节省了空间与模型训练时间。知识图谱在知识推理以及多源异质知识的整合提取方面显得尤为重要，通过学习知识图谱中已有事实三元组实体之间的语义关联进而推理出新的事实并将其添加到图谱中，促进了人工智能及其应用的发展^[45]。

如图 7 所示，通过相似实体在同一空间中相互靠近的

原则，只需要分析 Adam Ant 周围的实体便可推知他的职业以及其他的一些信息。即，在 Adam Ant 的周围相近的实体都是与音乐有关联的实体，则可推理出此人的职业必定与音乐有关。此外，为下游关系抽取、智能问答、信息检索、个性化智能推荐等任务发挥了必不可少的枢纽作用。例如，Apple 的 Siri、百度的小度、微软的 Cortana 等智能聊天机器人可以处理客户的请求或为用户提供帮助。从而帮助用户推荐附近的餐厅，回答简单的事实问题，或者管理日历活动等一系列日常任务。

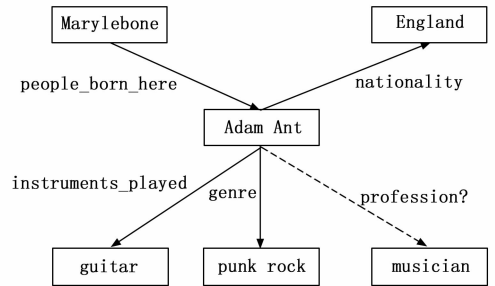


图 7 FB15K 中实体及其邻居节点

由表 5 所示，从智能问答、系统推荐、信息检索以及医药应用四个方面总结了当前知识图谱嵌入技术的典型应用案例^[46-60]。

表 5 知识图谱嵌入模型应用案例总结

应用场景	主要结论
智能问答	改造记忆网络，完成简单问答；提出条件聚焦神经网络，减少搜索空间；生成用户有好的自然答案；实现问题与知识图谱的双向交互；构造变分推理网络、认知图、稀疏表示回答多跳提问
系统推荐	编码用户一项图的高阶连通性，感知时间敏感和话题敏感的新闻，通过交互式会话向用户推荐高质量的项目
信息检索	使用关系记忆网络编码潜在依赖，提出显式语义排序，提升学术搜索引擎的理解能力
医药应用	通过编码、检索和重新解释提升了疾病分类准确性，生成稳健的医疗报告；构建知识驱动的分类器评估自杀风险

4 结束语

在近十年间，知识表示学习有了很大的发展，同时也提出了许多基于知识表示学习的方法。本文介绍了知识图谱的概念性知识，包括系统地讨论了知识图谱链接预测的研究现状、框架分析以及当前典型的应用场景。面向链接预测的知识图谱嵌入模型旨在提高知识图谱链接预测准确率，增强嵌入模型的表达性。同时，大规模知识图谱具有重要的人工智能应用前景。例如，在军事应用方面构建军用无人系统领域故障知识图谱用以智能搜索以及辅助决策；在目标检测控制系统中引入知识图谱用以多目标的关联判别；在航空航天方面，利用知识图谱设计雷达场景识别系

统用以空间目标的场景识别。在未来研究中, 应注重对面向链接预测的知识图谱嵌入模型的研究, 更好的进行大规模知识图谱补全, 从而促进人工智能应用的发展。

参考文献:

- [1] 杨晓晖, 孙 莹. 基于知识图谱的社交网络用户行为研究进展 [J]. 河北大学学报 (自然科学版), 2021, 41 (1): 77-86.
- [2] 李 刚, 李银强, 王洪涛, 等. 电力设备健康管理知识图谱: 基本概念、关键技术及研究进展 [J]. 电力系统自动化, 2022, 46 (3): 1-13.
- [3] 张 翔, 杨伟杰, 刘文文, 等. 人工智能时代知识图谱表示学习方法体系 [J]. 科技导报, 2021, 39 (22): 94-110.
- [4] BOLLACKER K D, EVANS C, PARITOSH P, et al. Freebase: A Collaboratively Created Graph Database for Structuring Human Knowledge [C] //Proceedings of the SIGMOD, Vancouver, BC, Canada, 2008 (8): 1247-1250.
- [5] BOLLACKER K, EVANS C, PARITOSH P, et al. Freebase: a collaboratively created graph database for structuring human knowledge [C] //Proceedings of the 2008 ACM SIGMOD international conference on Management of data. 2008: 1247-1250.
- [6] LEHMANN J, ISELE R, JAKOB M, et al. DBpedia: a large-scale, multilingual knowledge base extracted from Wikipedia [J]. Semantic Web, 2015, 6 (2): 167-195.
- [7] 陈晓军, 向 阳. STransH: 一种改进的基于翻译模型的知识表示模型 [J]. 计算机科学, 2019, 46 (9): 184-189.
- [8] 雷 涛, 杨亚宁, 唐佳璐, 等. 雷达场景识别系统设计 [J]. 计算机测量与控制, 2021, 29 (10): 158-163.
- [9] 石 川, 王睿嘉, 王 啸. 异质信息网络分析与应用综述 [J]. 软件学报, 2022, 33 (2): 598-621.
- [10] 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.
- [11] 张 翔, 杨伟杰, 刘文文, 等. 人工智能时代知识图谱表示学习方法体系 [J]. 科技导报, 2021, 39 (22): 94-110.
- [12] MA J, QIAO Y, HU G, et al. ELPKG: a high-accuracy link prediction approach for knowledge graph completion [J]. Symmetry 2019, 11 (9), 1096-1097.
- [13] BORDES A, USUNIER N, GARCIA-DURAN A, et al. Translating embeddings for modeling multi-relational data [C] //Proceedings of the 27th Annual Conference on Neural Information Processing Systems, Lake Tahoe. Red Hook: Curran Associates, 2013 (12): 2787-2795.
- [14] WANG Z, ZHANG J, FENG J, et al. Knowledge graph embedding by translating on hyperplanes [C] //Proceedings of the 28th AAAI Conference on Artificial Intelligence, Menlo Park, 2014 (7): 1112-1119.
- [15] LIN Y, LIU Z, SUN M, et al. Learning entity and relation embeddings for knowledge graph completion [C] //Proceedings of the 29th AAAI Conference on Artificial Intelligence, Austin, Menlo Park, 2015 (1): 2181-2187.
- [16] JI G, HE S, XU L, et al. Knowledge graph embedding via dynamic mapping matrix [C] //Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing, 2015: 687-696.
- [17] FAN M, ZHOU Q, CHANG E, et al. Transition-based knowledge graph embedding with relational mapping properties [C] //Proceedings of the 28th Pacific Asia Conference on Language, Information and Computation, Phuket, 2014 (12): 328-337.
- [18] XIAO H, HUANG M, YU H, et al. TransA: an adaptive approach for knowledge graph embedding [Z]. computer science, 2015.
- [19] 朱丽雅, 张 珺, 洪 亮, 等. 数字人文领域的知识图谱: 研究进展与未来趋势 [J]. 知识管理论坛, 2022, 7 (1): 87-100.
- [20] LIN Y, LIU Z, LUAN H, et al. Modeling relation paths for representation learning of knowledge bases [C] //Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, 2015: 705-714.
- [21] XIE R, LIU Z, JIA J, et al. Representation learning of knowledge graphs with entity descriptions [C] //Proceedings of the AAAI Conference on Artificial Intelligence, Arizona, 2016, 1143-1147.
- [22] XU J, QIU X, CHEN K, et al. Knowledge graph representation with jointly structural and textual encoding [C] //Proceedings of the international joint conference on artificial intelligence, 2017: 1318-1324.
- [23] GUPTA S, KENKRE S, TALUKDAR P. CaRe: open knowledge graph embeddings [C] //Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, 2019: 378-388.
- [24] SHI B, WEINGER T. Open-world knowledge graph completion [C] //Proceedings of the 32th AAAI Conference on Artificial Intelligence, 2018: 1957-1964.
- [25] WU Y, WANG Z. Knowledge graph embedding with numeric attributes of entities [C] //Proceedings of the Third Workshop on Representation Learning for NLP, 2018: 132-136.
- [26] AN B, CHEN B, HAN X, et al. Accurate text-enhanced knowledge graph representation learning [C] //Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2018: 745-755.
- [27] DETTMERS T, MINERVINI P, STEMETORP P, et al. Convolutional 2d knowledge graph embeddings [C] //Proceedings of the AAAI Press: Palo Alto, CA, USA, 2017: 1811-1818.
- [28] NGUYEN D Q, NGUYEN T D, NGUYEN D Q, et al. A novel embedding model for knowledge base completion based on convolutional neural network [C] //Proceedings of the NAACL-HLT, New Orleans, 2018 (6): 327-333.
- [29] BALAZEVIC I, ALLEN C, HOSPEDALES T M. Hypernetwork knowledge graph embeddings [C] //Proceedings of the

- ICAN, Munich, Germany, 2019 (11): 553–565.
- [30] VASHISHTH S, SANYAL S, NITIN V, et al. Composition-based multi-relational graph convolutional networks [C] // Proceedings of the ICLR, Addis Ababa, Ethiopia, 2020 (4): 26–30.
- [31] SHANG C, TANG Y, HUANG J, et al. End-to-end structure-aware convolutional networks for knowledge base completion [C] // Proceedings of the AAAI Press, Palo Alto, 2019: 3060–3067.
- [32] JAGVARAL B, LEE W, ROH J S, et al. Path-based reasoning approach for knowledge graph completion using CNN-BiLSTM with attention mechanism [J]. Expert Systems with Applications, 2020: 142–143.
- [33] 冯子桓, 梁 循, 牛思敏. 大数据时代的社交网络舆情主题图谱研究 [J]. 电子科技大学学报 (社科版), 2022, 24 (2): 19–28.
- [34] 韩亚楠, 刘建伟, 罗雄麟. 概率主题模型综述 [J]. 计算机学报, 2021, 44 (6): 1095–1139.
- [35] MAHDISOLTANI F, BIEGA J A, SUCHANEK F M. YAGO3: a knowledge base from multilingual wikipedias [C] // Proceedings of the CIDR, Asilomar, CA, USA, 2015 (1): 4–7.
- [36] MITCHELL T, COHEN W, HRUSCHKA E, et al. Never-ending learning [J]. Communications of the ACM, 2018, 61 (5): 103–115.
- [37] GARDNER M, MITCHELL T M. Efficient and expressive knowledge base completion using subgraph feature extraction [C] // Proceedings of the 2015 Conference on Empirical Methods in Natural Language, Lisbon, Portugal, 2015 (11): 1488–1498.
- [38] 余传明, 张贞港, 孔令格. 面向链接预测的知识图谱表示模型对比研究 [J]. 数据分析与知识发现, 2021, 5 (11): 29–44.
- [39] 彭 敏, 黄 婷, 田 纲, 等. 聚合邻域信息的联合知识表示模型 [J]. 中文信息学报, 2021, 35 (5): 46–54.
- [40] MIKOLOV T, SUTSKEVER I, CHEN K, et al. Distributed representations of words and phrases and their compositionality [C] // Proceedings of the NIPS, Lake Tahoe, 2013: 3111–3119.
- [41] WANG P, LI S, PAN R. Incorporating GAN for negative sampling in knowledge representation learning [C] // Proceedings of the 32nd AAAI Conference on Artificial Intelligence, Louisiana, 2018: 2005–2012.
- [42] 方 阳, 赵 翔, 谭 真, 等. 一种改进的基于翻译的知识图谱表示方法 [J]. 计算机研究与发展, 2018, 55 (1): 12.
- [43] GOODFELLOW I, POUGET-ABADIE J, MIRZA M, et al. Generative adversarial nets [J]. In Advances in Neural Information Processing Systems, 2014, 2672–2680.
- [44] CAI L, WANG W Y. KBGAN: adversarial learning for knowledge graph embeddings [Z]. 2017.
- [45] 饶官军, 古天龙, 常 亮, 等. 基于相似性负采样的知识图谱嵌入 [J]. 智能系统学报, 2020, 15 (2): 218–226.
- [46] BORDES A, USUNIER N, CHOPRA S, et al. Large-scale simple question answering with memory networks [Z]. Computer Science, 2015.
- [47] HE S, LIU C, LIU K, et al. Generating natural answers by incorporating copying and retrieving mechanisms in sequence-to-sequence learning [C] // Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, 2017: 199–208.
- [48] CHEN Y, WU L, ZAKI M J. Bidirectional attentive memory networks for question answering over knowledge bases [C] // Proceedings of the NAACLHLT. 2019: 2913–2923.
- [49] DING M, ZHOU C, CHEN Q, et al. Cognitive graph for multi-hop reading comprehension at scale [C] // Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, 2019: 2694–2703.
- [50] COHEN, WILLIAM W, et al. Scalable neural methods for reasoning with a symbolic knowledge base [C] // Proceedings of the International Conference on Learning Representations, 2019.
- [51] WANG H, ZHANG F, ZHAO M, et al. Multitask feature learning for knowledge graph enhanced recommendation [C] // Proceedings of the International Conference on World Wide Web Conference, 2019: 2000–2010.
- [52] WANG X, HE X, CAO Y, et al. Kgat: Knowledge graph attention network for recommendation [C] // Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2019: 950–958.
- [53] WANG H, ZHANG F, XIE X, et al. DKN: Deep knowledge-aware network for news recommendation [C] // Proceedings of the 2018 world wide web conference, 2018: 1835–1844.
- [54] ZHOU K, ZHAO W X, Bian S, et al. Improving conversational recommender systems via knowledge graph based semantic fusion [C] // Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2020: 1006–1014.
- [55] DAI Q, T D, PHUNG D. A relational memory-based embedding model for triple classification and search personalization [C] // Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, 2020: 3429–3435.
- [56] XIONG C, POWER R, CALLAN J. Explicit semantic ranking for academic search via knowledge graph embedding [C] // Proceedings of the 26th international conference on world wide web, 2017: 1271–1279.
- [57] LI C Y, LIANG X, HU Z, et al. Knowledge driven encode, retrieve, paraphrase for medical image report generation [C] // Proceedings of the AAAI Conference on Artificial Intelligence. 2019, 33 (1): 6666–6673.
- [58] 喻凡坤, 胡超芳, 罗晓亮, 等. 无人系统故障知识图谱的构建方法及应用 [J]. 计算机测量与控制, 2020, 28 (10): 66–71.
- [59] 郭 策, 高跃清, 沈宇婷, 等. 基于知识学习的多目标关联检测与识别方法 [J]. 计算机测量与控制, 2021, 29 (11): 201–206.
- [60] GAUR M, ALAMBO A, SAIN J P, et al. Knowledge-aware assessment of severity of suicide risk for early intervention [C] // Proceedings of the World Wide Web Conference, 2019: 514–525.