

引用格式:梁杨,崔鑫,张骋,赵海英.宋代石刻纹样知识体系挖掘与知识图谱构建研究[J].中国传媒大学学报(自然科学版),2022,29(04):41-49.

文章编号:1673-4793(2022)04-0041-09

宋代石刻纹样知识体系挖掘与知识图谱构建研究

梁杨^{1*},崔鑫³,张骋³,赵海英²

(1. 四川泸县宋代石刻博物馆;2. 北京邮电大学人工智能学院,北京 100876;
3. 北京邮电大学计算机学院,北京 100876)

摘要:以泸县宋代石刻纹样为载体,挖掘宋代石刻纹样知识体系,并构建了宋代石刻纹样知识图谱。在划分泸县宋代石刻种类,解析35类可移动文物核心元数据及结合专家知识的基础上,引入纹样元数据和纹样知识体系,对泸县宋代石刻纹样进行了深度挖掘,并搭建了一套基于泸县宋代石刻纹样的知识图谱构建路径:知识图谱设计、标签化(关键词提取)、结构化(命名实体识别、关系提取)、知识图谱存储、可视化,为传统文化数据的标签化和结构化提供解决方案。

关键词:泸县宋代石刻;知识体系挖掘;知识图谱构建

中图分类号:TP391 文献标识码:A

Research on knowledge system mining and knowledge map construction of stone carving patterns in Song Dynasty

LIANG Yang^{1*}, CUI Xin³, ZHANG Cheng³, ZHAO Haiying²

(1. Song Dynasty Stone Carving Museum in Luxian County, Sichuan Province, Sichuan 646000, China;
2. Artificial Intelligence Institute, Beijing University of Post and Telecommunication, Beijing 100876, China;
3. Beijing University of Post and Telecommunication, Beijing 100876, China)

Abstract: This paper takes the Song Dynasty stone carving patterns in Luxian County as the carrier, excavates the knowledge system of Song Dynasty stone carving patterns, and constructs the knowledge map of Song Dynasty stone carving patterns. On the basis of the classification of the stone carvings of the Song Dynasty in Luxian County, the analysis of the core metadata of 35 types of movable cultural relics and the combination of expert knowledge, pattern metadata and pattern knowledge system are introduced to deeply excavate the stone carvings of the Song Dynasty in Luxian County. Furthermore, a set of knowledge graph construction paths based on the Song Dynasty stone carvings in Luxian County is built in this paper: knowledge graph design, tagging (keyword extraction), structuring (named entity recognition, relationship extraction), knowledge graph storage, visualization, which provide solutions for labeling and structuring traditional culture data.

Keywords: Song Dynasty stone carvings in Luxian County; knowledge graph excavation; knowledge graph construction

基金项目:揭榜挂帅重点研发课题(课题编号:2021YFF0901701);基于泸县宋代石刻文物艺术价值体系与知识图谱构建预研究(项目批准号:SCWW2021A03).

通讯作者(*为通讯作者):梁杨(1986-),女,泸县宋代石刻博物馆馆长,主要从事文物博物研究,文物保护技术研究,博物馆运营管理等。Email:372448042@qq.com.

1 引言

四川泸县宋代石刻博物馆位于四川省泸州市泸县玉蟾山麓温泉度假区。博物馆总建筑面积达1.2万平方米,展陈面积达5000平方米。现有馆藏文物14000余件,囊括石质、字画、玉器、陶器等34个种类。国家珍贵文物550件,其中一级文物120件,二级文物145件,三级文物285件,在全国尤其是西南地区的县区实属罕见。

泸县宋代石刻可分为四神、武士、伎乐、侍仆和其他类别。四神、人物、建筑、家具、以及各种花鸟灵兽等图案,雕刻细腻,精湛巧妙,独具特色^[1],充分反映了南宋时期西南的世俗生活和社会文化状况,是当时南宋社会生活的再现,也是一幅南宋社会历史画卷,为我们研究宋代历史提供了非常丰富的实物资料。

纹样是泸县宋代石刻中比较具有象征性的艺术特征之一。建筑构件设计纹样、门窗设计纹样、家居用品纹样、人物服饰和头饰纹样等,均受宋代的文化观念影响,纹饰特征逐渐与外来风格相脱离,形成了本土化的艺术特色。宋代石刻的植物形态、动物造型和各种相关纹饰都十分生动,并且纹样种类繁多,雕刻精美,蕴含两宋特有的社会文化气息。

泸县宋代石刻,作为南宋西南一段波澜壮阔历史的再现,是我们的宝贵财富,有大量的内容值得我们探究。本文构建了宋代石刻纹样知识体系,有助于博物馆人员梳理对应知识以及方便展览。其次,本文使用统计机器学习的方法对宋代石刻纹样进行标签化和结构化处理,并对结构化数据进行存储和可视化展示,探寻了石刻纹样知识图谱的构建技术路径。

本论文主要分为两个部分,一是泸县宋代石刻纹样知识体系挖掘,二是泸县宋代石刻纹样知识图谱构建。

2 泸县宋代石刻纹样知识体系挖掘

泸县宋代石刻可分为四神、武士、伎乐、侍仆和其他类别。四神是指中国古代的四灵:青龙、白虎、朱雀和玄武。泸县宋代石刻中,四神数量比较多,它们形态各异,造型别致,雕刻艺术精湛。四神在宋墓中的雕刻位置一般是:左青龙、右白虎、前朱雀、后玄武,它们代表方位、信仰、吉祥等意思。武士在石刻中主要充当镇墓俑的角色,有保障功能的象征意义。伎乐类包括器乐演奏、舞蹈、戏剧和乐官等,表现的是当时的文化生活娱乐等内容。侍仆是指为主人服务的

侍仆,包括女侍和男侍侍仆类,其中女侍造型各异,生动逼真。比侍仆地位稍高的是仆从。其他类包括宗教类的飞天、人物故事、动植物、花卉和族谱等内容。动植物一般为门框装饰,多为压地隐起的浅浮雕,涉及动物的图案有狮、虎、羊、鹿、凤凰等,植物图案有荷花、牡丹、菊花、芙蓉、月季、水仙等;还有石质族谱与墓志铭。墓志铭出土有多件,其中一件是族谱与墓志铭的结合。族谱记载了多方面的内容,其中还涉及到道教,是研究道教等方面的重要物证,非常珍贵。

泸县宋代石刻中的纹样主要来自四神类石刻、其他类石刻。四神类石刻中的青龙石刻有龙纹,白虎石刻有虎纹。其他类石刻中有动物纹、植物纹,如精美绝伦的卷草花、莲瓣纹等^[2]。这些纹样种类繁多,雕刻精美,蕴含两宋特有的社会文化气息。图1、图2、图3、图4分别是南宋高浮雕青龙石刻、南宋高浮雕白虎石刻、南宋浅浮雕菊花石刻、南宋高浮雕荷花石刻。



图1 南宋高浮雕青龙石刻^[13]



图2 南宋高浮雕白虎石刻^[13]



图3 南宋浅浮雕菊花石刻^[13]

图4 南宋高浮雕荷花石刻^[13]

文物分为可移动文物与不可移动文物,石刻属于可移动文物。其中可移动文物描述元数据核心元数据包括文物类型、名称、文物识别号、所在位置、创作、材质、工艺技法、计量、描述、题识/标记、主题、考古发掘、级别、现状、来源、权限、展览/借展史、数字对象、相关文物、相关知识。具体如表1所示。

表1 可移动文物描述元数据核心元数据

序号	元素名称	定义	元素修饰词
1	文物类型	在正式的分类架构下,依据类似的特征如质地、功用等将文物归类	国家文物局普查分类 本地分类
2	名称	经审核认定的文物科学、准确、规范的名称	原名 其他名称
3	文物识别号	文物的识别编号	总登记号 其他本地号
4	所在位置	文物的当前所在位置,可以是文物所在机构,或者是文物所在行政区划	地理名称 管理机构名称 入藏日期
5	创作	文物或其主要部件的创作、设计、制造等活动信息	创作者 创作时间 创作地点
6	材质	构成文物主体材料的物质成分	
7	工艺技法	文物的制造技术、过程或方法	
8	计量	对文物进行测量所得到的尺寸、面积、体积、数量等信息	尺寸 质量

序号	元素名称	定义	元素修饰词
			数量
9	描述	以自由行文的格式描述文物的相关信息,特别是其他元素未涵盖的信息	
10	题识/标记	对镶嵌、贴、盖印、写、铭刻或附着于作品上之部份的区别或辨识描述	
11	主题	用以识别、描述和解释文物本身及其蕴含特征的术语或短语	
12	考古发掘	文物被发掘或者发现的环境	出土时间 出土地点 发掘者
13	级别	经审核认定的文物的级别	
14	现状	对文物的完残状况或者成套文物的完缺状况的具体描述	完残程度 保护优先等级
15	来源	文物入藏现收藏机构之前有关来源、传承等方面信息的描述	来源单位/个人 来源方式 入馆日期
16	权限	有关文物复制、展出或使用上的限制	
17	展览/借展史	文物被公开展示的历史记录,包括展示、借展,以及非正式展览	展览名称 策展者 展览地点 展览时间
18	数字对象	文物的数字对象,用于识别和展现文物的数字图像、声音、动画和影片以及人机互动、仿真等	数字对象识别号 数字对象关系类型 数字对象文件格式 数字对象文件日期 数字对象创建者 数字对象所属机构 数字对象权限 数字对象描述 数字对象链接
19	相关文物	描述和文物相关的文物,及文物间的关系	相关文物识别号 关系类型 相关文物链接 相关文物藏址
20	相关知识	与文物相关的人物、事件、考古、研究、术语等知识	相关知识出处 相关知识链接

纹饰不属于35类可移动文物,但是纹饰存在于大部分文物之上,参考中华人民共和国文物保护行业标准《纹饰类文物元数据规范》,在可移动文物描述元数

据核心元数据的基础之上,增加了纹样位置、纹样工艺、纹样拓片、纹样线图、纹样类型、纹样位置、纹样颜色、纹样特征、纹样意义等元素修饰词。如表2所示。

表2 纹饰描述元数据核心元数据

元素	元素修饰词
纹饰	纹饰位置
	纹饰工艺
	纹饰拓片
	纹饰线图
彩绘纹饰	纹饰类型
	纹饰位置
	纹饰颜色
	纹饰特征
	纹饰线图
	纹饰意义

针对石刻中的纹样,根据泸县石刻博物馆工作人员的建议,制定了如下的宋代石刻纹样知识体系。纹样实体由以下属性所描述:纹样名称、民族、年代、地域、构型、色彩、材质、寓意、来源。如表3所示。

表3 宋代石刻纹样知识体系

属性	关键词
纹样名称	龙纹、凤纹、卷云纹、卷叶纹样
民族	汉族、壮族、满族
年代	清朝、明朝、南宋
地域	中原、山东、四川
构型	二方连续、四方连续、波纹式
色彩	羽蓝色、水蓝色、灰色
材质	石刻、苏绣
寓意	辟邪消灾、幸福美满
来源	南宋高浮雕青龙石刻、南宋浅浮雕白虎石刻

3 泸县宋代石刻纹样知识图谱构建

泸县宋代石刻纹样知识图谱构建由以下几个步骤组成:知识图谱 Schema 设计、标签化(关键词提取)、结构化(命名实体识别、关系提取)、知识图谱存储、可视化。

3.1 知识图谱 Schema 设计

知识图谱 Schema 设计贯穿整个知识图谱的构建过程,目的是抽象出领域内的概念层次结构,定义每个概念的相关属性及概念间的关系。对于通用领域知识图谱,通常只需要宽泛地定义 Schema 或者直接使用 OpenKG 等开发知识图谱的结构,甚至采用“无 Schema”模式,直接将数据结构化 SPO 的三元组结

构即可。

但特定领域的应用,对知识的精确性要求较高。因此需要构建领域数据的 Schema 模式,包括定义数据的概念、类别、关联、属性约束等。

泸县石刻纹样知识图谱是面向特定领域的知识图谱,需要结合专家知识进行知识图谱 Schema 设计。根据第1节制定的泸县宋代纹样知识体系,给出图1的知识图谱 Schema(如图5所示),以及其对应的表4。

表4 宋代石刻知识图谱 Schema 表

实体	属性
纹样	纹样名称
纹样	民族
纹样	年代
纹样	地域
纹样	色彩
纹样	材质
纹样	来源
纹样	构型

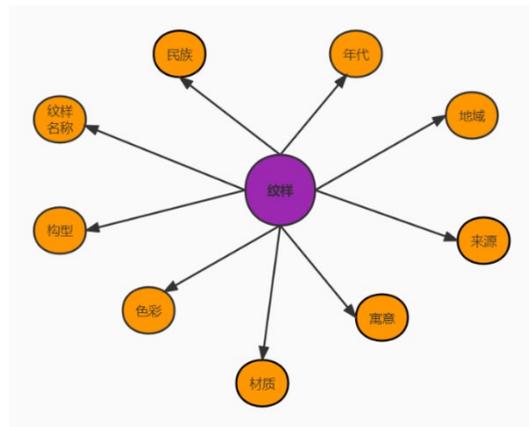


图5 宋代石刻知识图谱 Schema

3.2 标签化

标签是最能代表它们本身特质的一个词语,通过提取标签,可以更好地描述一段文字。通过关键词提取的技术手段对石刻纹饰的文字描述进行标签化,将文本中具有代表性的词或短语提取出来,用来表示文章的关键内容。

常见的关键短语抽取方法分为有监督(Supervised)和无监督(Unsupervised)。整体抽取流程则分为2个步骤:(1) Candidate Generation,得到候选短语集合;(2) Keyphrase Scoring,对候选短语进行打分。

无监督的方法由于其不需要数据标注及其普适性,得到了大范围的应用。基于统计的无监督方法主要有词频-逆向文件频率(Term Frequency-Inverse Document

Frequency,TF-IDF)^[3]以及YAKE;基于图网络的无监督方法主要有TextRank以及对应的一些变体;基于嵌入的无监督方法主要有EmbedRank。

虽然需要花费很多精力进行数据标注,但有监督方法在各个特定任务和数据集上,通常能够取得更好的效果。传统的有监督方法主要有KEA以及条件随机场(Conditional Random Fields,CRF)算法;基于深度学习的有监督方法主要包括CopyRNN以及BERT-KPE。

针对石刻纹样领域已标注数据比较少的特点,采用无监督方法进行关键词提取。经过实验对比后,基于统计的TF-IDF方法关键词提取效果更好,效率也更高。所以本节采用TF-IDF来获取实体摘要信息中关键词的频率。

TF-IDF的思想是一个词对文章的贡献可以依据文章中包含这个词的数量及其在语言集合中出现的频率来估算。TF-IDF最大的特点是保留文章里的重要词语,忽略常见但无关紧要的词语。

TF(term frequency):指的是某一个给定的词语在该文件中出现的次数,这个数字通常会被归一化(一般是词频除以该文件总词数),以防止它偏向长的文件。TF的计算如公式(1)所示。

$$TF_w = \frac{N_w}{N} \quad (1)$$

其中 N_w 是某一文本中词条 w 出现的次数, N 是该文本总词条数。

IDF(Inverse Document Frequency):反映了一个词在所有文本(整个文档)中出现的频率,如果一个词在很多的文本中出现,那么它的IDF值应该低,而反过来如果一个词在比较少的文本中出现,那么它的IDF值应该高。一个词语的重要性随着它在文件中出现的次数成正比增加,但同时会随着它在语料库中出现的频率成反比下降。IDF的计算公式如公式(2)所示。

$$IDF_w = \log\left(\frac{Y}{Y_w + 1}\right) \quad (2)$$

其中 Y 是语料库中的文档总数, Y_w 是包含词条 w 的文档数,分母数加一是为了避免未出现在任何文档中从而导致分母为0的情况。

TF-IDF的计算如公式(3)所示。

$$TF-IDF_w = TF_w \times IDF_w \quad (3)$$

从公式(2)、(3)便可以看出,某一特定文件内的高词语频率,以及该词语在整个文件集合中的低文件频率,可以产生出高权重的TF-IDF。因此,TF-IDF倾

向于过滤掉常见的词语,保留重要的词语。

在实验部分,泸县石刻博物馆专家对100段石刻相关文字进行关键词标注,采用TF-IDF算法得到的准确率为87%。下面是实验的一个示例。其中一段文字如下所示:

“南宋高浮雕青龙石刻(2002年泸县牛滩镇滩上村出王)长144厘米,宽53厘米,厚16.5厘米。龙头回首向右,张口露齿,龙须刻成锥形状并向后倾,龙头上长鬣后飘,双角略呈竹节状。背呈鱼鳍状,身刻鳞甲,尾弯曲略呈S形。四爪张开踏在云纹之上,呈飞奔状。龙尾上方刻有火焰宝珠。”

专家对这段文字进行关键词标注,得到的排名前5关键词如下:“高浮雕”、“石刻”、“竹节状”、“鱼鳍状”、“云纹”;通过TF-IDF得到的排名前5关键词结果如下:“南宋”、“高浮雕”、“石刻”、“鱼鳍状”、“云纹”。

3.3 结构化之命名实体识别

在线博物馆及作为开放领域知识共享平台的在线百科,都为石刻纹样知识图谱构建提供了规模巨大且类型多样的数据资源。本节从开放领域物文本的特点出发,从在线博物馆和在线百科网站等开放数据中抽取石刻文物实体,为石刻纹样知识图谱的构建提供基础。结合已有的四川泸县石刻博物馆提供的石刻结构化信息进行实体抽取。

早期的实体抽取方法中,基于知识的方法较常见,这些方法基于词典和领域知识。但是,由于领域和语言的规则特殊性以及词典的不完备性,这类实体抽取方法在召回率上往往很低,同时,需要领域专家对知识资源进行构建和维护。为了解决知识系统的弊端,许多学者采用统计机器学习模型,机器学习通过训练示例的输入并学习预期输出实现预测,从而取代人工设计的规则和词典。典型的最大熵模型中,随样本数量的增加而增长的约束条件会导致计算量增大;隐马尔可夫模型只关注当前词语特征而忽略了上下文,使特征选择受限;CRF考虑到全局所有特征根据上下文序列计算全局最优值。然而,CRF实现的实体抽取任务依赖于人工获取特征模板,大量的特征获取费时费力,且这些人工特征一般都会面向某一特定领域,具有一定的领域局限性。

神经网络模型,特别是近年来在各个领域的分类任务中表现优秀的深度学习模型,引入预训练语言模型替换这些人工构建的特征向量来实现词表征,然后使用标注数据对模型进行训练,解决了人工模板的弊

端,实现了实体抽取的自动化。但是传统的循环神经网络存在梯度消失/爆炸,对距离较长依赖处理并不擅长。长短时记忆网络(Long Short-Term Memory, LSTM)^[5]通过精巧地设计“门”控方式治理了长期依赖的梯度问题,但是LSTM只能依据前一时刻的时序信息来预测下一时刻的输出。双向长短时记忆循环神经网络(Bi-directional LSTM, BiLSTM)^[6]通过前馈和反馈网络,不但利用前馈网络参照上文,也会利用反馈网络同时兼顾下文,承上启下的捕获上下文特征,但是,BiLSTM预测每个标签的过程是独立的,未利用上文已预测的标签,所以不能像CRF一样对全局序列进行优化处理。为解决BiLSTM和CRF单独使用的弊端,许多学者将二者进行结合来实现实体抽取,以获得相对优异的实体抽取效果。

针对文物领域标注数据缺乏以及文本构词具特

殊性的问题,本节提出基于自训练算法的半监督实体抽取方法,使用ELMo(Embedding from Language Model)^[4]模型动态获取石刻纹样实体特征,设计BiLSTM和CRF模型的结合来实现文本特征提取和实体标签的获取,设计双重标注样本选择策略,选择置信度高的样本参与模型训练,提高模型的性能,以实现使用少量标注和大量无标注文物实体数据获得有效的石刻纹样实体抽取结果。

ELMo使用多层BiLSTM来训练语言模型,然后通过线性组合不同LSTM层的word vectors,得到最终的word embedding vectors。线性组合公式如公式(4)所示:

$$ELMo_k^{task} = E(R_k; \Theta^{task}) = \gamma^{task} \sum_{j=0}^L s_j^{task} h_{k,j}^{LM} \quad (4)$$

公式(4)的直观解释如图6所示。

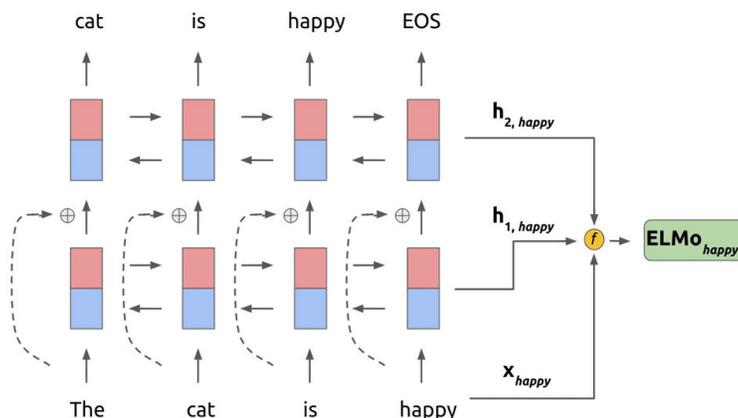


图6 ELMo公式直观解释

为解决石刻纹样实体标注数据少,数据构词特殊性的问题,提出基于自训练的半监督实体抽取模

型如图7所示,主要包括模型预训练和样本选择两个阶段。

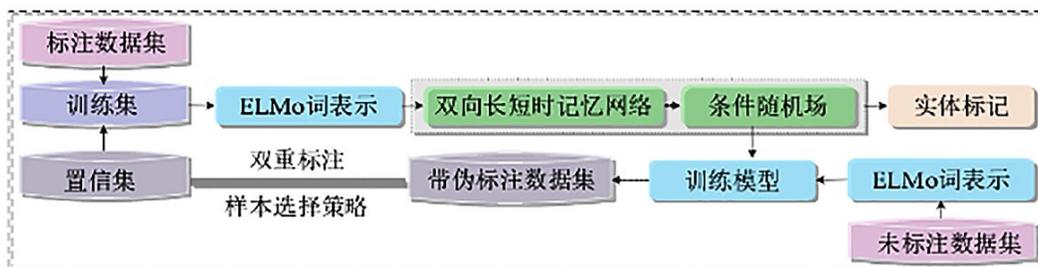


图7 自训练半监督实体抽取模型

第一阶段,将少量的标注数据输入到预训练的模型中,使用预先训练好的ELMo语言表示模型对输入序列生成词嵌入,接着,将生成的词嵌入用于训练BiLSTM-CRF模型得到训练模型,其中,BiLSTM用于

上下文特征的提取,CRF用于实体预测,同时,用训练模型对大量的无标注数据进行标签预测,得到无标注数据的伪标记。

第二阶段,对于第一阶段预测出的伪标记数据,使

用基于双重标注样本选择策略进行样本选择,选择高置信度的样本,形成置信集合,然后将该置信集合加入标注样本数据集,更新训练集,并对模型进行下一轮训练,在下一轮迭代中对模型参数进行更新。BiLSTM-CRF模型训练和高置信度样本选择这两个过程循环进行,直到无标注数据集为空或达到一定迭代次数时停止。

实验选择最具代表性的四类实体,石刻文物名称实体、朝代实体、出土地点实体、纹样名称实体,采用准确率、召回率、F1值进行评价,实验结果如表4所示:

表4 结构化之命名实体识别实验结果

	石刻文物名称实体	朝代实体	出土地点实体	纹样实体
精确率	83.21%	87.55%	86.21%	87.56%
召回率	83.10%	88.34%	85.10%	87.78%
F1值	84.37%	87.21%	84.37%	88.07%

可以看到,石刻文物名称实体的抽取效果并不是很好,主要是因为石刻文物实体名称长度存在界限模糊的问题,比如“南宋高浮雕青龙石刻”是石刻文物实体,但是算法会抽取“青龙石刻”作为石刻实体。

3.4 结构化之关系提取

通过石刻纹样实体抽取任务,获得一系列形式上独立且未关联的、离散的石刻纹样实体,要深入挖掘石刻纹样实体之间的内在联系,需要进一步挖掘石刻纹样实体间的关系,从而获得石刻纹样间的事实知识,实现石刻纹样实体间的语义链接。本节基于石刻纹样实体抽取的结果,从在线博物馆和在线百科等开放数据中抽取石刻纹样实体间的关系,为石刻纹样知识图谱构建提供有效三元组。

石刻纹样数据在提供石刻纹样实体的同时,还蕴含了实体之间的联系,通过挖掘实体间的关系可以得到实体语义信息,为实现实体间的语义链接,形成网状的知识结构提供基础,有助于理解石刻纹样背后的知识,扩展石刻纹样知识的广度和宽度。在石刻纹样实体抽取的前提下,从海量的网络石刻信息中自动提取潜在知识的第二步是关系抽取,它是信息提取和构建知识图谱的基础且是不可逾越的步骤,为提取有效三元组提供保障。

在关系抽取任务中,研究者使用了许多不同的模型,最早的是基于特征的方法,通过抽取高效且高识别度的显著特征实现关系实例的分类。Kambhatla等人^[7]考虑运用句子中的词语和结构信息训练最大熵模型,解决了标注数据稀缺问题,以及实体检测模块在提取实体之间语义关系时产生的错误问题。Zhou等人^[8]加入了额外

的特征,系统地探索了基于特征的关系提取方法,他们在句法方面结合了基本短语组块信息,并使用语义信息来提高性能。虽然这些方法对关系提取任务有一定的影响,但是特征的抽取效果在一定程度上决定了这类方法的整体性能,且由于人为误差,致使在实验中出现误差升级现象。

针对特征方法中特征抽取对模型的影响,研究者使用核函数计算相似度来实现两个实例关系的判断,从而避免了特征工程的缺点。Mooney等人^[9]提出一种基于子序列核推广的方法提取实体间语义关系,这种核方法使用三种类型的子序列模式识别两个实体之间的关系。Khayyamian等人^[10]提出了一种广义卷积核,这种广义核具有更强的灵活性和可定制性,可以方便地用于系统地生成更有效的特定语法子核,利用这些语法子核来提取关系。卷积核根据两个权重函数的不同定义产生不同的子核,用于提取关系,但是,基于核函数的方法具有较高的时间复杂度,对于渐趋增长的大数据显得捉襟见肘,且基于互联网提供的大量数据,需要从这些原始数据中提取有效特征,但传统的方法缺乏以原始形式处理自然数据的能力。

近年来,神经网络模型能够避免人工参与而捕获原始数据的特征,减轻了人工标注的繁重工作,性能也优于人工抽取特征和核函数的传统方法。最早被应用于关系抽取的神经网络模型是卷积神经网络CNN,递归神经网络RNN是另一种常用的选择,双向长短期记忆BiLSTM是为了改进RNN的特征获取而提出的。早期的深度学习模型将关系抽取任务视为一个多分类问题,在每个句子中添加一个关系类标签。Santos等人^[11]提出了一个分类关系的方法,该方法使用卷积层为文本生成一个分布式向量表示,并为每个类创建一个分数,将分布式向量表示与类表示进行比较。但是,CNN模型不能很好地表示高层部分之间精确的空间关系。Sorokin等人^[12]利用LSTM共同捕获存在关系的嵌入,通过结合上下文关系和目标关系的表示生成最终的预测结果。以上这些系统可以有效预测一个句子中单个实体对的一个关系。

现有关系抽取研究虽取得一定的进展,但石刻纹样数据中存在多重关系以及关键词具稀疏性,对现有方法提出挑战。本节提出一种基于词注意力机制的胶囊网络关系抽取方法,命名为Word-Attention Capsule net, WAtt-Capsnet。首先,设计字、词、词性和位置组合嵌入作为模型输入用于捕获词语的语义特征和上下文信息及语序特征。其次,利用双向长短期记忆模型获取底层

特征,利用胶囊网络捕获更高层次的语义特征。最后,为了抽取实体间多重关系,解决关系稀疏性的难题,提出一种基于词注意力机制的动态路由算法,考虑不同词的强度贡献,通过增加信息词的权重系数,迭代修正连接强度来解决关键词稀疏问题,以此有效实现石刻纹样

领域多重关系的抽取。

本节提出的模型是原始胶囊网络的一个变形,模型的体系结构包括四个层次:输入表示层、低层特征提取层、基于词注意力机制的动态路由胶囊网络特征聚类层和关系预测层。如图8所示。

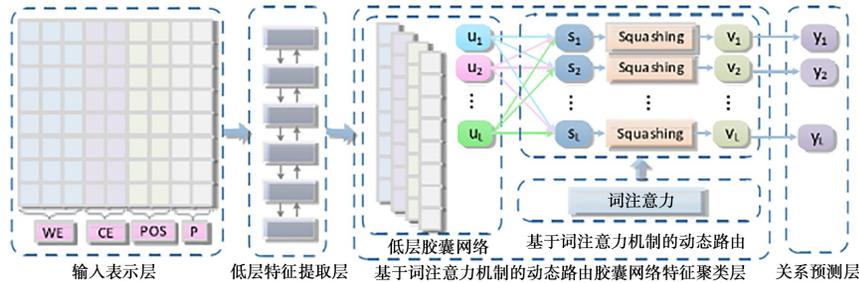


图8 关系抽取模型

(1).输入表示层。模型的输入表示是融合四个嵌入的组合,分别是字嵌入、词嵌入、词性嵌入和位置嵌入,以获取每个字和词的语义信息的同时,获取句子的上下文信息以及词语在句子中的语序信息。

(2).低层特征提取层。对于输入层的词表示,经过BiLSTM模型的前向网络和后向网络获取全局序列信息,以此来捕获句子的低层特征。

(3).基于词注意力机制的动态路由胶囊网络特征聚类层。采用带有词注意力动态路由的胶囊网络减少不相关词产生的噪声,实现高层特征的提取,以此捕获实体间多重关系。

(4).关系预测层。利用边际损失函数来预测实体间的关系。

实验选择最具代表性的四类关系,出土时间、出土地点、对应纹样、文物朝代,采用准确率、召回率、F1值进行评价,实验结果如表5所示:

表5 结构化之关系抽取实验结果

	出土时间	出土地点	对应纹样	文物朝代
精确率	84.35%	91.45%	83.34%	90.03%
召回率	84.90%	91.47%	84.21%	89.95%
F1值	83.77%	90.88%	83.59%	89.37%

可以看到“文物朝代”关系抽取效果很好,主要是因为数据集中大多是南宋的石刻文物。“出土地点”抽取效果很好,主要是因为该关系有明显的模式:“出土于”。“出土时间”跟“对应纹样”关系抽取效果不好,主要是因为并不是每一段文本数据中都有对应的“出土时间”跟“对应纹样”的信息。

3.5 知识图谱存储

一般情况下,我们使用数据库查找事物间的联系的时候,只需要短程关系的查询(两层以内的关联)。当需要进行更长程的,更广范围的关系查询时,就需要图数据库的功能。图数据库(Graph database)指的是以图数据结构的形式来存储和查询数据的数据库。Neo4j是目前用的最多的图数据库,世界数据库排行榜上排名21位。Neo4j属于原生图数据库,其使用的存储后端是专门为图结构数据的存储和管理进行定制和优化的,在图上互相关联的节点在数据库中的物理地址也指向彼此,因此更能发挥出图结构形式数据的优势。知识图谱中,知识的组织形式采用的就是图结构,所以非常适合用Neo4j进行存储。

3.6 可视化

根据以上步骤得到的三元组数据,通过Echarts可视化框架对存储在Neo4j图数据库中的知识图谱数据进行展示。其中前端界面主要使用Echarts框架,后端使用Python的Django Web框架,数据库部分使用MySQL数据库以及Neo4j图数据库。

4 总结

宋代石刻纹样蕴含着宋代人文意志和艺术风格,是进行历史和文化研究的宝贵资源。本文通过研究泸县宋代石刻纹样,深度挖掘其知识体系,结合知识图谱相关技术,搭建了一套基于泸县宋代石刻纹样的知识图谱构建路径,为传统文化数据的标签化和结构化提供了新的解决方案。

参考文献(References):

- [1] 李雅梅,张春新. 川南泸县南宋墓葬鸟兽石刻的象征意义[J]. 文艺研究, 2009, (1): 152-153.
- [2] 李能萍. 四川泸县宋代石刻的艺术价值研究[J]. 大观, 2021, (8): 90-91.
- [3] 施聪莺,徐朝军,杨晓江. TFIDF算法研究综述[J]. 计算机应用, 2009, 29(z1): 167-170+180.
- [4] Peters M E, Neumann M, Iyyer M, et al. Deep contextualized word representations [C]//Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), Association for Computational Linguistics, 2018: 2227-2237.
- [5] Shi X, Chen Z, Wang H, et al. Convolutional LSTM Network: a machine learning approach for precipitation nowcasting[C]//NIPS-2015, MIT Press, 2015.
- [6] Huang Z, Wei X, Kai Y. Bidirectional LSTM-CRF models for sequence tagging[DB/OL]. arXiv:1508.01991.
- [7] Kambhatla N. Combining lexical, syntactic, and semantic features with maximum entropy models for extracting relations[C]//ACLdemo '04: Proceedings of the ACL 2004 on Interactive poster and demonstration sessions, 2004:22-es.
- [8] Zhou G, Su J, Zhang J, et al. Exploring various knowledge in relation extraction [C]//Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL'05), 2002: 427-434.
- [9] Bunescu R C, Mooney R J.. Subsequence kernels for relation extraction [C]//International Conference on Neural Information Processing Systems, MIT Press, 2005: 171-178.
- [10] Khayyamian M, Mirroshandel S A, Abolhassani H. Syntactic tree-based relation extraction using a generalization of Collins and Duffy convolution tree kernel [C]//Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Companion Volume: Student Research Workshop and Doctoral Consortium, 2009: 66-71.
- [11] Santos C, Xiang B, Zhou B. Classifying relations by ranking with convolutional neural networks [C]//Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing, 2015: 626-634.
- [12] Sorokin D, Gurevych I. Context-aware representations for knowledge base relation extraction [C]//Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, 2017: 1784-1789.
- [13] 四川泸县宋代石刻博物馆[EB/OL]. 2022-07-22. <http://lxm.vip.gumaor.com>.

编辑:王谦

(上接第32页)

- deep features [DB/OL]. arXiv preprint arXiv:1707.04916, 2017.
- [3] Chen T, Wang S, Chen S. Deep multimodal network for multi-label classification [C]//2017 IEEE International Conference on Multimedia and Expo (ICME), IEEE, 2017:955-960.
- [4] Srivastava N, Salakhutdinov R R. Multimodal learning with deep Boltzmann machines [J]. The Journal of Machine Learning Research, 2014, 15(1): 2949-2980.
- [5] Jiang Q Y, Li W J. Deep cross-modal hashing [C]//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2017: 3232-3240.
- [6] Kulkarni G, Premraj V, Ordonez V, et al. Babytalk: Understanding and generating simple image descriptions [J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2013, 35(12): 2891-2903.
- [7] 赵海英,高子惠,邓恋,侯小刚,李宁. 基于图文混排的传统服饰图像以文标图算法[J]. 图学学报, 2021, 42(03): 398-405.
- [8] Ben-Younes H, Cadene R, Cord M, et al. Mutan: Multimodal tucker fusion for visual question answering [C]//Proceedings of the IEEE International Conference on Computer Vision, IEEE, 2017: 2612-2620.
- [9] Fukui A, Park D H, Yang D, et al. Multimodal compact bilinear pooling for visual question answering and visual grounding [DB/OL]. arXiv preprint arXiv: 1606.01847, 2016.
- [10] Dong X L, Hajishirzi H, Lockard C, et al. Multi-modal information extraction from text, semi-structured, and tabular data on the web [C]//Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, IEEE, 2020: 3543-3544.
- [11] Yu J, Jiang J, Yang L, et al. Improving multimodal named entity recognition via entity span detection with unified multimodal transformer [C]. Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, 2020:3342-3352.

编辑:王谦