

李靖, 李泽荃, 石福泰, 等. 基于概率融合算法的煤矿事故隐患文本知识实体抽取研究[J]. 矿业科学学报, 2024, 9(6): 1007-1016. DOI: 10.19606/j.cnki.jmst.2024915

LI Jing, LI Zequan, SHI Futai, et al. Textual knowledge entity extraction of hidden dangers in coal mine accidents based on probabilistic fusion algorithm[J]. Journal of Mining Science and Technology, 2024, 9(6): 1007-1016. DOI: 10.19606/j.cnki.jmst.2024915

基于概率融合算法的煤矿事故隐患文本 知识实体抽取研究

李靖^{1,2}, 李泽荃³, 石福泰⁴, 郝强⁵

1. 青海师范大学国家安全与应急管理学院, 青海西宁 810016;

2. 中国矿业大学(北京)能源与矿业学院, 北京 100083;

3. 华北科技学院, 河北廊坊 065201;

4. 华亭煤业集团有限责任公司, 甘肃平凉 744100;

5. 华能煤炭技术研究有限公司, 北京 100070

摘要:针对煤矿事故隐患文本数据的非结构化特性,基于煤矿事故隐患文本数据集,通过分析隐患描述文本数据的特征及隐含信息,结合事故隐患传播规律设计了适用于煤矿事故隐患描述文本的知识实体标注类型并使用 Brat 工具进行标注,构建用于知识实体抽取模型的数据集;提出一种基于动态权重融合的 BERT-IDCNN-CRF 模型,并引入基于牛顿冷却定律的概率融合算法。结果表明:引入概率融合算法后,动态权重融合的 BERT-IDCNN-CRF 在隐患文本知识实体抽取任务中表现最佳,其精度、召回率与 F_1 值分别提升了 8.93%、5.28%、7.51%,显著提高了模型的预测准确性和稳定性,并具有良好的适应性。

关键词:煤矿事故隐患;知识实体抽取;K 折交叉验证;概率融合

中图分类号:TD 77+1

文献标志码:A

文章编号:2096-2193(2024)06-1007-10

Textual knowledge entity extraction of hidden dangers in coal mine accidents based on probabilistic fusion algorithm

LI Jing^{1,2}, LI Zequan³, SHI Futai⁴, HAO Qiang⁵

1. School of National Safety and Emergency Management, Qinghai Normal University, Xining Qinghai 810016, China;

2. School of Energy and Mining Engineering, China University of Mining and Technology-Beijing, Beijing 100083, China;

3. North China University of Science and Technology, Langfang Hebei 065201, China;

4. Huating Coal Industry Group Co., Ltd., Pingliang Gansu 744100, China;

5. Huaneng Coal Technology Research Co., Ltd., Beijing 100070, China

Abstract: Given the unstructured nature of text data related to hidden dangers in coal mine accidents, extracting latent knowledge is crucial for constructing a knowledge graph of hidden dangers in coal mine accidents. This study proposes annotation types for knowledge entities to describe hidden dangers in coal mine accidents by analyzing the characteristics and latent information in the texts of hidden dangers

收稿日期:2024-06-11 修回日期:2024-10-18

基金项目:华能集团总部科技项目(HNKJ20-H33);青海省基础研究计划(2024-ZJ-904)

作者简介:李靖(1987—),男,山西定襄人,博士,讲师,主要从事安全生产信息化、智慧矿山等方面的研究工作。E-mail:lijing6816@163.com

通信作者:李泽荃(1983—),男,山西长治人,博士,副教授,主要从事煤矿智能化等方面的研究工作。E-mail:lzquancumtb@126.com

based on their propagation patterns. Using the Brat annotation tool, we annotated the text data related to hidden dangers of coal mine accidents to construct a dataset for knowledge extraction model. We proposes a BERT-IDCNN-CRF model based on dynamic fusion and introduced a probabilistic fusion algorithm based on Newton’s law of cooling. The results indicate that with the incorporation of the probabilistic fusion algorithm, the dynamically weighted BERT-IDCNN-CRF model achieved the best performance in the task of knowledge entity extraction from hidden danger texts. Its precision, recallrate, and F_1 -score improved by 8.93%, 5.28%, and 7.51%, respectively, significantly enhancing the model’s prediction accuracy and stability, while demonstrating excellent adaptability.

Key words: hidden dangers in coal mine accidents; knowledge entity extraction; K -fold cross-validation; probabilistic fusion

随着煤矿信息化程度的日益增强和煤矿安全生产监管体系的不断完善^[1],煤矿企业和监管部门都越来越重视煤矿事故隐患排查治理工作,积极管理并存储煤矿事故隐患文本数据。文本数据中蕴含着大量的煤矿事故隐患特征知识。如何科学地抽取隐患特征知识对于煤矿事故隐患的排查治理、煤矿安全生产水平的提高至关重要,也是提升智能化煤矿数据治理能力的重要途径^[2]。

煤矿隐患文本数据的知识实体抽取是构建煤矿隐患领域知识图谱的重要基础工作,旨在从未经标注的隐患文本数据中提取实体知识,并归类到预先设定的实体类别。知识实体抽取方法主要分为基于规则和词典、基于统计及基于深度学习3类。基于规则和词典方法是最早的实体抽取方法,主要依赖于经验和领域知识^[3-4];基于统计方法主要依赖数据驱动的学习方式,通过从大量标注数据中学习模型^[5-7];基于深度学习方法是随着深度学习技术崛起而发展起来的,能够处理大规模数据,捕获复杂的模式和上下文关系。Google 研究团队于2018年引入了基于双向Transformer的预训练语言模型,为实体识别带来了突破性进展^[8]。鹿晓龙^[9]提出了融合ALBERT预训练语言模型的BILSTM-CRF模型,充分考虑了煤矿安全文本中存在的一词多义的现象,并进一步提升了煤矿安全文本实体识别的准确率。赵彭彭^[10]在ALBERT-IDCDA-CRF模型的基础上提出FS-Transformer-CRF算法,通过合并简化编码器中的结构,减少训练时间成本消耗的同时提升了模型性能。

上述研究方法在文本实体识别的精确度和召回率上表现不佳,主要归因于模型过于简单或对特定应用领域的过度依赖,大多数研究主要面向开放领域或者语料较为丰富的垂直领域,如医学领域或军事领域等。本文提出一种面向煤矿事故隐患领域的深度学习知识实体识别算法,融合了具有多层动态权重的BERT(Bidirectional Encoder Represen-

tations from Transformers)、IDCNN(Iterated Dilated Convolutional Network)和CRF(Conditional Random Field)模型的优势,并引入牛顿冷却定律模型进行优化,在保留BERT模型主要优点的同时,有效地缩减模型训练所需的参数数量,显著提升模型识别的效果及性能。

1 煤矿事故隐患文本知识实体抽取框架

煤矿事故隐患文本知识实体抽取中包含煤矿事故隐患文本标注、模型调优及最优模型的预测评估3个步骤。为有效实现煤矿事故隐患文本知识实体抽取任务,识别隐患文本中定义的实体信息并用于领域知识图谱的构建,采用基于Brat标注的煤矿事故隐患文本知识实体抽取框架。整体框架由煤矿事故隐患文本标注层、输入层、模型层与输出层组成,如图1所示。

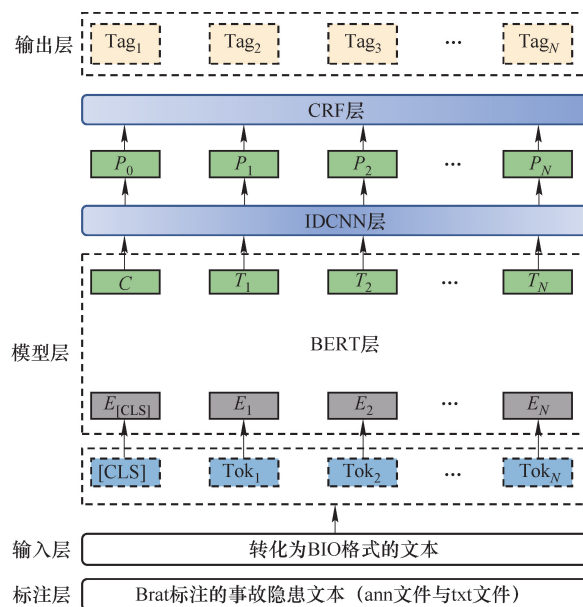


图1 煤矿事故隐患文本知识实体抽取框架

Fig. 1 Knowledge entity extraction framework for hidden dangers in coal mine accidents

(1) 标注层。在对煤矿事故隐患文本标注之前,必须先为煤矿事故隐患描述文本知识提取任务设立明确的标准和规范。确定需要标注的实体类型后,通过 Brat 标注工具完成实体的标注及审核,建立煤矿事故隐患标注的原始语料。

(2) 输入层。输入层将标注层中的原始语料统一按照 BIO 标记法执行自动标注数据的转换,形成 BIO 格式文本,作为模型层的输入,为后续特征提取与序列标注提供数据基础。

(3) 模型层。模型层从下到上共有 2 层,即 BERT 层和 IDCNN 层。BERT 层主要对输入的文本数据进行预处理和编码,结合输入的词向量、段落向量和位置向量,通过多层注意力机制进行交互和整合,最终输出高质量的语义表示;IDCNN 层主要用于从文本数据中提取更加丰富和复杂的特征,通过卷积操作得到表示文本的向量。

(4) 输出层。输出层包括了 CRF 层和模型的最终输出结果。CRF 层作为输出层的关键部分,负责对隐患文本中的实体信息进行精确标注,确保模型能够正确地识别和标注文本中的实体信息。在实验训练得到最优模型后,将未标注实体的隐患文本数据输入模型,模型的输出层会生成相应的实体标注结果,进一步为领域知识图谱的构建奠定数据基础。

2 煤矿事故隐患文本数据语料标注及预处理

在开展煤矿事故隐患知识提取任务过程中,依赖已有的煤矿事故隐患描述,通过制定标注规范并应用相关工具进行手工标注形成训练语料集。对于已标注的文本数据,在输入 BERT 模型后转换生成相应的嵌入表示向量,进而为 IDCNN 层与 CRF 层提供煤矿事故隐患的文本向量输入。

2.1 隐患文本特征分析及实体知识定义

在构建实体数据集前,需要定义煤矿事故隐患文本中的实体类别。煤矿事故隐患传播规律揭示了隐患动态传播机制、触发模型及耦合模式为核心的关键要素,不仅为探索隐患传播的内部机制奠定了基础,也为从隐患描述文本中确定知识实体类别提供了有力支撑。根据历史样本数据,煤矿事故隐患描述文本具有如下 3 个特点。

(1) 隐患描述文本在格式上呈现出一致性,通常涉及人员、设备、环境、制度等具体隐患的简洁描述。

(2) 大多数事故隐患描述文本存在标点等特

殊字符,可以将文本划分为多个子句,包含了重要的隐患信息。

(3) 煤矿事故隐患描述文本中通常包含了人的不安全行为、物的不安全状态、环境的不安全因素以及管理上的缺陷等相关内容。

因此,参照煤矿事故隐患定义^[11],同时结合煤矿事故隐患规律中关键要素所涉及事故隐患网络节点,设定了 10 类知识实体(表 1)。

表 1 煤矿事故隐患领域知识实体定义

Table 1 Definition of knowledge entities in hidden dangers in coal mine accidents

序号	实体	标签	示例
1	人员	Person	瓦斯检查工、支护工、探放水工、抽水工、安监部长、爆破工等
2	组织机构	Organization	掘进队、采煤区队、煤质科、地测部门等
3	设备设施	Equipment	带式输送机、综掘机、电缆、局部通风机、防火铁门、托辊、锚索、风筒、架空乘人装置、灭火器、通风系统等
4	地点	Place	回风巷、掘进工作面、变电所、临时水仓、运输巷、采区轨道下山等
5	部位	Position	带式输送机头、进风侧、末端、迎头后方、交叉口等
6	制度	Regulation	煤矿安全规程、机电管理制度、回风巷掘进作业规程、密闭检查制度、区域综合防突措施等
7	人的行为	Action	检测弄虚作假、未检查瓦斯、检查方法错误、记录不规范、未分析、未悬挂、未按作业规程要求、不能熟练使用、未签字、未携带甲烷检测报警仪等
8	物的状态	State	底鼓、顶板垮落、伞檐、喷浆体开裂、煤壁片帮、初撑力、顶板破碎等
9	环境因素	Factor	瓦斯超限、瓦斯涌出、无照明、有滴水、无警示牌、漏水严重等
10	管理因素	Management	不符合规定、未设置、未建立、未装设、未留设、未组织开展、检查结果无定期审查、瓦斯含量报告单、探放水记录、瓦斯日报表、采掘工程平面图、密闭墙管理台账、矿井专项防突设计、隐患排查治理情况、水灾事故应急演练报告、施工组织设计等

2.2 标注工具

文本标注是自然语言处理中的一项重要任务,主要是针对文本中出现的实体和关系等知识内容进行识别和分类。标注人员需要根据预定义的标注规则进行标注,确保标注结果符合标准,让后续的自然语言处理任务可以利用标注结果进行不同任务的模型训练与评估。目前有许多流行的文本标注工具,如 DeepDive^[12]、Doccano^[13]、Prodigy^[14]、YEDDA^[15]和 Brat^[16]等,见表2。

表2 文本标注工具特点对比

Table 2 Comparison of text annotation tools

工具名称	是否开源	主要功能	优点	缺点
DeepDive	开源	支持关系抽取	具有处理大规模、非结构化数据的处理能力,适合大规模的关系抽取任务,自定义性强且具有可扩展性	使用该工具需要一定的技术和计算资源支持,适用门槛相对较高
Doccano	开源	支持情感分类、命名体识别、序列标注等任务	可多人合作标注,可添加较多的标签数,具有很友好的用户界面	需要将其部署在服务器上,且需要一定的技术支持,对于一些安全性要求比较高的项目,可能存在风险
Prodigy	收费	支持实体标注、分类标注、情感标注	支持多种标注类型且具有强大的自定义功能,标注精度和效率方面表现优异	具有商业化的模式,定价较高
YEDDA	开源	支持块、实体、事件的标注	轻量级标注工具;支持快捷键方式进行文本的手动标注,界面直观友好	标注最大标签数有限制
Brat	开源	支持实体、关系、事件及属性的标注	轻量级标注工具;功能较全面,支持多用户协作标注,可自定义快捷键,学术界使用较多	只能用于Linux和Macos平台

综合考虑各类标注工具的开源性、可扩展性以及便捷性,在进行煤矿事故隐患文本标注的过程中选择 Brat 工具。使用 Brat 标注文本时,将每一个被标注存储文本的 txt 格式文件,生成一一对应存储标注结果的 ann 格式文件。ann 格式文件中包含的标注信息及其与原文本的关联,可以被 Brat 或其他工具进行读取和解析。具体字段信息包括标注的实体序号、实体类型、实体在 txt 格式文本中的起始位置和结束位置、实体本身等,如对“北山副井井下三水平消防材料库未按要求配备消防器材”隐患描述文本进行标注后的 ann 格式文件见表3。

表3 ann 文件内容

Table 3 ann file content

实体序号	实体类型	实体起始位置	实体结束位置	实体本身字段
T_1	Place	3	5	副井
T_2	Place	10	15	消防材料库
T_3	Management	15	21	未按要求配备
T_4	Equipment	21	25	消防器材

2.3 标注数据的预处理

利用 Brat 工具完成知识实体的标注任务后,每一条隐患描述文本都有对应的 ann 格式文件,此时还需要通过 BIO 标注法将 ann 格式文件进行自动标注,作为模型层的输入数据。通过将 Brat 工具标注的文本转换为自动标注好的字符后,获得煤矿事故隐患文本知识实体抽取的输入数据。

3 煤矿事故隐患文本知识实体抽取模型构建

将动态权重融合的 BERT 模型引入到 IDCNN 模型,将输出的预测序列输入到 CRF 模型,构成一个基于动态权重融合的 BERT-IDCNN-CRF 模型。为充分利用 BERT 的多层表示能力,使用动态权重融合的方法整合不同层次的 BERT 编码器的特征。具体地,使用一种基于加权平均的动态权重融合方法,根据不同层次 BERT 编码器的表征能力及文本特征,动态地调整不同层次 BERT 编码器的权重,以实现煤矿事故隐患文本知识的自动抽取。

3.1 动态权重融合的 BERT 模型

BERT 模型是由 DEVLIN 等提出的一种预训练语言表征模型,具备以下特点:①首次提出使用 Masked 语言模型和下一句预测两个无监督预测任务对文本数据进行预训练;②采用双向 Transformer 作为编码器;③完全采用了自注意力机制和全

连接层对输入文本进行建模。其基本结构^[5]如图 2 所示。

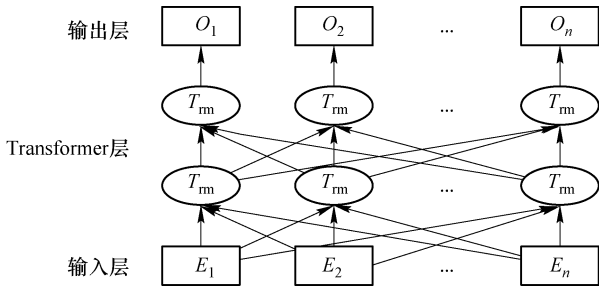


图 2 BERT 模型基本结构

Fig. 2 Architecture of the BERT model

动态权重融合的 BERT 模型是在使用 BERT 模型对文本进行编码时,将其 12 层 Transformer 编码器生成的表示赋予一个权重(w),该权重值可以通过训练来确定。此外,还需要通过最后一层的全连接层将 768 维度降至 512 维,结合其他层获得输出。对 BERT 模型不同输出层进行加权平均,得到文本向量。BERT 模型各层权重融合公式如下:

$$V_{\text{final}} = \sum_{i=1}^N W_i \times V_i \quad (1)$$

式中, V_{final} 为最终的表示层; V_i 为第 i 层生成的表示; W_i 为第 i 层的权重; N 为模型的总层数。

动态权重融合的 BERT 模型结构如图 3 所示。

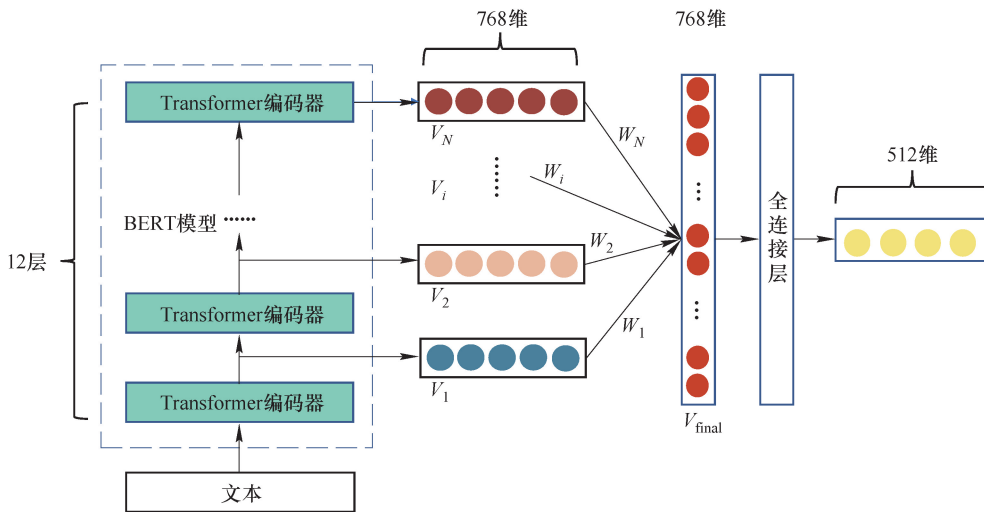


图 3 动态权重融合的 BERT 模型结构

Fig. 3 Schematic of BERT model structure of dynamic weight fusion

由图 3 可看出,每个 Transformer 编码器层都负责不同层次的文本信息特征提取,浅层 Transformer 编码器层主要对局部信息进行建模,深层 Transformer 编码器层对更抽象和全局的语义进行建模。由于每个编码器层都会生成一组词向量表示,随着层数的增加,计算资源和内存的消耗也会相应增加。同时,过多的层可能会引入冗余的信息,甚至导致过拟合。因此,在模型层数时,选择了输出计算代价相对较小的模型的最后 5 层。

3.2 IDCNN 模型

膨胀卷积是一种在卷积神经网络中使用的技术,旨在扩大卷积核的感受野,同时避免增加计算量和模型参数。IDCNN 模型结合了膨胀卷积的优势,能在维持模型参数数量和计算速度不变的前提下扩展卷积核的识别域,从而更好地处理长序列文本数据^[17]。在卷积神经网络中,卷积核在连续的输入区域上执行操作以提取特性。但在膨胀卷积中,通过引入膨胀因子,卷积核在处理时会跳过特

定间隔的输入数据,这使得在不改变卷积核尺寸的情况下,膨胀卷积能够覆盖更宽阔的输入区域,增强感受野。与卷积神经网络相比,膨胀卷积的结构更为复杂,但在卷积过程中的跳跃性质为神经网络的设计者提供了更多的灵活性。

3.3 CRF 模型

IDCNN 模型可以学习单词或短语的上下文信息,从而输出概率最高的标签结果,但无法学习不同单词之间的关系,输出的标签缺乏逻辑性,且输出结果可能较为混乱。因此,引入 CRF 模型,在做预测标签时,通过加入约束条件保证词标签的正确性。

CRF 模型的操作过程为,特征函数 $f(x)$ 是事先定义好的训练模型,使用给定的数据训练模型并确定权重参数 λ_k ,确定模型后用于实现序列注释。CRF 模型是一种判别式条件概率分布模型,可以用 $P(Y|X)$ 表示,其中 X 是输入变量,代表标记的观察序列,而 Y 是输出序列,代表与 X 对应的标签

序列。X 是观察序列 (x_1, \dots, x_n) , Y 是隐藏状态序列 (y_1, \dots, y_n) , 每个 (x_i, y_i) 对是线性链中最大的团, 并且满足:

$$P(y_i | x, y_1, y_2, \dots, y_n) = P(y_i | x, y_{i-1}, y_{i+1}) \quad (2)$$

给定一个预设的观察序列 x , 通过 CRF 模型求解隐态序列 y 的方程如下:

$$P(y | x) = \frac{1}{Z(x)} \prod_i \exp \left[\sum_k \lambda_k f_k(y_{i-1}, y_i, x, i) \right] \\ = \frac{1}{Z(x)} \exp \left[\sum_{i,k} \lambda_k f_k(y_{i-1}, y_i, x, i) \right] \quad (3)$$

式中, i 为节点的当前位置; k 为当前的特征函数, 每个特征函数被赋予 1 个权重 λ_k 。

因为状态序列和 2 个标记之间的关系有限, 所以为 CRF 模型定义了连续特征函数, 也就是转移特性函数和状态特性函数:

$$P(y | x) = \frac{1}{Z(x)} \exp \left[\sum_{i,j} \lambda_j t_j(y_{i-1}, y_i, x, i) + \sum_{i,l} \mu_l s_l(y_i, x, i) \right] \quad (4)$$

$$Z(x) = \sum_y \exp \left[\sum_{i,j} \lambda_j t_j(y_{i-1}, y_i, x, i) + \sum_{i,l} \mu_l s_l(y_i, x, i) \right] \quad (5)$$

式中, $Z(x)$ 为用于归一化, 形成概率值; t_j 为 i 处的转移特征, 相应的权重为 λ_j ; s_l 为 i 处的状态特征, 相应的权重为 μ_l ; j, l 分别为转移特征和状态特征函数的编号。

通常, 特征函数 t_j 和 s_l 的取值 1 或 0, 当满足特征条件时的取值 1, 否则为 0, 其公式如下:

$$t_j(y_{i-1}, y_i, x, i) = \begin{cases} 1 \\ 0 \end{cases} \quad (6)$$

$$s_l(y_{i-1}, y_i, x, i) = \begin{cases} 1 \\ 0 \end{cases} \quad (7)$$

鉴于状态特征函数 s_l 的权重为 μ_l , 转移特征 t_j 的权重为 λ_j , 则权重参数 $Z(x)$ 的 CRF 模型就得到了。训练 CRF 模型时, 使用式 (8) 给每个词打分, 条件概率模型 $P(Y|X)$ 则根据最大似然估计计算。在实际预测过程中, 对于 1 个给定的观察序列, 利用维特比算法计算最大标签序列:

$$\text{score} = \sum_{i,k} \lambda_k f_k(y_{i-1}, y_i, x, i) \quad (8)$$

4 实验与结果分析

4.1 数据集划分

以某煤矿事故隐患文本数据集为例, 选取煤矿事故隐患历史数据中隐患描述字段中的文本数据作为煤矿事故隐患知识实体抽取语料集。为确保样本数据的代表性和可靠性, 同时充分反映不同年度样本数据的分布特征, 采用分层抽样手段从语料集中选取了 2 000 条数据。鉴于实验数据集规模相对较小, 使用较大的测试集比例可以保证测试数据的充分性和代表性, 以更准确地评估模型的性能。因此, 将实验数据集划分为 80% 的训练集和 20% 的测试集。按照所提出的标注方法, 构建了煤矿隐患文本知识实体抽取任务的实验数据集。

4.2 模型评估方法及指标选取

4.2.1 K 折交叉验证法

K 折交叉验证法是将实验训练集划分成 K 份, 将其中的 K-1 份作为训练集, 剩余的 1 份作为验证集, 循环 K 次训练。这种划分方式是将全部的实验训练数据都作为验证集, 最后通过 K 份的平均得到最终模型验证的评价指标值。以 5 折交叉验证为例, 数据集划分如图 4 所示。 $M_1 \sim M_5$ 是不同轮次训练得到的模型评价指标值, 最终模型的评价指标值是 5 个轮次训练的平均值。

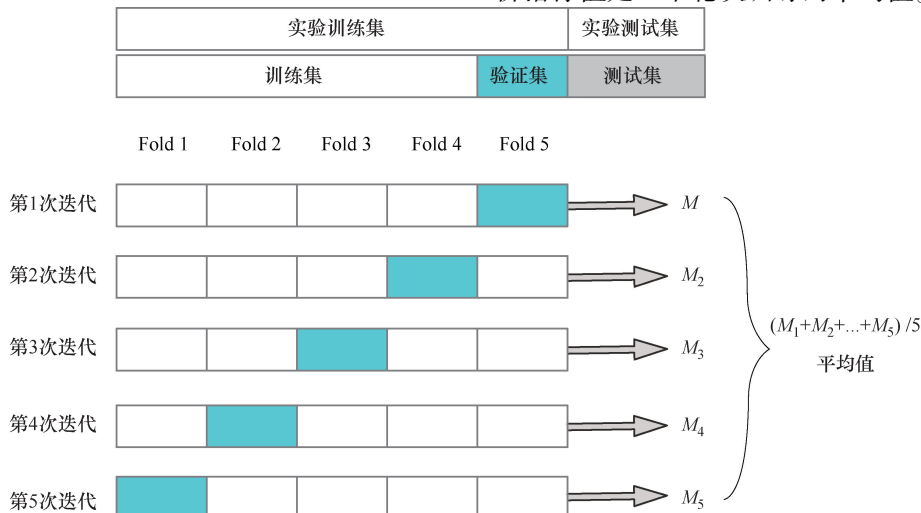


图 4 5 折交叉验证数据集划分

Fig. 4 Division of 5-fold cross-validation dataset

4.2.2 基于牛顿冷却定律的概率融合算法

根据牛顿冷却定律,物体的冷却速度与其当前温度和室温之间的温差成正比:

$$T'(t) = -\alpha [T(t) - H] \quad (9)$$

式中, $T(t)$ 为温度的时间函数; H 为室温, $^{\circ}\text{C}$; α 为室温与降温速率之间的比例关系,即冷却系数^[18]。

式(9)为微分方程,为了计算当前温度,求解该微分方程后可得 $T(t)$ 的函数表达式:

$$T = T_0 e^{-\alpha(t-t_0)} \quad (10)$$

式中, $t-t_0$ 为时间间隔, s ; T_0 为时间 t_0 时的温度, $^{\circ}\text{C}$; T 为当前温度, $^{\circ}\text{C}$ 。

所构建的模型在煤矿隐患文本知识实体抽取任务的实验数据集上经过 K 折交叉验证的训练后,得到 K 个模型,且具有 K 个不同的 F_1 值。通过对 F_1 值进行降序排列并计算 K 个模型所占的权重比值后,对各模型进行加权求和得到最优模型的参数,从而获得更优的评价指标值,进一步提升模型性能。可以把 F_1 值的排名想象成一个自然冷却的过程,即有如下假设:任一时刻, K 个模型中所有的 F_1 值,都有一个当前温度,温度最高的 F_1 值排在第一位;随着时间的流逝,所有模型的 F_1 值的温度都逐渐冷却。

基于上述假设,借鉴牛顿冷却定律原理,建立温度与时间之间的函数关系,构建一个权重指数式衰减的过程,计算不同模型的权重,最后对所有模型进行加权求和。模型融合计算的主要步骤如下。

(1)将 K 个模型的 F_1 值进行降序排列。

(2)计算不同模型的权重值:

$$\lambda_n = e^{-\alpha(n-1)}, \quad n = 1, 2, \dots, k \quad (11)$$

式中, λ_n 为第 n 个模型的权重值; α 为权重衰减系数,取值0.25; k 为模型的数量。

(3)对 K 个模型计算模型权重后,再进行加权求和获得融合后的模型。

4.2.3 评估指标选取

参考SIGHAN^[19]定义的评价指标来评估所提出的知识提取模型效果,具体包括精度 P (Precision)、召回率 R (Recall)和 F_1 值(F_1 Score)3个指标:

$$P = \frac{TP}{TP + FP} \times 100\% \quad (12)$$

$$R = \frac{TP}{TP + FN} \times 100\% \quad (13)$$

$$F_1 = \frac{2PR}{P + R} \times 100\% \quad (14)$$

式中,TP为模型正确预测出的知识实体正样本数

量;FP为模型错误预测出的知识实体正样本数量;FN为模型错误预测出的知识实体负样本数量。

4.3 实验环境及参数设定

煤矿事故隐患文本知识实体抽取实验分析研究主要基于PyTorch框架结构进行,实验环境为:Ubuntu 20.04.4 LTS操作系统,IntelXeon(R) Silver 4114 CPU,64GB内存,NVIDIA TITAN XP 12G×2显卡,CUDA Version 11.4,Python版本3.7,Torch版本1.13.1。选择CRoBERTa-ext(chinese_roberta_wwm_ext_L-12_H-768_A-12)作为实验分析的预训练语言模型。该模型是在chinese_roberta_L-12_H-768的基础上进行改进和扩展得到的,在大规模中文语料上进行预训练,具备提取丰富语义和上下文信息的能力^[20]。

K 折交叉验证折数的选取中,5、7和10折较为常用。实验选择 $K=7$ 。模型中其他主要参数的设置见表4。

表4 实验参数设定

Table 4 Setting of experimental parameters

参数名称	参数值	备注
Warmup_proportion	0.05	学习率预热比例
Keep_prob	0.90	保留率
Dropout	0.80	丢弃率
Decay_rate	0.85	学习率衰减率
Decay_step	200	学习率衰减步数
Train_epoch	30	训练轮次
Max_x_length	300	输入文本最大长度
Embed_learning_rate	5×10^{-5}	BERT模型微调学习率
Batch_size	8	每批次的训练样本数
IDCNN_filter_width	3	卷积核
IDCNN_dilation	1,1,2	膨胀宽度

4.4 结果分析

先后采用4组实验验证、评估所提出的模型在煤矿隐患文本知识实体抽取任务中的有效性和可行性。1~3组为在不引入基于牛顿冷却定律的概率融合算法的情况下,将BERT-BILSTM-CRF、BERT-IDCNN-CRF、动态权重融合的BERT-BILSTM-CRF和动态权重融合的BERT-IDCNN-CRF模型分别进行煤矿隐患文本知识实体抽取任务的性能对比。第4组是在1~3组实验的基础上,引入基于牛顿冷却定律的概率融合算法,比较分析各模型在测试集上的性能情况。

(1)实验1为4种模型在煤矿隐患文本知识实体抽取任务训练集上的平均每个轮次(Epoch)

的训练耗时(表5)与损失值的变化情况对比(图5)。

表5 训练耗时对比

Table 5 Comparison table of training time consumption

编号	模型	平均每一轮次训练时间/s
模型1	BERT-BILSTM-CRF	4
模型2	BERT-IDCNN-CRF	424
模型3	动态权重融合的 BERT-BILSTM-CRF	502
模型4	动态权重融合的 BERT-IDCNN-CRF	448

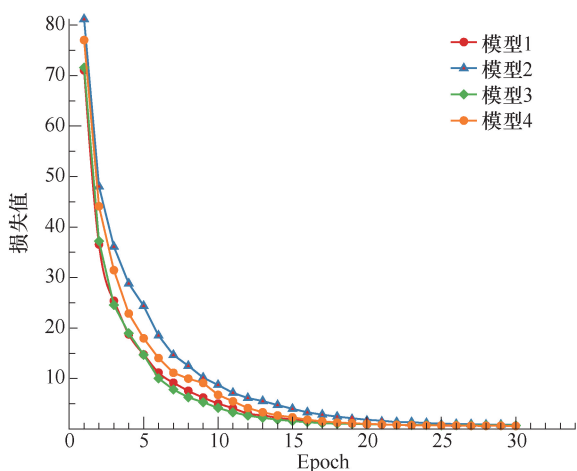


图5 各模型在训练集上的性能比较

Fig. 5 Comparison of each model's performance on training sets

由表5可知,模型架构和设计对训练时间会产生一定的影响,即引入动态权重融合的模型相对于未引入时在训练时间上有所增加,是因为动态权重融合机制需要更多的计算和迭代进行权重的调整和优化。引入IDCNN结构的模型训练时间相对较短,主要得益于IDCNN结构的并行特性和高效计算能力,使得模型能够更快地进行特征提取和参数更新。另外,从图5可看出,不同模型在30个Epoch内损失值基本保持不变,可以判定各模型在30个Epoch内已达到收敛稳定的状态。

(2) 实验2为4种模型在煤矿隐患文本知识实体抽取任务验证集上的 F_1 值对比(图6)。从图6可看出,随着训练步数的变化,模型4在验证集上的 F_1 值要优于其他模型。根据表5、图6可知,引入动态权重融合的模型4的训练时间增加及损失收敛速度减缓了,但是提升了知识实体抽取任务模型中的 F_1 值。

(3) 实验3为实验1中的4种模型在煤矿隐患

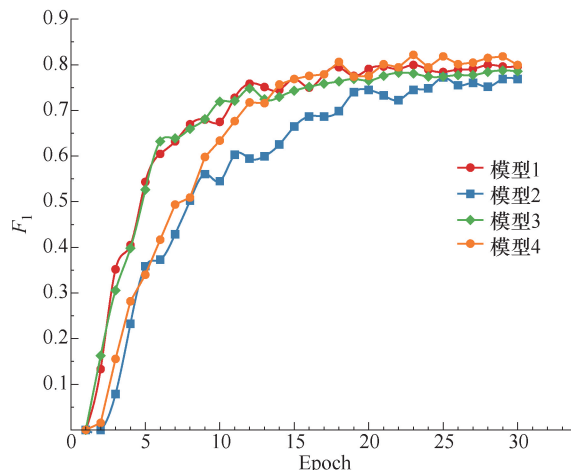


图6 各模型在验证集上的 F_1 值对比

Fig. 6 Comparison of the F_1 score of each model on validation sets

文本知识实体抽取任务测试集上的性能对比(表6)。在煤矿事故隐患领域同一数据集的知识实体识别任务中,相较于未引入动态权重融合的模型1,引入后的模型3的召回率虽然降低了0.13%,但精度、 F_1 值分别提升了0.35%、0.23%;同样地,相较于未引入动态权重融合的模型2,引入后的模型4的精度提升了0.52%,而召回率与 F_1 值分别降低了0.37%、0.26%。这可以说明引入动态权重融合的BERT模型在特征抽取能力方面更为出色。引入动态权重融合的模型能够充分利用BERT模型Transformer编码器的信息,并选择最佳的特征表示,更准确地捕获语义信息,进一步提升了知识实体抽取任务性能。

表6 实验结果比较

Table 6 Comparison of experimental results

编号	引入牛顿冷却定律的概率融合算法	模型	P/%	R/%	F_1 /%
1	否	BERT-BILSTM-CRF	85.65	78.04	77.60
2	否	BERT-IDCNN-CRF	85.59	78.69	78.14
3	否	动态权重融合的 BERT-BILSTM-CRF	86.00	77.91	77.83
4	否	动态权重融合的 BERT-IDCNN-CRF	86.11	78.32	77.88
5	是	BERT-BILSTM-CRF	93.60	83.36	84.29
6	是	BERT-IDCNN-CRF	94.05	83.85	84.80
7	是	动态权重融合的 BERT-BILSTM-CRF	93.06	83.60	84.40
8	是	动态权重融合的 BERT-IDCNN-CRF	95.04	83.60	85.39

(4) 实验4是将实验1中的4种模型分别引入基于牛顿冷却定律的概率融合算法,并对比了在测试集上的性能(表6)。相较于未引入概率融合算法的模型,引入了的模型的精度、召回率与 F_1 值都有较大的提升。动态权重融合的BERT-IDCNN-CRF在引入概率融合算法后,具有最佳的性能,精度、召回率与 F_1 值分别提升了8.93%、5.28%、7.51%。这进一步验证了基于牛顿冷却定律的概率融合算法可以提供更准确的概率估计信息,并综合考虑了多个模型的预测结果,从而提高模型性能;同时,动态权重融合的BERT-IDCNN-CRF与基于牛顿冷却定律的概率融合算法的组合不仅提高了模型的整体性能,还在一定程度上提高了模型的稳定性,使得模型在煤矿隐患文本数据集上都具有较好的适应性。

综上所述,动态权重融合的BERT-IDCNN-CRF引入基于牛顿冷却定律的概率融合算法模型在实验中显示出了较稳定的性能表现,在多个评估指标上表现最优,进一步验证了在煤矿事故隐患文本知识实体抽取方面的优势和可靠性。

5 结论

本文建立了煤矿事故隐患描述文本知识实体标注规范,并基于分层抽样和专业工具构建了高质量的标注数据集,得到的主要成果如下:

(1)建立了煤矿事故隐患描述文本知识实体标注规范。在采用分层抽样方法选取煤矿事故隐患文本原始数据后,使用Brat工具进行了文本数据标注,完成煤矿事故隐患文本标注数据集的构建。

(2)提出了一种基于动态权重融合的BERT-IDCNN-CRF模型。为评估模型性能,采用了 K 折交叉验证法,在有效避免模型过拟合的同时,能够保证模型的泛化能力。

(3)引入了基于牛顿冷却定律的模型融合算法,进一步提升了煤矿事故隐患文本知识实体抽取任务的识别效果。为煤矿事故隐患文本特征提取提供了一种新的解决思路,进一步提高了煤矿事故隐患文本数据挖掘能力,对于煤矿事故的预防和处理具有重要的实际意义。

参考文献

[1] 宋曦,丁文梅,宁云才,等. 煤矿安全生产管理体系优化研究——以陕西某煤矿为例[J]. 矿业科学学

报,2019,4(2):187-194.

SONG Xi, DING Wenmei, NING Yuncai, et al. Research on optimization of coal mine safety production management system—take a coal mine in Shaanxi as an example[J]. Journal of Mining Science and Technology, 2019, 4(2): 190-194.

[2] 王美君,谭章禄,李慧园,等. 智能化煤矿数据治理能力评估与提升策略研究[J]. 矿业科学学报, 2024, 9(1): 106-115.

WANG Meijun, TAN Zhanglu, LI Huiyuan, et al. Research on evaluation and promotion strategy of data governance capability for intelligent coal mines [J]. Journal of Mining Science and Technology, 2024, 9(1): 106-115.

[3] 汪诚愚,何晓丰,宫学庆,等. 面向上下位关系预测的词嵌入投影模型[J]. 计算机学报, 2020, 43(5): 868-883.

WANG Chengyu, HE Xiaofeng, GONG Xueqing, et al. Word embedding projection models for hypernymy relation prediction [J]. Chinese Journal of Computers, 2020, 43(5): 868-883.

[4] 侯运炳,陈袖龙,王雅先,等. 采场全生命周期的矿压危害事件知识图谱表示方法研究[J]. 矿业科学学报, 2024, 9(2): 295-303.

HOU Yunbing, CHEN Youlong, WANG Yaxian, et al. Research on knowledge graph representation method of mine pressure hazard events in stope whole life cycle [J]. Journal of Mining Science and Technology, 2024, 9(2): 295-303.

[5] MORWAL S, JAHAN N, CHOPRA D. Named entity recognition using hidden markov model (HMM) [J]. International Journal on Natural Language Computing, 2012, 1(4): 15-23.

[6] MCCALLUM A, FREITAG D, PEREIRA F. Maximum entropy markov models for information extraction and segmentation [J]. International Conference on Machine Learning, 2001, 17(1): 591-598.

[7] LAFFERTY J, MCCALLUM A, PEREIRA F, et al. Conditional random fields: probabilistic models for segmenting and labeling sequence data [C]. Proceedings of the Eighteenth International Conference on Machine Learning, Massachusetts, 2002, 1: 282-289.

[8] DEVLIN J, CHANG M W, LEE K, et al. BERT: pre-training of deep bidirectional transformers for language understanding [J]. arXiv: 1810.04805, 2019, 1-16.

[9] 鹿晓龙. 煤矿安全知识图谱构建技术研究 [D]. 徐州: 中国矿业大学, 2021.

LU Xiaolong. Study on construction technology of coal mine safety knowledge map [D]. Xuzhou: China Uni-

versity of Mining and Technology, 2021.

- [10] 赵彭彭. 基于深度学习的煤矿事故领域命名实体识别方法研究与应用[D]. 太原: 太原科技大学, 2021.
- ZHAO Pengpeng. Research and application of named entity recognition method in coal mine accident field based on deep learning[D]. Taiyuan: Taiyuan University of Science and Technology, 2021.
- [11] 李毅中. 安全生产事故隐患排查治理暂行规定[J]. 中华人民共和国国务院公报, 2008, 26: 44-47.
- Li Yizhong. Temporary provisions on the investigation and treatment of hidden trouble in safety production accidents[J]. Gazette of the State Council of the People's Republic of China, 2008, 26: 44-47.
- [12] ZHANG C. DeepDive: a data management system for automatic knowledge base construction[D]. Wisconsin: The University of Wisconsin-Madison, 2015.
- [13] NAKAYAMA H, KUBO T, KAMURA J, et al. Doccano: text annotation tool for human[EB/OL]. 2018. <https://github.com/doccano/doccano>.
- [14] HONNIBAL M, MONTANI I. Prodigy: a new tool for radically efficient machine teaching[EB/OL]. 2017. <https://explosion.ai/blog/prodigy-annotation-tool-active-learning>.
- [15] YANG J, ZHANG Y, LI L, et al. YEDDA: A lightweight collaborative text span annotation tool[J]. arXiv: 1711.03759, 2018, 5: 1-6.
- [16] STENETORP P, PYYSALO S, TOPIĆ G, et al. Brat: a web-based tool for NLP-assisted text annotation[C]//Demonstrations at the 13th conference of the european chapter of the association for computational linguistics, Avignon, France, 2012: 102-107.
- [17] YU F, KOLTUN V. Multi-scale context aggregation by dilated convolutions[C]//4th International Conference on Learning Representations, San Juan, Puerto rico, 2016: 1-13.
- [18] KONOVALENKO I, LUDWIG A, LEOPOLD H. Real-time temperature prediction in a cold supply chain based on Newton's law of cooling[J]. Decision Support Systems, 2021, 141: 113451.
- [19] TSENG Y H, LEE L H, CHANG L P, et al. Introduction to sighthan 2015 bake-off for chinese spelling check[C]//Proceedings of the Eighth SIGHAN Workshop on Chinese Language Processing, Beijing, China, 2015, 32-37.
- [20] 侯钰涛, 阿布都克力木·阿布力孜, 哈里旦木·阿布都克里木. 中文预训练模型研究进展[J]. 计算机科学, 2022, 49(7): 148-163.
- HOU Yutao, ABULIZI Abudukelimu, ABUDUKELIMU Halidanmu. Advances in chinese pre-training models[J]. Computer Science, 2022, 49(7): 148-163.

(责任编辑:张彩艳)