

文章编号: 2095-2163(2021)06-0014-06

中图分类号: TP391.1

文献标志码: A

基于 BERT 的突发事件文本自动标注方法

杨芷婷, 马汉杰

(浙江理工大学 信息学院, 杭州 310018)

摘要: 信息提取技术是自然语言处理技术的关键技术之一, 其中最主要的任务是事件元素提取。本文利用深度学习网络模型实现信息提取任务进行了深入研究。训练数据来源于上海大学构建的 CEC 已标注的语料库。相比于采用手工设立规则的识别方式和 BiLSTM 网络模型, 本文通过对数据进行预处理和搭建 BERT-BiLSTM-CRF 深度网络模型, 对文本数据训练实现标注, 在时间、报道时间、参与对象的识别准确率上均有所提升。

关键词: BERT; 中文突发事件; 自动标注; 信息提取

Automatic-annotation method for emergency text corpus based on BERT

YANG Zhiting, MA Hanjie

(School of Information Science and Technology, Zhejiang Sci-Tech University, Hangzhou 310018, China)

[Abstract] Information Extraction is one of the most important technology in Natural Language Process, which mainly job is extract the events element. This paper proposes a deep learning network method to solve this task. The training data comes from CEC corpus which was built by Shanghai University. In this experiment, compared with rule-based annotation method and Bi-LSTM network method, showing that using BERT+BiLSTM+CRF model can improve the efficiency of event extraction effectively.

[Key words] BERT; Chinese emergency event; automatic-annotation; information extraction

0 引言

自然语言处理技术(Natural Language Processing, NLP)是计算机科学、人工智能和语言学交叉的领域, 主要研究用计算机来处理、理解和应用人类语言, 达到人与计算机之间的有效通信。信息抽取为自然语言处理领域的一个重要研究方向。其中, 命名实体识别(Named Entity Recognition, NER)是信息抽取的基础任务, 其任务是从文本中识别出诸如人名、组织名、日期、时间、地点、特定的数字形式等内容, 并为之添加相应的标注信息, 为信息抽取后续工作提供便利^[1]。在实际自然语言处理任务中, 如社交媒体文本处理等, NER 作为上游任务在整个系统中起着举足轻重的作用。

随着网络信息的爆发式增长, 传统的文本分析手段已不适合处理海量突发事件信息, 机器学习和数据挖掘技术才是目前信息抽取任务处理过程中备受青睐的处理技术。随着评测会议, 如: MUC (Message Understanding Conference)^[2]、自动内容抽取(Automatic Content Extraction, ACE)^[3]的举办, 事件抽取技术取得了长足进展。2016 年 Peng 等人将

Chen 等人发表的 SOTA 中文分词系统^[4]与中文媒体 ER 模型结合^[5], 在实体识别训练过程中利用分词训练提供的输出参数训练, 使识别效果提高了 5%^[6]; Lample 在堆叠 LSTM 模型(S-LSTM)基础上, 结合基于字符的表示模型^[7]和词嵌入模型, 在多种语言上得到了较好的训练结果^[8]; 伟峰等人于 2019 年首先提出利用基于注意力机制^[9]的序列标注模型, 联合抽取句子级事件的触发词和实体, 与独立进行实体抽取和事件识别相比, 联合标注的方法在 F 值上提升了 1 个百分点; 武惠^[10]提出基于迁移学习和深度学习的 TrBiLSTM-CRF 模型, 采用实例迁移学习算法将源域知识迁移到目标域, 在小规模数据集上取得了较好的效果。

本文利用深度学习网络搭建学习模型, 以标注的中文事件语料数据为输入, 训练得到自动提取事件信息的网络模型, 该模型主要由 BERT-BiLSTM-CRF 组成。其中, BERT 预训练语言模型由谷歌人工智能团队提出, 能够较完整地保存文本语义信息; BiLSTM-CRF 是较为常见的序列标注模型, 在语音识别、词性标注、实体识别等领域应用广泛。在中文文本事件知识提取领域, 研究语义、推理和挖掘是提

作者简介: 杨芷婷(1997-), 女, 硕士研究生, 主要研究方向: 机器学习、自然语言处理、人工智能; 马汉杰(1982-), 男, 博士, 副教授, 硕士生导师, 主要研究方向: 机器视觉、数据挖掘、视频图像传输与处理。

收稿日期: 2021-02-21

取信息的主要手段之一,在程序开发过程中还需语言领域相关储备知识,对研究人员来说是个不小的挑战,而深度学习网络训练模型和大数据,使得中文信息处理发展向前推进了一大步。

1 相关工作

由于事件文本的特殊性,学者们对不同的事件语料库采用了不同的标注体系^[11],目前影响较大的事件标注语料库有 ACE 测评语料^[12]和 TimeBank 语料^[13],中文事件标注语料比较常用的是中文突发事件语料库(CEC)。CEC 语料库是由上海大学(语义智能实验室)所构建。CEC 语料库根据国务院颁布的《国家突发公共事件总体应急预案》分类体系,收集了 5 类(地震、火灾、交通事故、恐怖袭击和食物中毒)突发事件的新闻报道作为生语料,经过对

生语料进行文本预处理、文本分析、事件标注以及一致性检查等处理,最后将标注结果保存到语料库中,CEC 合计 332 篇^[14]。与 ACE 和 TimeBank 语料库相比,CEC 语料库的规模虽然偏小,但是对事件和事件要素的标注却最为全面。

CEC 标注语料采用 XML 格式,文本格式以 Title、ReportTime、Content、eRelation 标签依序组成,如图 1 所示。Content 标签内容为已标注事件元素的新闻文本。Event 标签主要包含的标签有:事件触发词(Denoter)、事件发生地点(Location)、时间(Time)、对象者(Participant)等。eRelation 标签定义事件之间的关系属性,有 5 种类型值:因果(causal)、伴随(accompany)、跟随(follow)、组成(composite)以及意念包含(thoughtContent)。

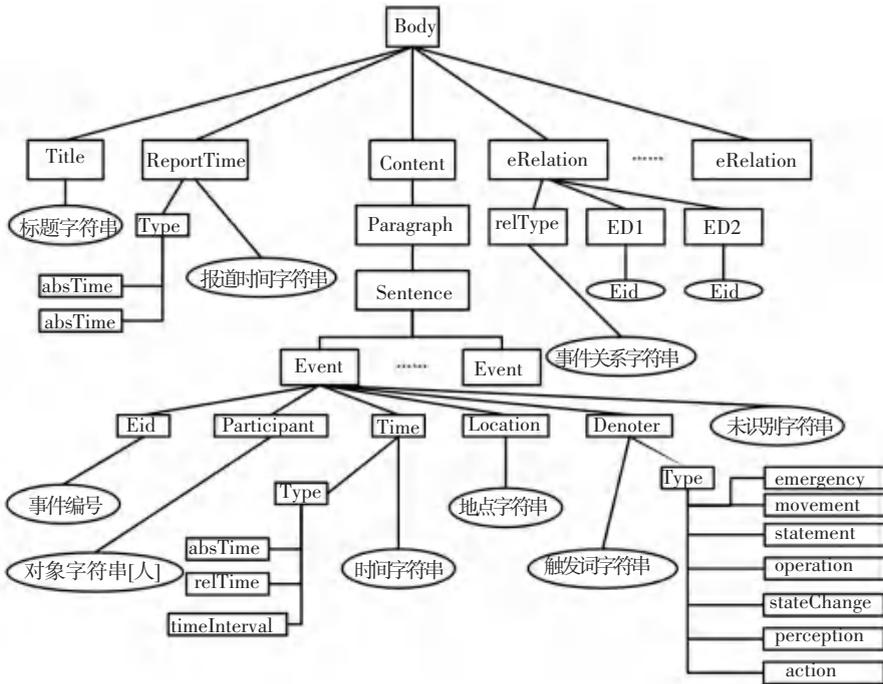


图 1 CEC 语料标签树

Fig. 1 CEC corpus label tree

中文实体识别与英文不同,是以分词为基础(word-based)的训练模型,识别效果与分词准确性相关。如:“南京市长江大桥”若错误地被分割为“南京/市长/江大桥”,则会影响该实体的识别效果。不少专家研究得出,在中文实体识别中使用基于字(character-based)深度识别模型优于基于词的模型。但由于中文的多义性和多态性,单纯依靠字特征将会丢失词语隐藏的信息。因此,如何将基于字的模型和基于词的模型混合得到更好的结果,是当前中文实体识别的一大研究热点^[15]。在数据不

足或特殊文本的情况下,引入语法结构特征和词性特征,通过编写规则,识别事件信息是常用的事件抽取方法^[16]。

本文主要针对的是中文文本的事件识别和事件元素提取研究。利用 BERT-BiLSTM-CRF 模型对 CEC 语料进行训练,提取出相关事件元素。在进行处理之前,需要对事件语料进行预处理,将标注的 xml 文件转换成可训练格式。事件触发词抽取任务要求正确识别触发词并判断触发词赋予正确的类型。对于实体识别,要求对文本中事件触发词

Transformer 编码单元如图 6 所示, 编码单元最主要的模块是自注意力部分计算公式:

$$Attention(Q, K, V) = Softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (8)$$

其中, Q, K, V 为输入向量分别乘以 W^Q, W^K, W^V 矩阵, 经过线性变化得到 d_k 为输入向量维度。经过 d_k 进行缩小之后通过 $softmax$ 归一化得到权重表示, 最后得到句子中所有词向量的带权和, 这样的词向量相较于传统词向量更加具有全局性。

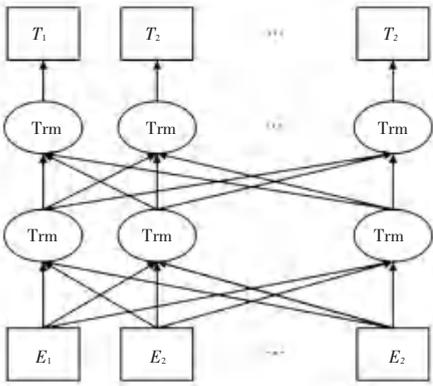


图 5 双向 Transformer 结构

Fig. 5 Bidirectional Transformer structure

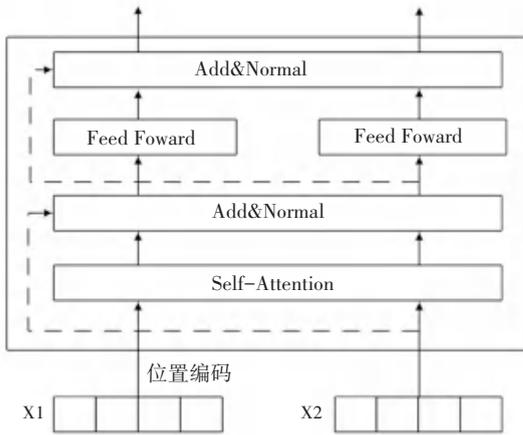


图 6 Transformer 编码单元

Fig. 6 Transformer encoder unit

此外, BERT 模型采用了多头注意力机制 (Multi-Head Attention), 多头自注意力计算过程分为 4 步:

- (1) 输入经过线性变换后生成 Q, K, V 3 个向量;
- (2) 进行分头操作 (假设原始向量维度为 512, 分成 8 个 head 后, 每个 head 维度为 64);
- (3) 每个 head 进行自注意力计算;
- (4) 最后将计算结果拼接起来。公式描述为:

$$head_i = Attention(QW_i^Q, KW_i^K, VW_i^V). \quad (9)$$

$$MultiHead(Q, K, V) = Concat(head_1, head_2, \dots, head_h)W^O. \quad (10)$$

为解决深度学习中的退化问题, 编码单元加入了残差网络和层归一化, 如式 11 所示:

$$Z = LayerNorm(Concat(MultiHead, Sublayer(X))). \quad (11)$$

之后将归一化结果使用 ReLU 做为激活函数, 运算如下:

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2. \quad (12)$$

该模型在预训练中主要包括 2 个任务: Masked 语言模型和 Next 句子预测。在训练过程中, 首先构造句子对。构造方法是在规模文本中, 选择具有上下文关系的句子对, 对其中 50% 的句子对进行随机替换, 使其不具有上下文关系, 然后在“Masked 语言模型”和“Next 句子预测”任务上进行训练, 捕捉词级别和句子级别的表示。

2.3 CRF 模型

CRF 模型^[20]由 Collobert 提出, 相较于 softmax 分类器, 能考虑标签序列的全局信息, 获得更优的标签序列。CRF 是在给定随机变量 X 的条件下, 随机变量 Y 的马尔科夫随机场, 在序列预测问题常用的是线性链马尔科夫随机场。对于观察序列 $x = (x_1, x_2, \dots, x_n)$ 和状态序列 $y = (y_1, y_2, \dots, y_n)$, 利用 $Softmax$ 归一化后的概率如式 (13) 所示:

$$P(y | x) = \frac{1}{Z(x)} \exp\left\{ \sum_{i,k} \lambda_k f_k(y_{i-1}, y_i, x, i) + \sum_{i,j} \mu_j h_j(y_i, x, i) \right\}. \quad (13)$$

其中, f_k 是转移特征函数; h_j 是状态特征函数; λ_k, μ_j 是对应的权值; $Z(x)$ 是归一化因子, 计算公式如下:

$$Z(x) = \sum_y \exp\left\{ \sum_{i,k} \lambda_k f_k(y_{i-1}, y_i, x, i) + \sum_{i,j} \mu_j h_j(y_i, x, i) \right\}. \quad (14)$$

模型在预测过程 (解码) 时, 使用动态规划的 Viterbi 算法来求解最优路径, 如式 (15) 所示:

$$y^* = \operatorname{argmax} P(y | x). \quad (15)$$

3 实验过程

3.1 实验数据

本文使用的是上海大学公开的中文突发事件语料库, 是一个一个小规模的事件语料库, 合计 332 篇。语料文本包括: 地震、火灾、交通事故、恐怖袭击、食物中毒等 5 类。其中各类文本数据统计见表 1。

表 1 CEC 语料文本数据统计

Tab. 1 CEC corpus text data statistics

类型	文章篇数	句子数	无事件句子数	事件数	触发词	事件元素
地震	62	401	41	1 002	1 002	2 461
火灾	75	433	39	1 216	1 216	2 935
交通事故	85	514	9	1 802	1 802	4 186
恐怖袭击	49	324	38	823	823	2 042
食物中毒	61	392	17	1 111	1 111	2 777
SUM	332	2 064	114	5 954	5 954	14 401

对于 CEC 语料的处理方式可以参考实体识别训练时使用的 BIO 三段标记法:对于每个需要识别的标签,将第一个字标记为“B-(实体类别)”,后续标记为“I-(实体类别)”,对于无关字,一律标记为 O。以下面一段 xml 格式文本为例做一说明:

```
<Event eid="e5">
```

```
<Participant sid="s5">学校</Participant>等单位也立即
```

```
<Denoter type="movement" did="d5">疏散</Denoter>
```

```
<Participantoid="o5" sid="s6">人群</Participant>到
```

```
<Location type="destination" lid="l5">安全场所</Location>。
```

```
</Event>
```

经过转换之后变为如下形式(文字斜杠后表示标注序列):

学/B-PP 校/I-PP 等/O 单/O 位/O 也/O 立/O 即/O 疏/B-DNT 散/I-DNT 到/O 安/B-LOC 全/I-LOC 场/I-LOC 所/I-LOC。/O

BIO 三段记法最大的优点是支持逐字标记,减少了系统因分词而产生的误差。标记好的数据有 O、B-PP、I-PP、B-LOC、I-LOC、B-TM、I-TM、B-DNT、I-DNT、B-OJ、I-OJ、B-RT、I-RT、X、[CLS] 和 [SEP] 共 16 大类;[CLS] 为句子开始标志,[SEP] 为句子结尾标志。对于每一类实体的识别效果,采用精确率(P)、召回率(R) 和 F 值作为模型性能的评价标准,具体计算公式如下:

$$P = \frac{\text{正确识别出的命名实体个数}}{\text{识别出的命名实体个数}} \times 100\%, \quad (16)$$

$$R = \frac{\text{正确识别出的命名实体个数}}{\text{标准结果中命名实体个数}} \times 100\%, \quad (17)$$

$$F1 = \frac{2PR}{P + R} \times 100\%. \quad (18)$$

3.2 模型参数

本文采用由谷歌人工智能团队开发的 Tensorflow 框架搭建模型,BERT 预训练语言模型默认采用 12 头注意力机制,每次读取序列长度为 128,预训练词长度为 768;训练批次为 16,优化器采用 Adam,学习率设置为 10^{-5} 。LSTM 隐藏单元设为 128,为解决梯度消失和爆炸问题,设置丢弃率为 0.5,采用梯度裁剪技术,clip 设置为 5。由双向 LSTM 网络输出得到的 256 维字向量,经过压缩为 16 维向量作为 CRF 层的输入。

3.3 实验结果及分析

将训练数据从 xml 格式转换为适合的训练数据后,得到标签统计见表 2。

表 2 训练数据标签统计

Tab. 2 Train data labels statistics

标签	个数
Denoter	5 953
Time	1 406
Location	1 658
Participant	3 309
ReportTime	332
Object	2 013

由于训练数据量较少,将训练数据按照 8:2 分割为训练集和测试集。训练得到的网络模型在测试集上各类实体的识别率与 BiLSTM+CRF 神经网络识别率对比结果见表 3。

表 3 实验识别结果

Tab. 3 Label element recognition experiment results

事件元素	本文方法			BiLSTM+CRF		
	P	R	$F1$	P	R	$F1$
Denoter	0.76	0.80	0.78	0.63	0.60	0.61
Location	0.57	0.65	0.61	0.51	0.55	0.53
Object	0.51	0.56	0.53	0.33	0.23	0.27
Participant	0.59	0.70	0.64	0.65	0.56	0.60
Time	0.79	0.85	0.82	0.79	0.73	0.76
合计	0.67	0.74	0.70	0.61	0.56	0.58

由此可知,相较于使用训练好的维基百科字向量+BiLSTM+CRF模型,使用BERT模型得到的训练结果在各个元素上皆优于该模型,尤其在对象元素识别上得到了近20%的提高。再对比文献[24],采用人工语法规则自动标注得到的实验结果见表4。

表4 人工语法规则自动标注结果

Tab. 4 Experiment results of using grammar rule

命名实体	精确率	召回率	F1值
Denoter	0.74	0.89	0.81
Location	0.64	0.67	0.66
Participant	0.57	0.50	0.57
ReportTime	0.94	0.94	0.94
Time	0.68	0.80	0.74

可以看出,使用机器学习的方法在时间、报道时间、参与对象3类实体的识别准确率、召回率和F1值上均有所提高,在发生地点、触发词的识别率上稍有降低。由此表明,在利用BERT-BiLSTM-CRF模型基础上,确实可以提高部分实体识别的精确率,避免了对实验文本的语法规则和人工实现过滤规则等耗费时力的操作,但在触发词等实体类的识别上稍显劣势,这也是今后需要研究改进的地方。

4 结束语

本文利用BERT-BiLSTM-CRF深度学习模型,对CEC语料库进行自动化标注,提高了标注效率。与传统手工标注方法相比极大的提高标注速度,即使在识别准确率不高的情况下也可人工调整,有利于大规模语料标注工作。对比BiLSTM-CRF网络模型,在事件各个要素识别上都取得较为明显的优化。本文实验模型还存在改进的地方,如在无明显规则的事件触发词、事件参与对象等实体识别的效果并不理想,这是由于事件对象短语在事件句中并没有较为明显的规律特征,需要结合中文语法特征进一步发掘有效识别规则,有待进一步研究。

参考文献

[1] 张晓艳,王挺,陈火旺. 命名实体识别研究[J]. 计算机科学, 2005(4):44-48.

[2] Wikipedia: Message understanding Conference[EB/OL]. 2013-12-27. https://en.wikipedia.org/wiki/Message_Understanding_Conference

[3] Wikipedia: Message Understanding Conference[EB/OL]. 2013-12-27. https://en.wikipedia.org/wiki/Automatic_content_extraction

[4] CHEN X, X QIU, ZHU C, et al. Long Short-Term Memory Neural Networks for Chinese Word Segmentation [C]// Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. 2015.

[5] PENG N, DREDZE M. Named Entity Recognition for Chinese Social Media with Jointly Trained Embeddings[C]// Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. 2015.

[6] PENG N, DREDZE M. Improving Named Entity Recognition for Chinese Social Media with Word Segmentation Representation Learning[C]// Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers). 2016.

[7] LING W, T LUÍS, MARUJO L, et al. Finding Function in Form: Compositional Character Models for Open Vocabulary Word Representation[J]. Computer Science, 2015:1899-1907.

[8] LAMPLE G, BALLESTEROS M, SUBRAMANIAN S, et al. Neural Architectures for Named Entity Recognition [C]// Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics; Human Language Technologies. 2016.

[9] 仲伟峰,杨航,陈玉博,等. 基于联合标注和全局推理的篇章级事件抽取[J]. 中文信息学报,2019,33(9):88-95,106.

[10] 武惠,吕立,于碧辉. 基于迁移学习和BiLSTM-CRF的中文命名实体识别[J]. 小型微型计算机系统,2019,40(6):1142-1147.

[11] 赵军,刘康,周光有. 等开放式文本信息抽取[J]. 中文信息学报,2011,25(6):98-110.

[12] The automatic content extraction (ACE) program-tasks, data, and evaluation [C]// International Conference on Language Resources & Evaluation. 2004.

[13] PUSTEJOVSKY J, HANKS P, R SAURÍ, et al. The TimeBank corpus[J]. proceedings of corpus linguistics, 2003.

[14] 付剑锋. 面向事件的知识处理研究[D]. 上海大学, 2010.

[15] JIN Y, XIE J, GUO W, et al. LSTM-CRF Neural Network with Gated Self Attention for Chinese NER [J]. IEEE Access, 2019(99):1.

[16] 张建权. 基于CNN和BiGRU-attention的互联网敏感实体识别方法[J]. 网络安全技术与应用,2020(4):61-65.

[17] DEVLIN J, CHANG M W, LEE K, et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding [J]. 2018.

[18] 李妮,关焕梅,杨飘,等. 基于BERT-IDCNN-CRF的中文命名实体识别方法[J]. 山东大学学报(理学版),2020,55(1):102-109.

[19] VASWANI A, SHAZEER N, PARMAR N, et al. Attention is all you need [C]// Advances in Neural Information Processing Systems. 2017: 5998-6008.

[20] LAFFERTY J. Conditional random fields: Probabilistic models for segmenting and labeling sequence data [C]// Proceedings of the 18th International Conference on Machine Learning, 2001. Morgan Kaufmann, 2001: 282-289.

[21] 刘炜,王旭,张雨嘉,等. 一种面向突发事件的文本语料自动标注方法[J]. 中文信息学报,2017,31(2):76-85.