



(12) 发明专利

(10) 授权公告号 CN 113705218 B

(45) 授权公告日 2023.03.21

(21) 申请号 202111031194.8

G06F 40/30 (2020.01)

(22) 申请日 2021.09.03

(56) 对比文件

(65) 同一申请的已公布的文献号
申请公布号 CN 113705218 A

CN 112084381 A, 2020.12.15

CN 112000792 A, 2020.11.27

CN 113051887 A, 2021.06.29

(43) 申请公布日 2021.11.26

CN 111626056 A, 2020.09.04

CN 111159336 A, 2020.05.15

CN 111444726 A, 2020.07.24

(73) 专利权人 四川大学

地址 610065 四川省成都市武侯区一环路
南一段24号

魏优等.基于深层语境词表示与自注意力的
生物医学事件抽取.《计算机工程与科学》.2020,
第42卷(第09期),

严红等.基于深度神经网络的法语命名实
体识别模型.《计算机应用》.2019,第39卷(第5
期),1288-1292.

(72) 发明人 陈兴蜀 蒋梦婷 袁磊 刘朋
黄铁脉 廖志红 宋可儿 冯科
王海舟 王文贤 罗永刚

审查员 彭一

(74) 专利代理机构 成都禾创知家知识产权代理
有限公司 51284

专利代理师 刘凯

(51) Int. Cl.

G06F 40/279 (2020.01)

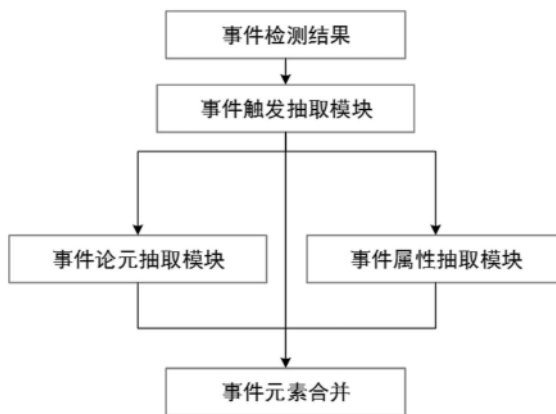
权利要求书2页 说明书7页 附图4页

(54) 发明名称

基于字符嵌入的事件元素网格化抽取方法、
存储介质及电子装置

(57) 摘要

本发明公开了一种基于字符嵌入的事件元素网格化抽取方法、存储介质及电子装置,方法包括以下步骤:首先构建事件元素抽取基础模型,包括基于BERT的信息预学习表示层、字编码嵌入层、BiGRU双向门控循环神经网络层、自注意力层和CRF条件随机场输出层,并将基础模型按功能细化分为事件触发词抽取、事件论元抽取和事件属性抽取3个网格模块;然后分别对事件触发词抽取模型、事件论元抽取模型和事件属性抽取模型进行抽取优化;最后采用训练得到的事件元素抽取模型对测试数据进行事件元素抽取的预测。本发明方法在事件元素抽取任务中表现良好,获得了较高的准确率。



1. 一种基于字符嵌入的事件元素网格化抽取方法,其特征在于,包括以下步骤:

步骤1:构建事件元素抽取基础模型;

所述基础模型为字符嵌入神经网络深度学习模型,包括基于BERT的信息预学习表示层、字编码嵌入层、BiGRU双向门控循环神经网络层、自注意力层和CRF条件随机场输出层;其运行步骤如下:

步骤1.1:基于BERT的信息预学习表示层对样本数据上下文语义特征进行预学习,生成突发元事件域的文本表示模型;

步骤1.2:字编码嵌入层利用训练好的BERT语言模型中生成的语义表示向量输入到BiGRU双向门控循环神经网络层;

步骤1.3:BiGRU双向门控循环神经网络层提取输入序列的上下文依赖的长距离深层特征;步骤1.4:自注意力层对BiGRU双向门控循环神经网络层学习到的深层特征进行加权变换,以突出文本序列中重点词汇信息;

步骤1.5:CRF条件随机场输出层将触发词提取转化为序列标注任务,以解决传统词向量转化为字向量后上下文标注信息问题;

将所述基础模型按功能细化分为事件触发词抽取、事件论元抽取和事件属性抽取3个网格模块,以用于后续步骤根据不同模型的功能特性,分别对模型进行优化;

步骤2:对事件触发词抽取模型进行抽取优化:从一个事件句抽取一个主事件触发词作为事件触发词,多余事件触发词作外部特征,辅助表征主事件;将标注数据中的所有事件触发词作为知识库,作为事件触发词抽取的先验特征;抽取句子中与知识库内事件触发词匹配的触发词,并标注出来,与根据句子BERT语义表示模型获得的字符嵌入向量拼接;并将事件类型向量拼接到字符嵌入向量中;事件触发词抽取任务的目标向量由事件触发词的抽取结果表示,对应事件触发词的标注结果;

步骤3:对事件论元抽取模型进行抽取优化:在原文本BERT语义特征之上,将文本中所有字符到事件触发词的相对距离作为文本结构特征,事件触发词本身的相对距离为0;并将事件主体与客体联合,事件时间与地点联合,采用两个独立的模型进行提取;事件论元抽取任务的目标向量对应事件论元的抽取结果标注;

步骤4:对事件属性抽取模型进行抽取优化:定义事件属性为事件时态和事件极性,模型输出转变为多分类问题,将所述基础模型中CRF条件随机场输出层置换构造两个分类器;将事件触发词及左右两端进行动态池化获得的特征作为全局特征,与根据句子BERT语义表示模型获得的字符嵌入向量拼接,并采用十折交叉验证方法进行优化;

步骤5:采用步骤1-步骤4训练得到的事件要素抽取模型对测试数据进行事件元素抽取结果的预测。

2. 根据权利要求1所述的基于字符嵌入的事件元素网格化抽取方法,其特征在于,所述步骤1.3中,BiGRU双向门控循环神经网络层同时包含一个正向的 $\overline{GRU}(x_0, x_1, \dots, x_t, \dots, x_{n-1})$ 和一个反向的 $\overline{GRU}(x_{n-1}, x_{n-2}, \dots, x_t, \dots, x_0)$,正向GRU捕捉0:t的上文特征信息 a_t ,反向GRU捕捉t:n-1的下文特征信息 a'_t ,通过拼接所捕获的上下文特征信息,获得句子的上下文信息 y_t ,如公式(1)-(3)所示:

$$\overline{GRU}(x_0, x_1, \dots, x_t, \dots, x_{n-1}) = (a_0, a_1, \dots, a_t, \dots, a_{n-1}) \quad (1)$$

$$\overrightarrow{GRU}(x_{n-1}, x_{n-2}, \dots, x_t, \dots, x_0) = (a'_{n-1}, a'_{n-2}, \dots, a'_t, \dots, a'_0) \quad (2)$$

$$y_t = [a_t, a'_t] \quad (3)$$

式中, x_t 表示词序列特征向量; a_t 表示正向GRU捕捉0:t的上文特征信息; a'_t 表示反向GRU捕捉t:n-1的下文特征信息; y_t 表示获得句子的上下文信息;

所述加权变换公式(4)所示:

$$e_{ij} = a(W \overrightarrow{y_i}, W \overrightarrow{y_j}) \quad (4)$$

式中, e_{ij} 表示句子j的特征对句子i的重要性; a 表示注意力机制; W 表示共享参数的线性变换权重矩阵; y_i 和 y_j 分别表示获得的句子i和句子j的上下文信息。

3. 根据权利要求1所述的基于字符嵌入的事件元素网格化抽取方法, 其特征在于, 所述步骤2中, 在事件触发词抽取任务中, 假设词 w_i 的事件触发词类型目标向量为 $[tri_0, tri_1, tri_2, \dots, tri_j, \dots, tri_n]$, 则 tri_j 的设置方式如公式(5)所示:

$$tri_j = \begin{cases} 0, & \text{其他} \\ 1, & \text{如果 } w_i \text{ 的事件触发词类别标签ID为 } j \end{cases} \quad (5)$$

所述步骤3中, 在事件论元抽取任务中, 假设词 w_i 的事件论元类型目标向量为 $[arg_0, arg_1, \dots, arg_j, \dots, arg_n]$, 其中 arg_j 的设置方式如公式(6)所示:

$$arg_j = \begin{cases} 0, & \text{其他} \\ 1, & \text{如果 } w_i \text{ 的事件论元类别标签ID为 } j \end{cases} \quad (6)$$

4. 一种存储介质, 其特征在于, 所述存储介质中存储有计算机程序, 其中, 所述计算机程序被设置为运行时执行所述权利要求1至3任一项中所述的方法。

5. 一种电子装置, 包括存储器和处理器, 其特征在于, 所述存储器中存储有计算机程序, 所述处理器被设置为通过所述计算机程序执行所述权利要求1至3任一项中所述的方法。

基于字符嵌入的事件元素网格化抽取方法、存储介质及电子装置

技术领域

[0001] 本发明涉及事件抽取技术领域,具体涉及一种基于字符嵌入的事件元素网格化抽取方法、存储介质及电子装置。

背景技术

[0002] 信息抽取技术是把关注的非结构化数据信息从海量文本数据中抽取出来,并转换为结构化的数据。通过信息抽取技术,可以过滤低价值的信息内容,快速获得精准和高质量信息。事件是信息的一种重要表达形式,信息抽取领域的重点研究方向即为事件抽取。该研究中的权威学会ACE (Automatic Content Extraction) 对事件抽取作了明确定义,事件抽取要求将文本数据中表征事件信息的非结构化数据转化为结构化、可存储使用的精准知识。

[0003] 当今社会,在网络新闻媒体上实时推送着各类大大小小的热点事件。面对日益增长的海量互联网信息,快速定位到公众讨论的具体事件变得至关重要。这不仅可以帮助舆情监管人员快速定位到具体事件,了解事件的具体要素,还可以将事件抽取结果提供给其他自然语言处理任务,以进行更深入的分析应用。因网络和社会因素影响,事件抽取技术研究在国内外研究热度逐年升高。

发明内容

[0004] 针对上述问题,本发明的目的在于提供一种基于字符嵌入的事件元素网格化抽取方法、存储介质及电子装置,事件元素网格化抽取是在事件检测任务的基础上,将模型细化为事件触发词抽取、事件论元抽取和事件属性抽取3个网格模块,各网格模块既联合共享基础模型事件语义信息,又独立优化各自元素的抽取性能。技术方案如下:

[0005] 一种基于字符嵌入的事件元素网格化抽取方法,包括以下步骤:

[0006] 步骤1:构建事件元素抽取基础模型;

[0007] 所述基础模型为字符嵌入神经网络深度学习模型,包括基于BERT的信息预学习表示层、字编码嵌入层、BiGRU双向门控循环神经网络层、自注意力层和CRF条件随机场输出层;其运行步骤如下:

[0008] 步骤1.1:基于BERT的信息预学习表示层对样本数据上下文语义特征进行预学习,生成突发元事件域的文本表示模型;

[0009] 步骤1.2:字编码嵌入层利用训练好的BERT语言模型中生成的语义表示向量输入到BiGRU双向门控循环神经网络层;

[0010] 步骤1.3:BiGRU双向门控循环神经网络层提取输入序列的上下文依赖的长距离深层特征;

[0011] 步骤1.4:自注意力层对BiGRU双向门控循环神经网络层学习到的深层特征进行加权变换,以突出文本序列中重点词汇信息;

[0012] 步骤1.5:CRF条件随机场输出层将触发词提取转化为序列标注任务,以解决传统词向量转化为字向量后上下文标注信息问题;

[0013] 将所述基础模型细化为事件触发词抽取、事件论元抽取和事件属性抽取3个网格模块;

[0014] 步骤2:对事件触发词抽取模型进行抽取优化:从一个事件句抽取一个主事件触发词作为事件触发词,多余事件触发词作外部特征,辅助表征主事件;将标注数据中的所有事件触发词作为知识库,作为事件触发词抽取的先验特征;抽取句子中与知识库内事件触发词匹配的触发词,并标注出来,与句子BERT语义编码的输出的字符嵌入向量拼接;并将事件类型向量拼接字符嵌入向量中;事件触发词抽取任务的目标向量由事件触发词的抽取结果表示,对应事件触发词的标注结果;

[0015] 步骤3:对事件论元抽取模型进行抽取优化:在原文本BERT语义特征之上,将文本中所有字符到事件触发词的相对距离作为文本结构特征,事件触发词本身的相对距离为0;并将事件主体与客体联合,事件时间与地点联合,采用两个独立的模型进行提取;事件论元抽取任务的目标向量对应事件论元的抽取结果标注;

[0016] 步骤4:对事件属性抽取模型进行抽取优化:定义事件属性为事件时态和事件极性,模型输出转变为多分类问题,将所述基础模型中CRF条件随机场输出层置换构造两个分类器;将事件触发词及左右两端进行动态池化获得的特征作为全局特征,与句子BERT语义编码的输出的字符嵌入向量拼接,并采用十折交叉验证方法进行优化;

[0017] 步骤5:采用步骤1-步骤4训练得到的事件要素抽取模型对测试数据进行事件元素抽取结果的预测。

[0018] 进一步的,所述步骤1.3中,BiGRU双向门控循环神经网络层同时包含一个正向的 $\overrightarrow{GRU}(x_0, x_1, \dots, x_t, \dots, x_{n-1})$ 和一个反向的 $\overleftarrow{GRU}(x_{n-1}, x_{n-2}, \dots, x_t, \dots, x_0)$, 正向GRU捕捉0:t的上文特征信息 a_t , 反向GRU捕捉t:n-1的下文特征信息 a'_t , 通过拼接所捕获的上下文特征信息,获得句子的上下文信息 y_t , 如公式(1)-(3)所示:

$$[0019] \quad \overrightarrow{GRU}(x_0, x_1, \dots, x_t, \dots, x_{n-1}) = (a_0, a_1, \dots, a_t, \dots, a_{n-1}) \quad (1)$$

$$[0020] \quad \overleftarrow{GRU}(x_{n-1}, x_{n-2}, \dots, x_t, \dots, x_0) = (a'_{n-1}, a'_{n-2}, \dots, a'_t, \dots, a'_0) \quad (2)$$

$$[0021] \quad y_t = [a_t, a'_t] \quad (3)$$

[0022] 式中, x_t 表示词序列特征向量; a_t 表示正向GRU捕捉0:t的上文特征信息; a'_t 表示反向GRU捕捉t:n-1的下文特征信息; y_t 表示获得句子的上下文信息;

[0023] 所述加权变换公式(4)所示:

$$[0024] \quad e_{ij} = a(W \overrightarrow{y}_i, W \overleftarrow{y}_j) \quad (4)$$

[0025] 式中, e_{ij} 表示句子j的特征对句子i的重要性; a 表示注意力机制; \cdot 表示共享参数的线性变换权重矩阵; y_i 和 y_j 分别表示获得的句子i和句子j的上下文信息。

[0026] 更进一步的,在事件触发词抽取任务中,假设词 w_i 的事件触发词类型目标向量为 $[tri_0, tri_1, tri_2, \dots, tri_j, \dots, tri_n]$, 则 tri_j 的设置方式如公式(5)所示:

$$[0027] \quad tri_j = \begin{cases} 0, & \text{其他} \\ 1, & \text{如果 } w_i \text{ 的事件触发词类别标签ID为 } j \end{cases} \quad (5)$$

[0028] 所述步骤3中,在事件论元抽取任务中,假设词 w_i 的事件论元类型目标向量为 $[\arg_0, \arg_1, \dots, \arg_j, \dots, \arg_n]$,其中 \arg_j 的设置方式如公式(6)所示:

$$[0029] \quad \arg_j = \begin{cases} 0, \text{其他} \\ 1, \text{如果} w_i \text{的事件论元类别标签ID为} j \end{cases} \quad (6).$$

[0030] 一种存储介质,所述存储介质中存储有计算机程序,其中,所述计算机程序被设置为运行时执行上述的方法。

[0031] 一种电子装置,包括存储器和处理器,所述存储器中存储有计算机程序,所述处理器被设置为通过所述计算机程序执行所上述的方法。

[0032] 本发明的有益效果是:本发明利用事件元素抽取基础模型,分别对模型细化的事件触发词抽取、事件论元抽取和事件属性抽取3个网格模块进行抽取优化,各网格模块既联合共享基础模型事件语义信息,又独立优化各自元素的抽取性能,结果表明基于字符嵌入的事件元素网格化抽取模型在事件元素抽取任务中表现良好,在事件元素抽取任务中表现良好,获得了较高的准确率;此外,该模型后续可以开展更多的研究。

附图说明

[0033] 图1是本发明方法流程示意图。

[0034] 图2是本发明中建立的事件要素抽取基础模型示意图。

[0035] 图3是本发明分模块事件论元抽取对比实验结果示意图。

[0036] 图4是本发明分模块事件属性抽取优化方法对比分析实验结果示意图。

[0037] 图5是本发明事件元素抽取不同方法对比实验结果。

具体实施方式

[0038] 下面结合附图和具体实施方式对本发明做进一步详细的说明。一种基于字符嵌入的事件元素网格化抽取方法,包括以下步骤:

[0039] 步骤1:构建事件元素抽取基础模型;

[0040] 如图2所示,事件元素抽取基础模型主要包含基于BERT的信息预学习表示层、字编码嵌入层、BiGRU双向门控循环神经网络层、Self-attention自注意力层和CRF条件随机场输出层。

[0041] 使用BERT模型能够对样本数据上下文语义特征预学习,生成突发元事件域的文本表示模型。再利用训练好的BERT语言模型中生成的语义表示向量输入到BiGRU,利用BiGRU提取输入序列的上下文依赖的长距离深层特征。

[0042] 其中,BERT语言模型是一个著名的语言模型,是2018年10月由Google AI研究院提出的一种预训练模型。本发明仅使用Bert模型进行语义表示。

[0043] BiGRU双向门控循环神经网络层同时包含一个正向的 $\overrightarrow{GRU}(x_0, x_1, \dots, x_t, \dots, x_{n-1})$ 和一个反向的 $\overleftarrow{GRU}(x_{n-1}, x_{n-2}, \dots, x_t, \dots, x_0)$,正向GRU捕捉0:t的上文特征信息 a_t ,反向GRU捕捉t:n-1的下文特征信息 a'_t ,通过拼接所捕获的上下文特征信息,获得句子的上下文信息 y_t ,如公式1-3所示。

$$[0044] \quad \overrightarrow{GRU}(x_0, x_1, \dots, x_t, \dots, x_{n-1}) = (a_0, a_1, \dots, a_t, \dots, a_{n-1}) \quad (1)$$

$$[0045] \quad \overrightarrow{GRU}(x_{n-1}, x_{n-2}, \dots, x_t, \dots, x_0) = (a'_{n-1}, a'_{n-2}, \dots, a'_t, \dots, a'_0) \quad (2)$$

$$[0046] \quad y_t = [a_t, a'_t] \quad (3)$$

[0047] Self-attention自注意力层用于对BiGRU学习到的深层特征进行加权变换,突出文本序列中重点词汇信息,如公式4所示。最终使用CRF将触发词提取转化为序列标注任务,解决传统词向量转化为字向量后上下文标注信息。

$$[0048] \quad e_{ij} = a(W\vec{y}_i, W\vec{y}_j) \quad (4)$$

[0049] 步骤2:对事件触发词抽取模型进行抽取优化;

[0050] 一个事件描述句中可能存在多个事件触发词。在事件元素的抽取过程中,不仅要完成事件元素的抽取,还必须使事件元素和事件触发词对应。同时,一个事件描述句中信息元素有限,多个事件中存在主次关系。为了抽取出主要关注事件和更丰富的事件元素,一个事件句抽取一个主事件触发词作为事件触发词,多余事件触发词作外部特征,辅助表征主事件。将标注数据中的所有事件触发词作为知识库,类似于远程监督的方式,作为事件触发词抽取的先验特征。抽取句子中与知识库内事件触发词匹配的触发词标注出来,与句子BERT语义编码的输出的字符嵌入向量拼接。

[0051] 另外,事件元素的组成和事件的类型有很大的关系,如“突袭”等涉恐涉爆事件,由触发词的含义可知是两方发生冲突,一般在触发词的邻近位置会有冲突的双方;“地震”等重大灾情事件,由触发词的含义可知是某处有灾情表述,那么触发词邻近位置出现地点要素的可能性会很大。因此,事件元素抽取中事件类型具有重要语义线索,事件类型向量拼接到字符嵌入向量中。

[0052] 在事件触发词抽取任务中,目标向量是事件触发词的抽取结果表示,目标向量对应了事件触发词的标注结果。如表1所示,三种事件触发词标签长度,BIO标注模式分别是“B-Trigger”,“I-Trigger”和“Other”。

[0053] 表1事件触发词标注标记于含义

ID	标记	含义
0	Other	其他
1	B-Trigger	事件触发词的开始
2	I-Trigger	事件触发词的中间

[0055] 假设词 w_i 的事件触发词类型目标向量为 $[tri_0, tri_1, tri_2]$,其中 tri_j 的设置方式如公式(1)所示:

$$[0056] \quad tri_j = \begin{cases} 0, & \text{其他} \\ 1, & \text{如果 } w_i \text{ 的事件触发词类别标签ID为 } j \end{cases} \quad (1)$$

[0057] 步骤3:对事件论元抽取模型进行抽取优化;

[0058] 事件论元中的事件主体、事件客体、事件时间和事件地点四个元素在语义结构上

受到事件触发词的重要影响。为了获得事件论元元素在语句语义结构上的潜在特征,在原文本BERT语义特征之上,将文本中所有字符到事件触发词的相对距离作为文本结构特征,事件触发词本身的相对距离为0。并将事件主体与客体联合,事件时间与地点联合,采用两个独立的模型进行提取。

[0059] 在事件论元抽取任务中,目标向量对应了事件论元的抽取结果标注。各个事件元素标签类型及其含义如表2所示,九种事件触发词标签长度,BIO标注模式分别是“B-Subject”,“I-Subject”,“B-Object”,“I-Object”,“B-Time”,“I-Time”,“B-Location”,“I-Location”和“Other”。

[0060] 表2事件论元标注及含义

ID	标记	含义
0	Other	其他
1	B-Subject	主体论元的开始位置
2	I-Subject	主体论元的中间位置
3	B-Object	客体论元的开始位置
4	I-Object	客体论元的中间位置
5	B-Time	时间论元的开始位置
6	I-Time	时间论元的中间位置
7	B-Location	地点论元的开始位置
8	I-Location	地点论元的中间位置

[0062] 假设词 w_i 的事件论元类型目标向量为 $[\arg_0, \arg_1, \dots, \arg_j, \dots, \arg_8]$,其中 \arg_j 的设置方式如公式(2)所示:

$$[0063] \quad \arg_j = \begin{cases} 0, & \text{其他} \\ 1, & \text{如果 } w_i \text{ 的事件论元类别标签ID为 } j \end{cases} \quad (2)$$

[0064] 在事件论元抽取中,事件主体、事件客体、事件时间和事件地点元素分布差距较大,一个模型会导致此事件时间和事件地点两元素抽取效果较差。为了提升事件论元抽取中各事件元素的抽取效果,表3所示为是否采用分模块进行事件论元抽取的对比实验结果。

[0065] 表3分模块事件论元抽取对比实验结果

抽取方法	事件论元抽取结果		
	P	R	F
一个模块	0.6105	0.6667	0.6374
Sub&Obj+Tim&Loc	0.7579	0.7941	0.7756

[0067] Sub&Obj和Tim&Loc表示将四个事件元素拆分为两个事件论元对,独立的训练两个模型,进行事件论元的抽取。由图1可知,分模块进行事件论元的抽取,能够解决数据中论元分布不均的问题,有效提升论元抽取效果。

[0068] 步骤4:对事件属性抽取模型进行抽取优化;

[0069] 定义事件属性分为事件时态和事件极性,事件时态分为“过去”、“现在”、“将来”和“其他”,事件极性分为“肯定”、“否定”和“可能”。模型输出转变为多分类问题,所将基础模型CRF输出层置换构造两个分类器。分类器激活使用softmax多分类函数,损失函数为CrossEntropyLoss。

[0070] 表征事件时态和事件极性的词语大多存在事件触发词附近。相比与利用文本全局特征,设置事件触发词附近的池化窗口,提取相关的紧密局部特征,更有利于事件属性的抽取。采用将事件触发词及左右两端进行动态池化获得的特征作为全局特征,与句子BERT语义编码的输出的字符嵌入向量拼接。此外,为了提升模型泛化性能,考虑采用十折交叉验证方法进行优化。

[0071] 事件属性优化抽取中,加入了事件触发词左右动态池化特征和十折交叉验证两种优化方法。为了验证上述事件属性元素抽取模型中所在采用的优化方法的有效性,对比分析实验结果如表4所示。

[0072] 表4事件属性抽取优化方法对比分析实验结果

优化方法	事件属抽取结果		
	P	R	F
无优化	0.6304	0.6095	0.6198
Trigger 池化特征	0.7096	0.6701	0.6893
十折交叉验证	0.6781	0.6912	0.6846
Trigger 池化特征+十折交叉验证	0.7351	0.7025	0.7184

[0074] 由图4可知,相比于基础模型无优化的情况,添加触发词池化特征或者进行十折交叉验证都能提高事件属性的抽取效果;同时添加触发词池化特征和进行十折交叉验证能大幅度提升事件属性的抽取性能。经过分析,事件触发词左右动态池化特征利用触发词与事件属性潜在关系,有利于提升事件属性元素的抽取性能;十折交叉验证可以在一定程度上减小过拟合,在有限的数据中获取尽可能多的有效信息,缓解数据中元素分布不均的问题,提升模型的泛化能力。

[0075] 步骤5:采用步骤1-步骤4训练得到的事件要素抽取模型对测试数据进行事件元素抽取结果的预测。BiGRU-SATT-CRF为本发明提出的基于字符嵌入的事件元素抽取方法,实验结果如表5所示。

[0076] 表5事件元素抽取不同方法对比实验结果

实验方法	事件元素抽取结果		
	P	R	F
JRNN	0.6642	0.7061	0.6845
Croos-Event	0.6858	0.6894	0.6877
[0077] JointBeam	0.7140	0.6233	0.6656
dbRNN	0.7303	0.6981	0.7138
DMCNN	0.7554	0.6365	0.6909
JMEE	0.7638	0.7137	0.7379
BiGRU-SATT-CRF	0.7861	0.7282	0.7560

[0078] 从图5的实验结果可以看出,基于字符嵌入的神经网络事件元素抽取方法的实验结果均优于其他抽取方法,这说明了基于字符嵌入和分模块优化的神经网络方法在事件元素抽取任务中具有一定的优势。

[0079] 可将本发明方法编为程序代码,通过计算机刻度存储介质存储该代码,将程序代码传输给处理器,通过处理器执行本发明方法。

[0080] 本发明利用事件元素抽取基础模型,分别对模型细化的事件触发词抽取、事件论元抽取和事件属性抽取3个网格模块进行抽取优化,构造不同的特征向量和目标向量,结果表明基于字符嵌入的事件元素网格化抽取模型在事件元素抽取任务中表现良好。此外,该模型后续可以开展更多的研究。

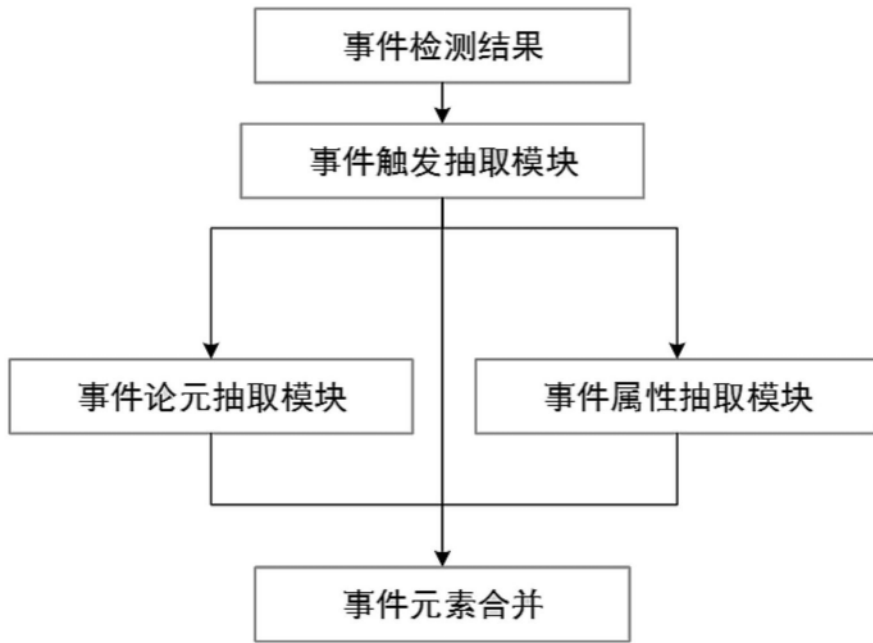


图1

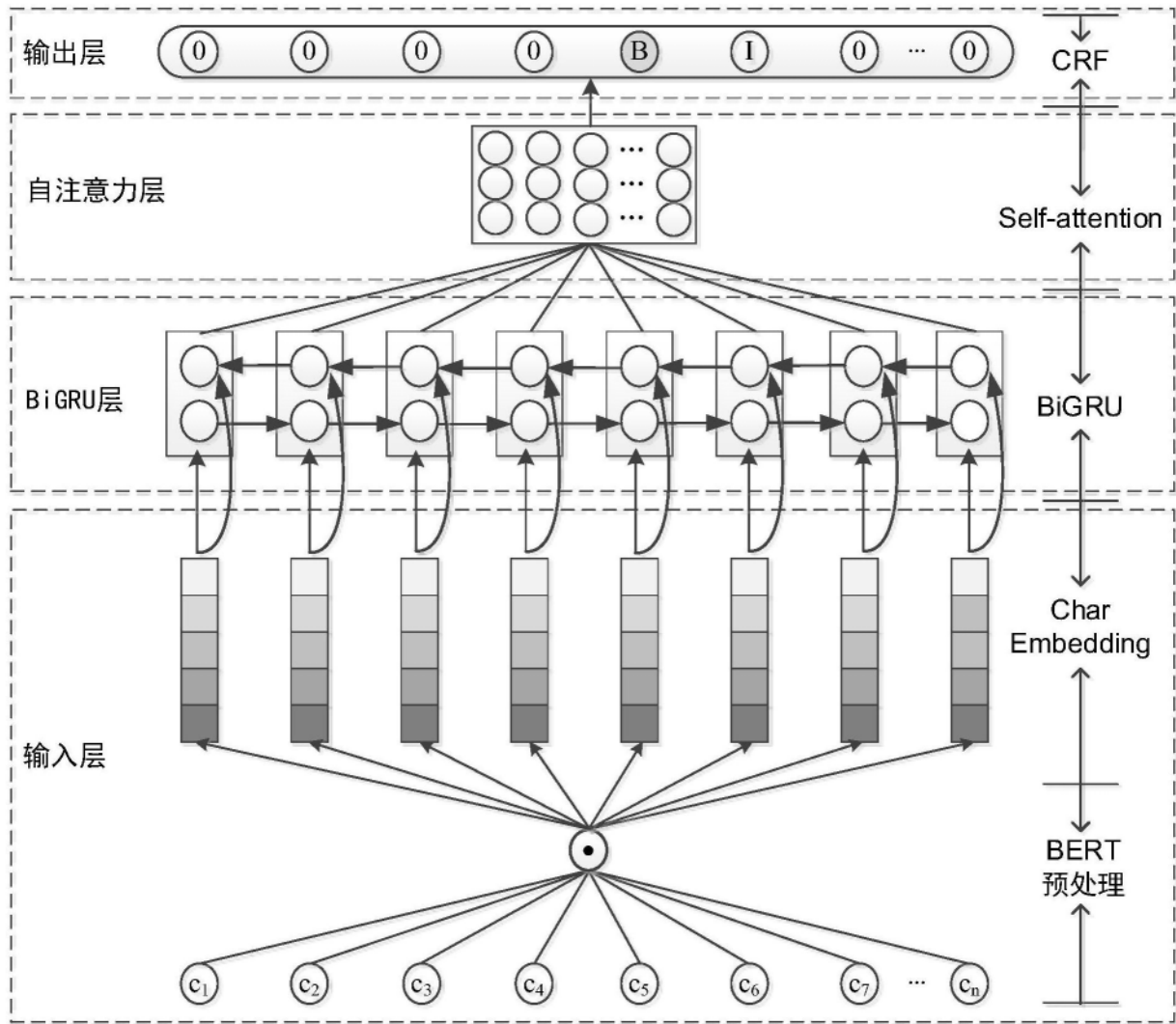


图2

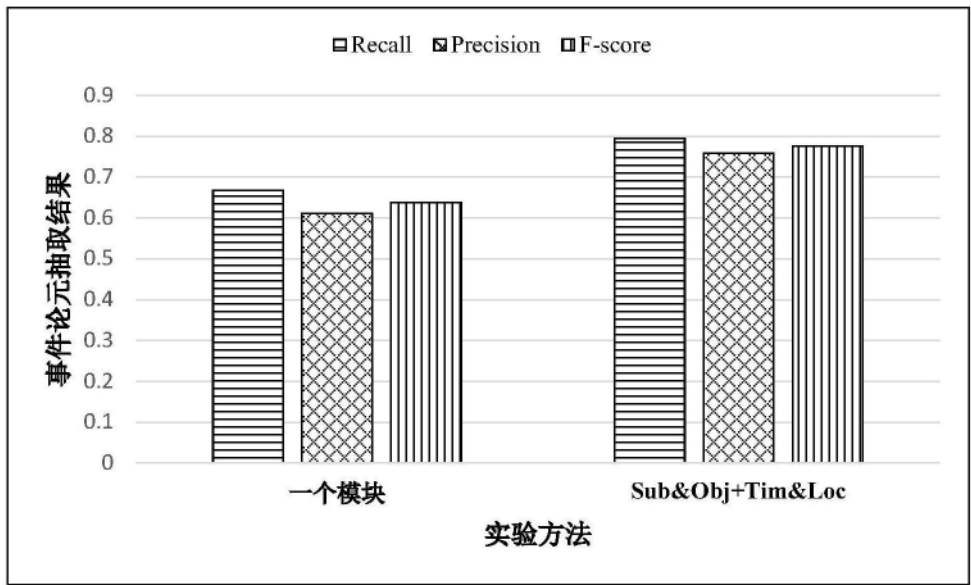


图3

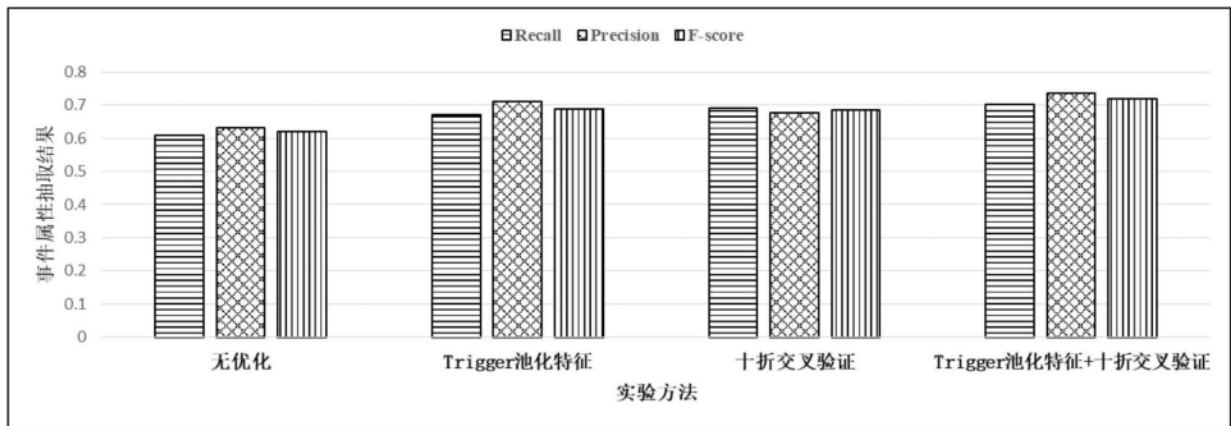


图4

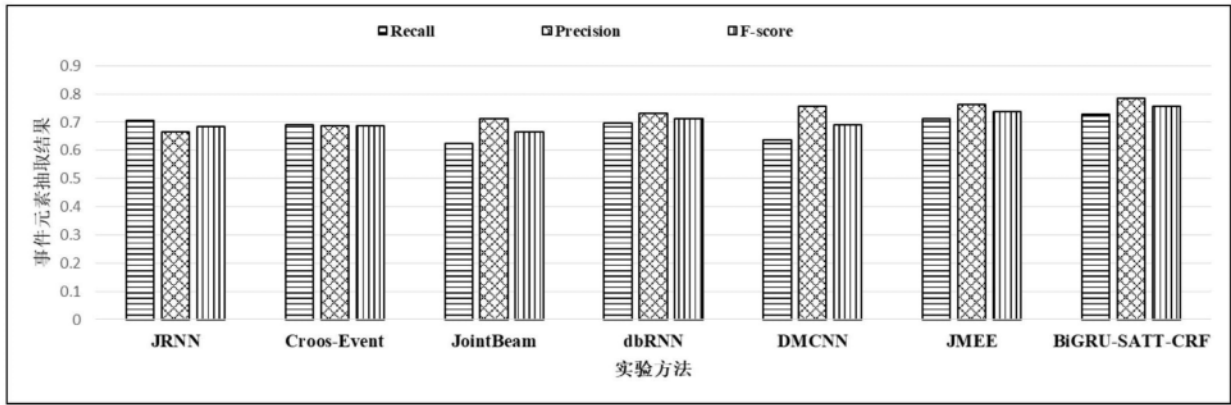


图5