



(12) 发明专利申请

(10) 申请公布号 CN 112836046 A

(43) 申请公布日 2021.05.25

(21) 申请号 202110039836.2

(22) 申请日 2021.01.13

(71) 申请人 哈尔滨工程大学

地址 150001 黑龙江省哈尔滨市南岗区南
通大街145号哈尔滨工程大学科技处
知识产权办公室

(72) 发明人 范贺添 申林山 黄少滨 李熔盛
吴汉瑜 谷虹润

(51) Int. Cl.

G06F 16/35 (2019.01)

G06F 40/295 (2020.01)

G06K 9/62 (2006.01)

G06N 3/04 (2006.01)

G06N 3/08 (2006.01)

G06F 16/36 (2019.01)

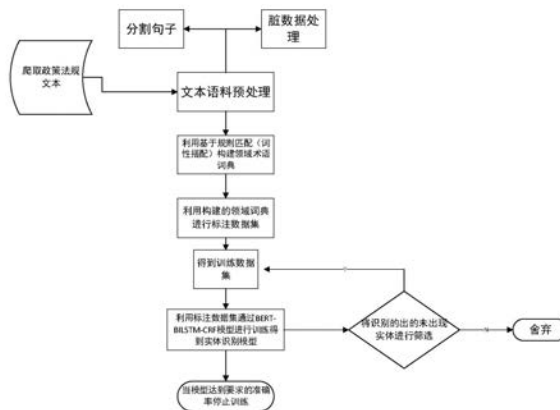
权利要求书1页 说明书7页 附图2页

(54) 发明名称

一种四险一金领域政策法规文本实体识别方法

(57) 摘要

本发明属于命名实体识别技术领域,具体涉及一种四险一金领域政策法规文本实体识别方法。本发明通过预训练语言模型BERT得到每个字符基于上下文特征动态生成的具有上下文语义信息的字向量,通过word2vec中的skip-gram模型得到每个词语的词向量,将具有上下文语义信息的字向量和其所在的词的词向量利用维度拼接的方式进行特征融合,得到联合字词向量,不仅可以弥补少量标注样本特征不足和字符语义提取不充分的问题,还为字向量补充了词级短语信息,从而在一定程度上提高模型的准确率。本发明可以有效解决四险一金领域命名实体识别任务中标注数据不足以及识别精度不高的问题。



1. 一种四险一金领域政策法规文本实体识别方法,其特征在于,包括以下步骤:

步骤1:输入待识别的文本;构建四险一金领域实体分词和标注词典,基于四险一金领域实体分词和标注词典预训练语言模型BERT;

步骤2:对待识别的文本进行分词处理;

步骤3:取部分分词后的待识别的文本构建训练集,其余分词后的待识别的文本组成测试集;根据四险一金领域实体分词和标注词典对训练集中分词后的待识别的文本进行标注;

步骤4:将训练集中标注好的词语切分为单个的汉字,并根据该词语的实体类别以及每个汉字在实体中出现的位置进行进一步的BIO实体边界标记;

步骤5:将标注好的训练集输入到预训练语言模型BERT中,得到每个字符基于上下文特征动态生成的具有上下文语义信息的字向量 $W_i^{charbert}$;

步骤6:将标注好的训练集输入word2vec中的skip-gram模型中训练,得到每个词语的词向量;

步骤7:将具有上下文语义信息的字向量 $W_i^{charbert}$ 和其所在的词的词向量 W_i^{word} 利用维度拼接的方式进行特征融合,得到融合后的字词组合向量 W_i ;

$$W_i = \begin{bmatrix} W_i^{charbert} \\ W_i^{word} \end{bmatrix}$$

步骤8:将训练集中融合后的字词组合向量 W_i 输入至双向长短期记忆网络Bi-LSTM和条件随机场CRF模型进行训练,得到实体识别和分类模型;

首先将训练集中融合后的字词组合向量 W_i 作为输入向量输入到双向长短期记忆网络Bi-LSTM中采集到有效的上下文信息,再利用条件随机场CRF模型作为解码器对模型进行解码,即通过转移概率来得到每个字符最优的标记序列,从而为实体打上类别标签,实现实体识别和分类;

步骤9:将测试集输入到训练好的实体识别和分类模型中,得到待识别文本的实体识别结果;所述的实体识别结果包含实体、实体的起始位置、实体的终止位置、实体的类别标签。

一种四险一金领域政策法规文本实体识别方法

技术领域

[0001] 本发明属于命名实体识别技术领域,具体涉及一种四险一金领域政策法规文本实体识别方法。

背景技术

[0002] 随着社会发展,制度逐步健全完善,我国坚持基本保障制度的作用越来越凸显。因此对四险一金领域的问答系统和知识图谱构建等研究工作具有重要的意义。命名实体识别(Named Entity Recognition,NER)作为知识图谱的重要基本单元是知识图谱构建和补全的核心技术。是指识别文本中具有特定意义的实体,主要包括人名、地名、机构名、专有名词等。因此,在构建四险一金领域知识图谱时,识别出四险一金领域相关的专业术语和常用的命名实体(如机构名,地名等)也是具有重要意义的。

[0003] 传统的命名实体识别方法主要有基于规则匹配的算法和基于机器学习的算法两大类。然而传统的机器学习模型(如CRF条件随机场)虽仍是NER的主流模型的重要组成部分。它的目标函数不仅考虑输入的状态特征函数,而且还包含了标签转移特征函数。从而得到最优标记序列。)但是都存在着一个共同的缺点对于特征提取的要求比较高,需要选择对命名实体识别任务有影响的各种特征,并将这些特征组合成向量来表示文本中的词语并且需要事先对预处理的数据进行大量的人工标注才能训练出较好的效果,因此建模的成本很高。近年来,伴随着计算机算力的发展以及词的分布式表示(word embedding)的提出,深度学习方法逐渐被运用到命名实体识别任务中,神经网络成为可以高效处理许多NLP任务的模型主要表现在基于神经网络的深度学习方法具有很强的泛化性,为了使词语能包含更加全面的语义信息以及句法特征,今年学者们又提出了利用预训练语言模型进一步增强字向量的表示,其中最为突出的是由谷歌研究员Devlin等提出的BERT模型(Bidirectional Encoder Representations from Transformers),利用自注意力机制和Transformer编码器对大规模公开语料进行预训练,得到更具有上下文语义信息的字向量,pengM等利用该方法在通用领域的实体识别效果已经取得不错的效果。

发明内容

[0004] 本发明的目的在于提供一种四险一金领域政策法规文本实体识别方法。

[0005] 本发明的目的通过如下技术方案来实现:包括以下步骤:

[0006] 步骤1:输入待识别的文本;构建四险一金领域实体分词和标注词典,基于四险一金领域实体分词和标注词典预训练语言模型BERT;

[0007] 步骤2:对待识别的文本进行分词处理;

[0008] 步骤3:取部分分词后的待识别的文本构建训练集,其余分词后的待识别的文本组成测试集;根据四险一金领域实体分词和标注词典对训练集中分词后的待识别的文本进行标注;

[0009] 步骤4:将训练集中标注好的词语切分为单个的汉字,并根据该词语的实体类别以

及每个汉字在实体中出现的位置进行进一步的BIO实体边界标记；

[0010] 步骤5:将标注好的训练集输入到预训练语言模型BERT中,得到每个字符基于上下文特征动态生成的具有上下文语义信息的字向量 $W_i^{charbert}$;

[0011] 步骤6:将标注好的训练集输入word2vec中的skip-gram模型中训练,得到每个词语的词向量;

[0012] 步骤7:将具有上下文语义信息的字向量 $W_i^{charbert}$ 和其所在的词的词向量 W_i^{word} 利用维度拼接的方式进行特征融合,得到融合后的字词组合向量 W_i ;

$$[0013] \quad W_i = \begin{bmatrix} W_i^{charbert} \\ W_i^{word} \end{bmatrix}$$

[0014] 步骤8:将训练集中融合后的字词组合向量 W_i 输入至双向长短期记忆网络Bi-LSTM和条件随机场CRF模型进行训练,得到实体识别和分类模型;

[0015] 首先将训练集中融合后的字词组合向量 W_i 作为输入向量输入到双向长短期记忆网络Bi-LSTM中采集到有效的上下文信息,再利用条件随机场CRF模型作为解码器对模型进行解码,即通过转移概率来得到每个字符最优的标记序列,从而为实体打上类别标签,实现实体识别和分类;

[0016] 步骤9:将测试集输入到训练好的实体识别和分类模型中,得到待识别文本的实体识别结果;所述的实体识别结果包含实体、实体的起始位置、实体的终止位置、实体的类别标签。

[0017] 本发明的有益效果在于:

[0018] 本发明通过预训练语言模型BERT得到每个字符基于上下文特征动态生成的具有上下文语义信息的字向量,通过word2vec中的skip-gram模型得到每个词语的词向量,将具有上下文语义信息的字向量和其所在的词的词向量利用维度拼接的方式进行特征融合,得到联合字词向量,不仅可以弥补少量标注样本特征不足和字符语义提取不充分的问题,还为字向量补充了词级短语信息,从而在一定程度上提高模型的准确率。本发明可以有效解决四险一金领域命名实体识别任务中标注数据不足以及识别精度不高的问题。

附图说明

[0019] 图1为本发明的预训练语言模型BERT的模型图。

[0020] 图2为本发明的整体实施流程图。

[0021] 图3为本发明的实施例中实体标签描述表。

具体实施方式

[0022] 下面结合附图对本发明做进一步描述。

[0023] 本发明涉及一种四险一金领域政策法规文本实体识别方法,用于从四险一金领域政策法规文本中自动识别出与具有领域特性的命名实体,具体的说,从中央到地方政府发布的政策法规文本中识别出和四险一金领域相关的命名实体。

[0024] 现有的四险一金领域的命名实体识别存在以下问题:一是与通用领域不同,四险一金政策法规文本的实体具有特殊性,不但包含有大量专有的领域术语,在普通词库不一定包含这些领域术语;会出现大量名词组合的情况。二是四险一金领域也缺少公开的大规

模标注的数据集。

[0025] 针对以上问题本发明提出了利用基于规则词性搭配方法对四险一金领域词典的构建。利用领域词典对选取原始文本进行标注。不仅减少了大量的人工成本,也方便了后续快速扩充训练数据和对原始文本进行分词,标注等预处理工作。将BERT预训练作为字向量的特征层和通过Word2Vec模型对四险一金政策法规分词后文本中的词语特征进行提取,并训练成的词向量拼接得到的联合字词向量。不仅可以弥补少量标注样本特征不足和字符语义提取不充分的问题,还为字向量补充了词级短语信息。最后利用双向长短期记忆网络(Bi-LSTM)和条件随机场(CRF)对联合字词向量进行训练得到四险一金领域实体识别模型。

[0026] 面向四险一金领域,针对该领域中实体长度过长和存在词语嵌套造成的识别精度不高的问题,本发明提出一种基于预训练语言模型BERT的实体识别方法,该模型利用BERT模型增强政策法规中字符的语义表示并根据其所在上下文特征动态生成字向量,同时考虑到汉字不是中文语义的最基本单位,使用生成的动态字符向量与所在词的词向量拼接后得到组合向量作为Bi-LSTM-CRF模型输入,其中Bi-LSTM层进行编码和CRF层解码,最后标注出实体识别结果。

[0027] 一种四险一金领域政策法规文本实体识别方法,包括以下步骤:

[0028] 步骤1:输入待识别的文本;构建四险一金领域实体分词和标注词典,基于四险一金领域实体分词和标注词典预训练语言模型BERT;

[0029] 步骤2:对待识别的文本进行分词处理;

[0030] 步骤3:取部分分词后的待识别的文本构建训练集,其余分词后的待识别的文本组成测试集;根据四险一金领域实体分词和标注词典对训练集中分词后的待识别的文本进行标注;

[0031] 步骤4:将训练集中标注好的词语切分为单个的汉字,并根据该词语的实体类别以及每个汉字在实体中出现的位置进行进一步的BIO实体边界标记;

[0032] 步骤5:将标注好的训练集输入到预训练语言模型BERT中,得到每个字符基于上下文特征动态生成的具有上下文语义信息的字向量 $W_i^{charbert}$;

[0033] 步骤6:将标注好的训练集输入word2vec中的skip-gram模型中训练,得到和每个词语的词向量;

[0034] 步骤7:将具有上下文语义信息的字向量 $W_i^{charbert}$ 和其所在的词的词向量 W_i^{word} 利用维度拼接的方式进行特征融合,得到融合后的字词组合向量 W_i ;

$$[0035] \quad W_i = \begin{bmatrix} W_i^{charbert} \\ W_i^{word} \end{bmatrix}$$

[0036] 步骤8:将训练集中融合后的字词组合向量 W_i 输入至双向长短期记忆网络Bi-LSTM和条件随机场CRF模型进行训练,得到实体识别和分类模型;

[0037] 首先将训练集中融合后的字词组合向量 W_i 作为输入向量输入到双向长短期记忆网络Bi-LSTM中采集到有效的上下文信息,再利用条件随机场CRF模型作为解码器对模型进行解码,即通过转移概率来得到每个字符最优的标记序列,从而为实体打上类别标签,实现实体识别和分类;

[0038] 步骤9:将测试集输入到训练好的实体识别和分类模型中,得到待识别文本的实体

识别结果;所述的实体识别结果包含实体、实体的起始位置、实体的终止位置、实体的类别标签。

[0039] 实施例1:

[0040] 因为四险一金政策法规文本是通过网络爬虫获取,可能含有html标签以及一些乱码和表格符号,应该对原始文本采用utf-8编码格式进行统一编码,通过制定正则表达式去除空格等乱码字段。将预处理好的文本进行分词和词性标注。

[0041] 领域术语构成词的方式可以分为单词概念与词组型领域概念。单词型领域概念是由一个单词组成,所以其不能再被分割,是最小的独立词单元。而词组型领域概念是由两个或两个以上的单词构成,并不要求其中的单词一定是单词型领域概念,可以是其他词语。再对语料分词后统计发现,四险一金领域术语多集中在二元、三元和四元词组,通过对N-gram进行统计,选出一起出现频率较高的词组,通过分析和统计领域词语的特点,根据词性制定规则表和人工筛选,去除不符合规则的词语。利用构建好的词典,并借助Jieba分词器+用户字典的方式依据最大匹配原则对原始政策法规文本进行分词处理,并对分好词的文本进行实体类别自动标注等预处理工作;

[0042] 通过爬取以下几种知识,包括四险一金的司法案例、中央法律法规和地方法规规章相关和与四险一金领域百科词条。法规主要来源于北大法宝,百科词条主要来源于百度百科。通过基于规则词性搭配和部分人工帮助,利用政策法规文本为语料得到领域术语概念集合。(因为目前的中文分词工具虽然达到了较高的准确率,但是由于分词粒度细,对一些领域概念处理效果不好,如“基本养老保险费”,在经过分词后为“基本/养老保险费”,而“基本养老保险费”应该被看作一个术语实体,却被分成2个词,导致失去部分语义信息。)除领域专业术语外,本发明通过对政策法规中出现常用领域实体进行人工定义及归类。最后将四险一金领域实体总结5个类别(包括领域术语、地名、机构名、人名、法规名)进行类别标注从而构建四险一金领域实体分词和标注词典。

[0043] 利用构建好的词典并借助Jieba分词工具对原始的政策法规本进行分词和添加类别标记,本发明所用的语料包括四险一金领域(养老保险,工伤保险,医疗保险,失业保险,公积金)有关部门发布的司法案例,中央法律法规,以及地方法规规章共计25554篇文章作为实验语料,其中养老保险7704篇,失业保险1357篇,工伤保险1946篇,生育/医疗保险7749+996=8745篇,住房公积金2969篇。对语料中每个险种按原有比例共抽取1000篇。将标注好类别的词语切分为单个的汉字,并根据其实体类别以及在实体词中出现的位置进行进一步的BIOES实体边界标记,如“基本养老保险费”被标注为{基B-PRO}{本I-PRO}{养I-PRO}{老I-PRO}{保I-PRO}{险I-PRO}费I-PRO}。(PRO为实体标签)从标注的1000篇政策法规中70%用作训练集20%用作验证集,10%用作测试集。

[0044] 字级别特征利用预训练的BERT语言模型对输入的文本信息的字向量初始化,所获得的字向量记为序列 $X = (x_1, x_2, x_3, \dots, x_n)$ 可以利用上下文语义信息,解决传统字符向量不能根据语境表示为不同的特征向量的问题,从而可以更加有效提取文本中的语义特征。

[0045] 词级别的特征提取和表示,通过Word2Vec模型对四险一金政策法规分词后文本中的词语特征进行提取,并训练成词向量表示。

[0046] 通过使用维度拼接的方式对词向量和通过BERT模型得到字向量进行融合。

[0047] 利用双向长短期记忆网络 (Bi-LSTM) 和条件随机场 (CRF) 模型对得到的用于实体识别和分类的字词的联合特征向量进行训练, 最终得到可以实现对四险一金领域政策法规文本进行实体识别的模型, 并对得到的模型的F1值分别进行评估测试, 并应用于四险一金领域知识图谱构建。

[0048] bert字向量: 在训练语料上使用word2vec中的skip-gram模型训练得到字符向量将字符向量 W_i^{char} , 输入到预训练语言模型bert中得到具有上下文语义信息字向量 $W_i^{charbert}$

[0049] 词向量: 对于词向量 W_i^{word} 的获得, 本文首先使用jieba分词对中文文本进行分词, 然后使用skip-gram模型在分词后的语料训练得到。

[0050] 将通过BERT预训练语言模型得到的具有上下文语义信息的字向量和其所在的词得到的词向量利用维度拼接的方式进行特征融合, 最终得到的维度为二者之和的字词联合表示, 即为融合后的字词组合向量。

$$[0051] \quad W_i = \begin{bmatrix} W_i^{charbert} \\ W_i^{word} \end{bmatrix}$$

[0052] LSTM又称为长短期记忆网络, 作为循环网络RNN的一种变体也是一种序列模型, 它输入门, 遗忘门, 输出门, 选择性的传递时序信息, 从而有效的克服了普通RNN模型由于序列过长导致的梯度消失问题。LSTM结构结构可以形式化表示为:

$$[0053] \quad i_t = \sigma(x_t \cdot w_{xh}^i + h_{t-1} \cdot w_{hh}^i + b_h^i)$$

$$[0054] \quad f_t = \sigma(x_t \cdot w_{xh}^f + h_{t-1} \cdot w_{hh}^f + b_h^f)$$

$$[0055] \quad o_t = \sigma(x_t \cdot w_{xh}^o + h_{t-1} \cdot w_{hh}^o + b_h^o)$$

$$[0056] \quad \tilde{c}_t = \tanh(x_t \cdot w_{xh}^c + h_{t-1} \cdot w_{hh}^c + b_h^c)$$

$$[0057] \quad c_t = i_t \otimes \tilde{c}_t + f_t \otimes c_{t-1}$$

$$[0058] \quad h_t = o_t \otimes \tanh(c_t)$$

[0059] 其中 x_t 是t时刻的单元输入, i_t, f_t, o_t 分别表示t时刻的输入门, 遗忘门, 和输出门。 w 和 b 代表3种门的权重参数矩阵和偏置向量。 \tilde{c}_t 为当前时刻t的输入得到的中间状态用于更新当前时刻状态 c_t, h_t 为当前时刻输出。(σ为sigmoid激活函数, tanh为双曲正切激活函数) 所以通过双LSTM, 可以有效的采集到词语的上下文信息, 因此将每个组合嵌入的顺序传递的隐藏输出 \vec{r}_i 和逆序传递的隐藏输出 \vec{r}_i 拼接起来得到组合嵌入的最终隐层表示 $r'_i = [\vec{r}_i, \vec{r}_i]$ 。

[0060] 对于序列标记任务, 考虑相邻标签之间的相关性并对给定的句子联合解码出最佳的标签序列是十分有用的。例如对于带有BIO标签的NER任务中, “B-PER I-PER”是合法的序列, 但是“B-LOC, I-ORG”, “O, I-label”是非法标签序列。因为实体标签B-LOC后应接I-LOC而非“I-ORG”实体标签的首个标签应该是“B-”, 而非“I-”。而使用条件随机场 (CRF) 联合建模标签序列, 而不是单独解码每一标签, 可以有效的解决产生非法标签的问题。因此我们将编码层得到的组合嵌入的隐层表示 r'_i 输入到CRF层中根据所有可能的序列标签 y 给出最终序列概率:

$$[0061] \quad p(y|r'; W, b) \propto \prod_{i=1}^n \varphi_i(y_{i-1}, y_i, r')$$

$$[0062] \quad \varphi_i(y', y, r) = \exp(W_{y', y}^T r_i' + b_{y', y})$$

[0063] 实验选取的评价指标为F1值是通过准确率P以及召回率R计算得到的,具体计算公式如下所示:

$$[0064] \quad P = \frac{TP}{TP + FP} R = \frac{TP}{TP + FN} F1 = \frac{2P \times R}{P + R}$$

[0065] 其中,TP表示判定正确的正例,FP表示负例被判定为正例,FN表示正例被判定为负例。

[0066] 本发明所提的实体识别算法模型在Python 3.6.8、keras2.1.4和Tensorflow 1.14.0的环境下进行实验,训练集和测试集的batch_size为64,epoch为25,为了防止过拟合dropout率为0.2,sequence_length 100,提前停止条件:2个周期验证集准确率没有提升。BERT模型的预训练过程需要大量的算力才能实现,BERT预训练语言模型版本,其常用的两种模型参数如图所示,其中,L表示层数,H表示隐藏层,A是自注意力的头数。本实验使用BERT-Base-Chinese模型版本进行实验,此模型共有12层,隐含层为768,12个头,包含110M个参数。训练的第一步需要在每个批次输入64个句子,并每个字所在词训练得到的词向量与通过BERT模型得到字向量维度拼接得到联合特征表示。本发明在网络训练阶段选取adam函数作为优化器进行迭代训练,每轮训练通过不断降低误差,提高准确率训练模型的参数,首先将组合向量作为输入向量首先输入到BI-LSTM中采集到有效的上下文信息,最后再利用条件随机场作为解码器对模型进行解码,即通过转移概率来得到每个字符最优的标记序列,从而为实体打上类别标签,实现实体识别和分类。通过模型训练,最终可以得到该模型在验证集的准确率可达到93.8%,召回率为90.05%,F值为91.3%,在准确率这一项评价指标上明显优于仅使用字符向量 w_i^{char} 作为特征的模型的准确率为87.1%和仅使用bert字符向量作为特征(而未添加词向量作为词级短语补充)的模型的准确率为89.2%。

[0067] 模型测试阶段:用户通过输入待测试的句子,可以返回给用户json格式的结果,其中包含以下几个信息:识别并抽取出的实体(word)、实体的起始位置(start)、实体的终止位置(end)、实体的类别标签(type),每个类别标签所表示的实际含义可以参见图3。

[0068] 例如用户在控制台输入的待测试语句为“参加城乡居民社会养老保险人员就业后又参加企业职工基本养老保险的,可保留城乡居民社会养老保险关系,具体转移办法按照人力资源社会保障部、财政部《城乡养老保险制度暂行办法》”。识别的结果为: {'entities': [{'word': '城乡居民社会养老保险', 'start': 3, 'end': 12, 'type': 'PRO'}, {'word': '企业职工基本养老保险', 'start': 21, 'end': 30, 'type': 'PRO'}, {'word': '城乡居民社会养老保险关系', 'start': 35, 'end': 46, 'type': 'PRO'}, {'word': '人力资源社会保障部', 'start': 55, 'end': 63, 'type': 'ORG'}, {'word': '财政部', 'start': 64, 'end': 65, 'type': 'ORG'}] } {'word': '城乡养老保险制度暂行办法', 'start': 67, 'end': 78, 'type': 'LAW'}] }。

[0069] 本发明通过词性组合的方式预先建立四险一金领域词典和定义实体类别,并对词典中的实体进行标记,利用Jieba分词工具以及相关算法可以实现对原始的四险一金政策法规文本进行自动标注,从而得到一定规模的标注语料库,减少了人工进行数据标记的成本。在特征提取方面,本发明将BERT预训练作为字向量的特征层和通过Word2Vec模型对四

险一金政策法规分词后文本中的词语特征进行提取,并训练成的词向量拼接得到的联合字向量,不仅可以弥补少量标注样本特征不足和字符语义提取不充分的问题,还为字向量补充了词级短语信息,从而可以在一定程度上提高模型的准确率。本发明可以有效解决四险一金领域命名实体识别任务中标注数据不足以及识别精度不高的问题。

[0070] 以上所述仅为本发明的优选实施例而已,并不用于限制本发明,对于本领域的技术人员来说,本发明可以有各种更改和变化。凡在本发明的精神和原则之内,所作的任何修改、等同替换、改进等,均应包含在本发明的保护范围之内。

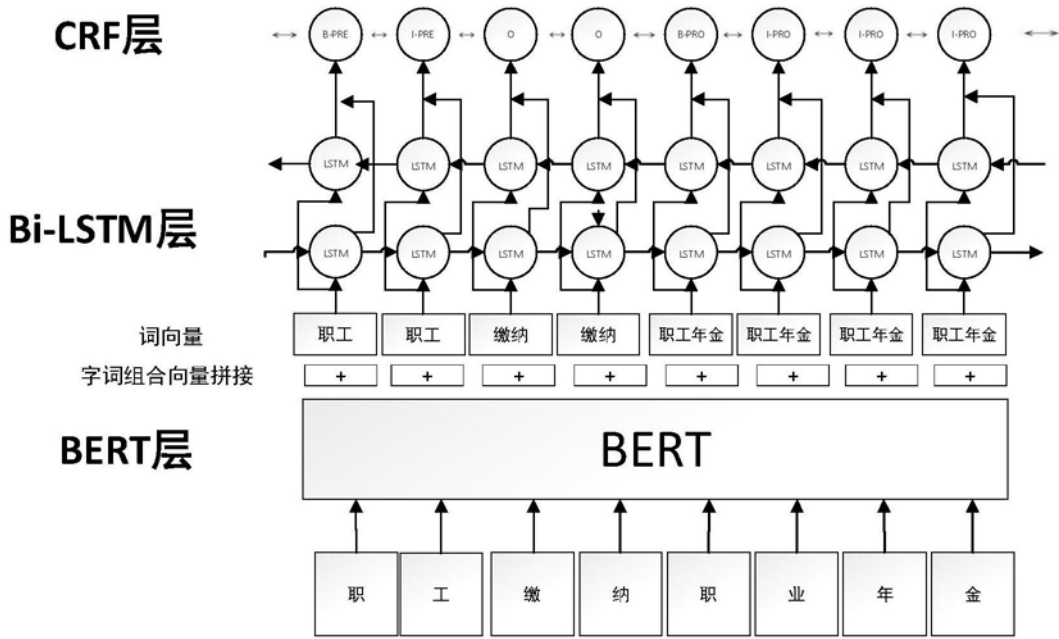


图1

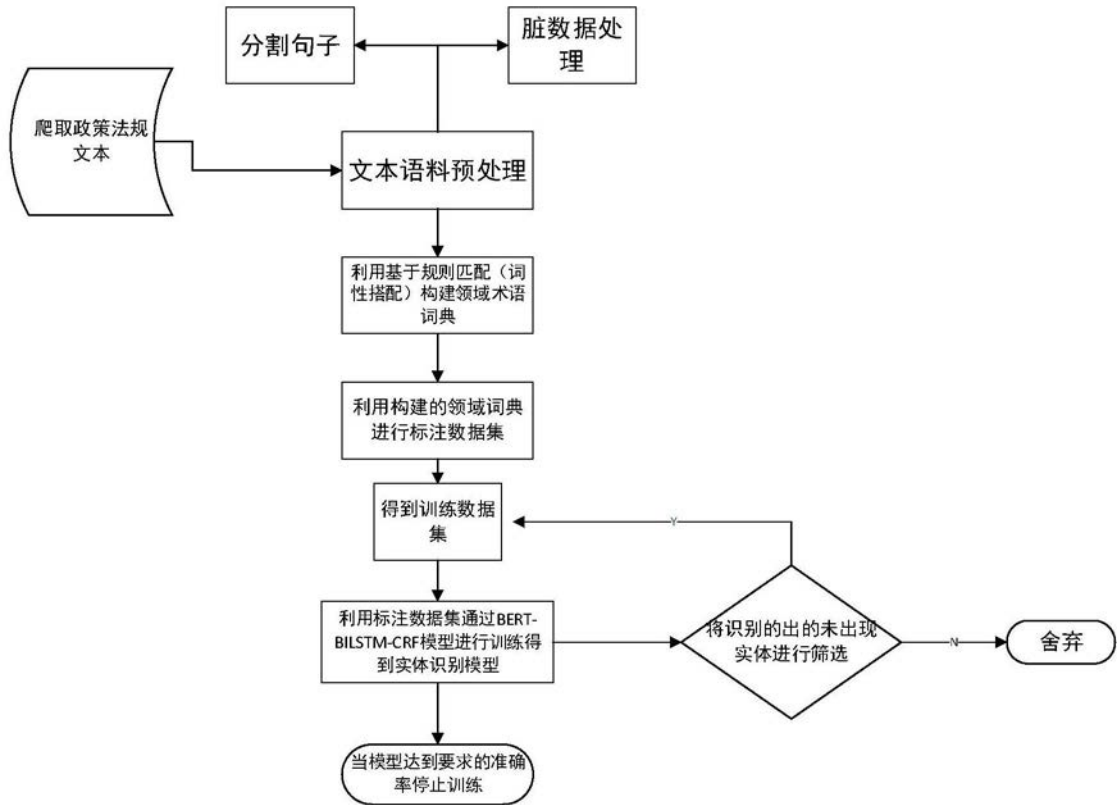


图2

实体标签	描述	例子
PRO	领域术语	工伤保险储备金
PER	人名	王宁
LOC	地名	黑龙江省, 合肥市
ORG	机构名	财政部
LAW	法规	工伤保险条例

图3