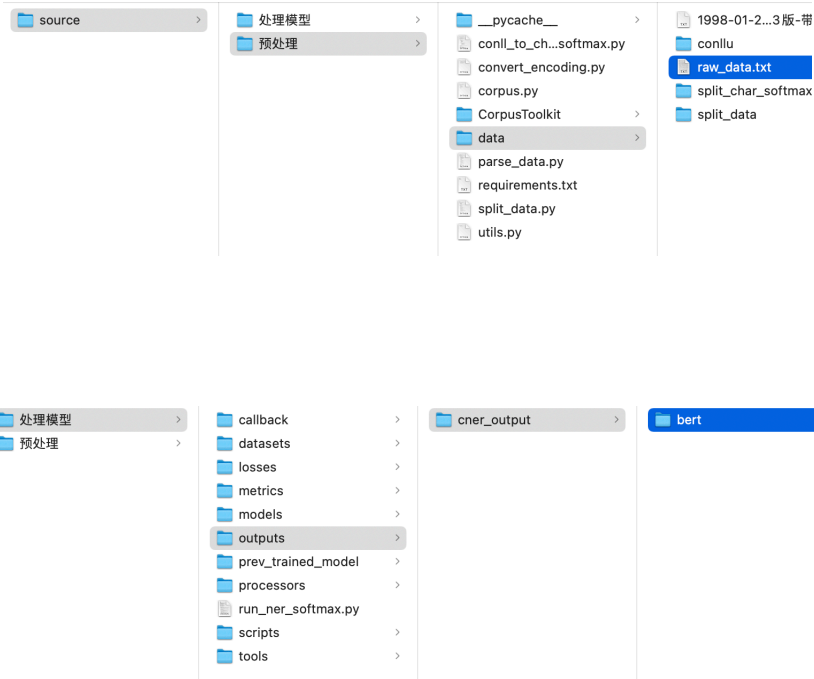


知识工程第一次大作业报告

姓名：赵炫皓 学号:1120200603

项目源码：



1 NER简介

NER（命名实体识别）又称作专名识别，是自然语言处理中常见且基本的一项任务。命名实体识别通常指的是文本中具有特殊意义或指代性强的实体，通常包括人名，地名，机构名，时间，专有名词等。NER系统就是从非结构化的文本中抽取上述实体，并且可以按照业务需求识别出更多类型的实体。

2 NER问题的分解

NER问题目标是从文本中抽取特定需求实体的文本片段。针对这个任务，通常使用基于规则的方法和基于模型的方法。由于基于规则的方法适用性差，维护成本高，因此一般只用于处理简单的NER问题。

2.2 基于模型的方法

从模型的角度来看，NER问题实际上是序列标注问题，序列标注问题指的是模型的输入是一个序列，包括文字，时间等，输出也是一个序列。针对输入序列的每一个单元，输出一个特定的标签。以中文分词任务举例，例如输入序列的每一个单元，输出一个特定的标签。源码中使用的是BIO标注方法，即B表示这个字是词的开始Begin，I表示这个字是词的中间部分。

2.2.1 softmax函数

softmax用于多分类过程中，它将多个神经元输出，映射到(0, 1)区间内，可以看成概率来理解，从而进行多分类。

假设我们有一个数组，V，Vi表示V中的第i个元素，那么这个元素的softmax值就是：

$$S_i = \frac{e^i}{\sum_j e^j}$$

softmax将原来的值通过softmax函数作用映射为（0，1）的值，而这些值的累和是1（满足概率的性质),那么可以将其理解为概率，在最后输出结点时选取概率最大的结点作为我们预测的目标。

当我们对分类的Loss进行改进的时候，我们需要通过梯度下降，每次优化一个step大小的梯度，这个时候要求Loss对每个权重矩阵求偏导，用链式法则。使用softmax函数后，梯度求导过程非常方便。

### 3 基于BERT + softmax模型打造中文NER系统

#### 3.1 明确标注目标

NER可以根据业务需要标注各种不同类型的实体，因此首先要明确需要抽取的实体类型。一般通用场景下，最常提取的是时间，任务，地点以及机构名。作业任务中提取了PERSON,LOCATION,ORGNAZATION三种实体。

#### 3.2 数据以及工具准备

明确任务后需要训练数据和模型工具。首先下载人民日报1998中文标注语料库，进行处理的环境配置如下：

```
1.1.0 =< PyTorch < 1.5.0
cuda=9.0
python3.6+
Click==7.0
joblib==0.13.2
nltk==3.4.5
numpy==1.17.2
pandas==0.25.1
ply==3.11
python-dateutil==2.8.0
pytz==2019.2
scikit-learn==0.21.3
scipy==1.3.1
six==1.12.0
tokenizer-tools==0.15.3
```

#### 3.3 数据的预处理

人民日报1998语料库下载完毕后，解压打开“199801.txt”这个文件，首先用预处理中python文件convert\_encoding.py将编码格式转换成UTF-8格式在data文件夹中得到raw\_data.txt，方便使用文字编辑器读取。我们需要的提取的实体是人名，地名，组织机构名，根据1998语料库的词性标记说明，对应的词性依次为nr、ns、nt。运行python文件parse\_data.py对语料库进行处理。处理的要点有四则：

1. 1998语料库标注人名时，将姓和名分开标注，因此需要合并姓名，比如 兰/nrf 红光/nrg 应该合并成 兰红  
光/nr
2. 中括号括起来的几个词表示大粒度分词，表意能力更强，需要将括号内内容合并，比如 [中央/n 人民/n 广  
播/vn 电台/n]nt 应该合并成 中央人民广播电台/nt
3. 时间合并,例如将“1997年/t 3月/t” 合并成“1997年3月/t”
4. 全角字符统一转为半角字符，尤其是数字的表示，比如： 1 9 9 8 年/t 应该转换成 1998年/t

运行parse.py文件，raw\_data.txt转换为conllu格式的文件data\_False-True-True-True-True-True-False.conllu  
通过语料预处理得到结果如图1，2所示。

19980101-01-001-001/m 迈向/v 充满/v 希望/n 的/u 新/a 世纪/n —/w 一九九八年/t 新年/t 讲话  
n (/w 附/v 图片/n 1/m 张/q ) /w  
19980101-01-001-002/m 中共中央/nt 总书记/n 、 /w 国家/n 主席/n 江/nr 泽民/nr  
19980101-01-001-003/m (/w 一九九七年/t 十二月/t 三十一日/t ) /w  
19980101-01-001-004/m 12月/t 31日/t , /w 中共中央/nt 总书记/n 、 /w 国家/n 主席/n 江/nr  
泽民/nr 发表/v 1998年/t 新年/t 讲话/n 《/w 迈向/v 充满/v 希望/n 的/u 新/a 世纪/n 》  
/w 。 /w (/w 新华社/nt 记者/n 兰/nr 红光/nr 摄/Vg ) /w  
19980101-01-001-005/m 同胞/n 们/k 、 /w 朋友/n 们/k 、 /w 女士/n 们/k 、 /w 先生/n 们/k :  
w  
19980101-01-001-006/m 在/p 1998年/t 来临/v 之际/f , /w 我/r 十分/m 高兴/a 地/u 通过/p  
[中央/n 人民/n 广播/vn 电台/n]nt 、 /w [中国/ns 国际/n 广播/vn 电台/n]nt 和/c [中央/n 电  
视台/n]nt , /w 向/p 全国/n 各族/r 人民/n , /w 向/p [香港/ns 特别/a 行政区/n]ns 同胞/n 、  
/w 澳门/ns 和/c 台湾/ns 同胞/n 、 /w 海外/s 侨胞/n , /w 向/p 世界/n 各国/r 的/u 朋友/n  
们/k , /w 致以/v 诚挚/a 的/u 问候/vn 和/c 良好/a 的/u 祝愿/vn ! /w  
19980101-01-001-007/m 1997年/t , /w 是/v 中国/ns 发展/vn 历史/n 上/f 非常/d 重要/a 的  
u 很/d 不/d 平凡/a 的/u 一/m 年/q 。 /w 中国/ns 人民/n 决心/d 继承/v 邓小平/nr 同志  
/n 的/u 遗志/n 、 /w 继续/v 把/p 建设/v 有/v 中国/ns 特色/n 社会主义/n 事业/n 推向/v 前  
进/v 。 /w [中国/ns 政府/n]nt 顺利/ad 恢复/v 对/p 香港/ns 行使/v 主权/n , /w 并/c 按照/p  
7w 一国两制/j ” /w 、 /w 高度/d 自治/v 的/u 方针/n 保持/v 香港  
ns 的/u 繁荣/an 稳定/an 。 /w [中国/ns 共产党/n]nt 成功/a 地/u 召开/v 了/u 第十五/m 次/q  
全国/n 代表大会/n , /w 高举/v 邓小平理论/n 伟大/a 旗帜/n , /w 总结/v 百年/m 历史/n , /w  
展望/v 新/a 的/u 世纪/n , /w 制定/v 了/u 中国/ns 跨/v 世纪/n 发展/v 的/u 行动/vn 纲领

图1.预处理前

19980101-01-001-001/m 迈向/v 充满/v 希望/n 的/u 新/a 世纪/n —/w 一九九八年新年/t 讲话/n ( /w 附/v 图片/n 1/m 张/q ) /w  
19980101-01-001-002/m 中共中央/nt 总书记/n 、 /w 国家/n 主席/n 江泽民/nr  
19980101-01-001-003/m (/w 一九九七年十二月三十一日/t ) /w  
19980101-01-001-004/m 12月31日/t , /w 中共中央/nt 总书记/n 、 /w 国家/n 主席/n 江泽民/nr 发表  
/v 1998年新年/t 讲话/n 《/w 迈向/v 充满/v 希望/n 的/u 新/a 世纪/n 》 /w 。 /w (/w 新华社/  
nt 记者/n 兰红光/nr 摄/Vg ) /w  
19980101-01-001-005/m 同胞/n 们/k 、 /w 朋友/n 们/k 、 /w 女士/n 们/k 、 /w 先生/n 们/k : /w  
19980101-01-001-006/m 在/p 1998年/t 来临/v 之际/f , /w 我/r 十分/m 高兴/a 地/u 通过/p 中央  
人民广播电台/nt 、 /w 中国国际广播电台/nt 和/c 中央电视台/nt , /w 向/p 全国/n 各族/r 人民/n  
/w 向/p 香港特别行政区/ns 同胞/n 、 /w 澳门/ns 和/c 台湾/ns 同胞/n 、 /w 海外/s 侨胞/n , /w  
/w 向/p 世界/n 各国/r 的/u 朋友/n 们/k , /w 致以/v 诚挚/a 的/u 问候/vn 和/c 良好/a 的/u  
祝愿/vn ! /w  
19980101-01-001-007/m 1997年/t , /w 是/v 中国/ns 发展/vn 历史/n 上/f 非常/d 重要/a 的/u 很  
/d 不/d 平凡/a 的/u 一/m 年/q 。 /w 中国/ns 人民/n 决心/d 继承/v 邓小平/nr 同志/n 的/u  
遗志/n , /w 继续/v 把/p 建设/v 有/v 中国/ns 特色/n 社会主义/n 事业/n 推向/v 前进/v 。 /w  
[中国政府/nt] 顺利/ad 恢复/v 对/p 香港/ns 行使/v 主权/n , /w 并/c 按照/p 7w 一国两制/j ” /w  
、 /w 高度/d 自治/v 的/u 方针/n 保持/v 香港/ns 的/u 繁荣/an 稳定/an 。 /w 中国共产党/nt 成功/a 地/u 召开/v 了/u 第十五/m 次/q 全国/n 代表大会/n , /w 高  
举/v 邓小平理论/n 伟大/a 旗帜/n , /w 总结/v 百年/m 历史/n , /w 展望/v 新/a 的/u 世纪/n , /w  
制定/v 了/u 中国/ns 跨/v 世纪/n 发展/v 的/u 行动/vn 纲领/n 。 /w

图2.预处理后

### 3.4 模型训练

根据NER任务需求以及BERT+softmax训练要求，训练模型需要3个步骤：1，确定标签体系 2，处理训练数据文件 3，训练模型

#### 3.4.1 确定标签体系

大部分情况下，标签体系越复杂准确度也越高，但相应的训练时间也会增加。本次处理使用BIOESX标签体系。

Tokens	IO	BIO	BMEWO	BMEWO+
昨	O	O	O	O
天	O	O	O	O
,	O	O	O	O_PERSON
李	I_PERSON	B_PERSON	B_PERSON	B_PERSON
晓	I_PERSON	I_PERSON	M_PERSON	M_PERSON
明	I_PERSON	I_PERSON	E_PERSON	E_PERSON
前	O	O	O	PERSON_O
往	O	O	O	O_LOCATION
上	I_LOCATION	B_LOCATION	B_LOCATION	B_LOCATION
海	I_LOCATION	I_LOCATION	E_LOCATION	E_LOCATION
。	O	O	O	LOCATION_O

图3.标签体系

### 3.4.2 数据标注处理

得到data\_False-True-True-True-True-True-False.conllu文件后，用split\_data.py把文件分为训练集，验证集，和测试集,分别为split\_data文件夹的dev.conllu,test.conllu,train.conllu。然后使用conll\_to\_char\_softmax.py文件将conllu格式文件用tokenizer\_tools库所自带的BILUOEncoderDecoder将数据转化为BILUO标注再转化为BIO标注。得到如图所示的标注文件。

```

train.char.bmes
767 太 · · 0
768 阳 · · 0
769 ( · · 0
770 中 · · 0
771 国 · · 0
772 画 · · 0
773 ) · · 0
774 ( · · 0
775 苗 · · B-PER
776 再 · · I-PER
777 新 · · I-PER
778 李 · · B-PER
779 翔 · · I-PER
780 ) · · 0
781
782 昆 · · B-LOC
783 明 · · I-LOC
784 : · · 0
785 实 · · 0
786 施 · · 0
787 优 · · 0
788 惠 · · 0

```

图4.BIO标注





```

03/28/2022 00:01:33 - INFO - root - Evaluate the following checkpoints: ['/Users/myfile/Desktop/BERT-NER-Pytorch-master/outputs/cner_output/bert']
03/28/2022 00:01:34 - INFO - root - Loading features from cached file /Users/myfile/Desktop/BERT-NER-Pytorch-master/datasets/cner/cached_soft-dev_bert-base-chinese_512_cner
03/28/2022 00:01:35 - INFO - root - ***** Running evaluation *****
03/28/2022 00:01:35 - INFO - root - Num examples = 171
03/28/2022 00:01:35 - INFO - root - Batch size = 24
[Evaluating] 8/8 [=====] 12.5s/step2
03/28/2022 00:03:15 - INFO - root -

03/28/2022 00:03:15 - INFO - root - ***** Eval results *****
03/28/2022 00:03:15 - INFO - root - acc: 0.9385 - recall: 0.9639 - f1: 0.9510
- loss: 0.0227
03/28/2022 00:03:15 - INFO - root - ***** Entity results *****
03/28/2022 00:03:15 - INFO - root - ***** LOC results *****
03/28/2022 00:03:15 - INFO - root - acc: 0.9141 - recall: 0.9744 - f1: 0.9433

03/28/2022 00:03:15 - INFO - root - ***** ORG results *****
03/28/2022 00:03:15 - INFO - root - acc: 0.9524 - recall: 0.8333 - f1: 0.8889

03/28/2022 00:03:15 - INFO - root - ***** PER results *****
03/28/2022 00:03:15 - INFO - root - acc: 0.9860 - recall: 0.9658 - f1: 0.9758

```

图6.BERT + Softmax训练结果

#### 4. 总结

本次任务基于BERT + Softmax模型，完成了对于人民日报1998版语料的NER。在本次任务中，首先对python语言进行复习，并熟悉pytorch和pandas等库，在此基础上理解了NLP任务中比较基本的NER任务的处理思路与处理方法。