

bert-bilstm-crf命名实体识别

1、运行环境

- python3.6
- tensorflow=1.10.0

2、模型文件

- 其中文件out-bert, output-wwm, output-wwm-ext是输出后得到的模型。

(1) bert文件夹

存放有bert模型, bert源代码。其中bert-wwmm, bert-wwm-ext, chinese_wwm_L-12-H-768_A-12.zip, 预训练好的模型存储到百度网盘中。

(2) db文件夹, 是一些常用的数据库处理方法、

(3) NERdata文件夹和NLPCC2016KBQA文件夹, 是使用的数据集, 将存储在百度网盘中。

(4) bert_lstm_ner.py文件

该文件夹包含数据参数, bert_lstm_ner模型, 该文件可以用于训练。

(5) data_process.py文件

用于语料库的处理, 全部处理成小于max_seq_length的序列, 这样可以避免解码出现不合法的数据或者在最后算结果的时候出现out of range 的错误。

(6) lstm_crf_layer.py文件

是BLSTM-CRF网络。

(7) terminal_predict.py文件

基于命令行的在线预测方法, 需要基本的参数配置, 以及预测模型的调用。

3、模型训练

模型训练需要, 基于bert_lstm_ner.py文件, 对模型进行训练。

其中运行的主要函数是main函数, 需要自己修改参数, 其运行方法主要如下。

```
1 task_name = "NER"
2 do_train = True
3 do_eval = True
```

```

4 do_predict = True
5 data_dir = "NERdata"
6 max_seq_length = 128
7 train_batch_size = 32
8 learning_rate = 2e-5
9 num_train_epochs = 3.0
10 output_dir = "./output_qa/result_dir_ner/"
11 tf.app.run()
12 # filter model
13
14 if FLAGS.filter_adam_var:
15     adam_filter(FLAGS.output_dir)
16

```

其中一些常用参数，需要训练的时候自己设置。

4、模型预测

基于terminal_predict.py文件，实现模型预测。其中需要设置模型路径，语料路径。其使用模型的方法如下：

```

1 sentence = '刚开增值税普通发票时，把发票号码弄错了怎么办？打印出来的号和实际发
票上的号不符，怎么办？'
2 predict(sentence)
3 def convert(line):
4     feature = convert_single_example(0, line, label_list, FLAGS.max_seq_leng
th, tokenizer, 'p')
5     input_ids = np.reshape([feature.input_ids],(batch_size, FLAGS.max_seq_le
ngth))
6     input_mask = np.reshape([feature.input_mask],(batch_size, FLAGS.max_seq_
length))
7     segment_ids = np.reshape([feature.segment_ids],(batch_size, FLAGS.max_se
q_length))
8     label_ids = np.reshape([feature.label_ids],(batch_size, FLAGS.max_seq_len
gth))
9     return input_ids, input_mask, segment_ids, label_ids
10
11
12 file = open('./question_entity_result.txt', mode='r', encoding='utf-8')
13 file_write = open('./question_entity_result1.txt', mode='a', encoding='u
tf-8')
14 sentences = file.readlines()
15 sentences_lsit = []
16 for line in sentences:
17     sentences_lsit.append(line)

```

```

18 from tqdm import tqdm
19
20 with graph.as_default():
21     print(id2label)
22     start = datetime.now()
23     for i in tqdm(range(int(len(sentences_lsit)/8))):
24         sentence = sentences_lsit[i*8]
25         sentence = tokenizer.tokenize(sentence)
26         # print('your input is:{}'.format(sentence))
27         input_ids, input_mask, segment_ids, label_ids = convert(sentence)
28
29         feed_dict = {input_ids_p: input_ids,
30                      input_mask_p: input_mask,
31                      segment_ids_p: segment_ids,
32                      label_ids_p: label_ids}
33         # run session get current feed_dict result
34         pred_ids_result = sess.run([pred_ids], feed_dict)
35         pred_label_result = convert_id_to_label(pred_ids_result, id2label)
36         #todo: 组合策略
37         result = strage_combined_link_org_loc(sentence, pred_label_result[0])
38
39         file_write.write(sentences_lsit[i * 8])
40         file_write.write(sentences_lsit[i * 8 + 1])
41         file_write.write(sentences_lsit[i * 8 + 2])
42         file_write.write(sentences_lsit[i * 8 + 3])
43         file_write.write(sentences_lsit[i * 8 + 4])
44         file_write.write(sentences_lsit[i * 8 + 5])
45         file_write.write('bert_lstm_crf_book:' + str([result])+'\n')
46         file_write.write(sentences_lsit[i * 8 + 6])
47         file_write.write(sentences[i * 8 + 7])
48

```

其中包括，单个文本的预测，以及批量的预测。