BERT NER 최종 발표

NLP209

프로젝트 주제

주제: BERT를 이용한 키워드 추출

목적1) BERT모델을 이용하여 BIO tagging data로 fine-tuning을 하여 NER task를

수행하도록 해 봄으로써 모델 이해

- 2) 기존의 biobert는 tensorflow로 이루어져 있어 코드를 이해하는데 어려움 多
- → 우리는 이를 keras로 수정하여

사용한 Data: BC5CDR-disease의 train, validation, test data

🤳 train - Windows 메모장	☐ train_dev - Windows 메모장 ☐ test - Windows 메모장
파일(F) 편집(E) 서식(O) 보기(V)	파일(F) 편집(E) 서식(O) 보기(V) 도 파일(F) 편집(E) 서식(O) 보기(V) .
degree O	Selegiline O Torsade B
of O	induced O de I
orthostatic B	postural B pointes I
hypotension I	hypotension I ventricular B
occurred O	in O tachycardia I
with O	Parkinson B during O
standing O	s I low O
. О	disease I dose O
	intermittent O
Orthostatic B	a O dobutamine O
hypotension I	longitudinal O treatment O
was O	on O

이틀간 진행한 내용

- data(BC5CDR-disease)를 input data, target data 형식으로 변환
 - 1. data를 읽어와 단어의 총 개수가 30을 넘지 않도록 data 자르기 (B, I 사이에 잘리지 않게) → [labels, sentence] 형식의 examples 생성

데이터를 table로 불러올때

이틀간 진행한 내용

- data(BC5CDR-disease)를 input data, target data 형식으로 변환
 - 2. 각각의 example들을 WordPieceTokenizer로 토큰화
 - → biobert의 Tokenization을 우리 코드 상에서 구현.
 - → 이를 이용하여 tokenize 진행

E LICEITOE	minuai	text = convert_to_unicode(text)	
■ README.md	Update README.md		
initpy	initial	output_tokens = [] for token in whitespace_tokenize(text):	
create_pretraining_data.py	initial	chars = list(token)	
		<pre>if len(chars) > self.max_input_chars_per_word: output_tokens.append(self.unk_token)</pre>	
download.sh	download script	continue	
extract_features.py	initial	Tokenization의	
modeling.py	initial	is_bad = False Wordpiecetokenizer	
modeling_test.py	initial	sub_tokens = [] 부분 while start < <mark>len</mark> (chars):	
optimization.py	initial	end = len(chars)	
	initial	cur_substr = None while start < end:	
optimization_test.py	IIIIIIdi	substr = "".join(chars[start:end])	
requirements.txt	fix requirements	if start > 0:	
run_classifier.py	initial	substr = "##" + substr if substr in self.vocab:	
E Tun_clussinci.py		cur_substr = substr	
run_ner.py	initial	break	
run_pretraining.py	initial	end -= 1	
		if cur_substr is None:	
run_qa.py	initial	is_bad = True break	
run_re.py	initial	sub_tokens.append(cur_substr)	
sample_text.txt	initial	start = end	
Biobert □ Biobert □	_ initial	if is_bad: output_tokens.append(self.unk_token)	
tokenization.py를 tokenization.py를 으리 코드에 올린 스	i 前进 己	else:	
구의 포트에 달린 구	- · · · · _	output_tokens.extend(sub_tokens)	
■ tokenization_test.py run_ner.py 코드 =	는정·	return output_tokens	

이틀간 진행한 내용

- data(BC5CDR-disease)를 input data, target data 형식으로 변환
 - 3. convert_single_example함수를 이용하여 input_ids, segment_ids, label_ids
 - → input_ids, segment_ids: train_x
 - → label_ids: train_y

```
def convert_single_example(ex_index, example, label_list, max_seq_length,tokenizer):
   #vocab=path+'/vocab.txt' #구드path로 변경!
                                                                                    - Mask 변수 삭제
  # wordtt= FullTokenizer(vocab) #class형태로 선언
   Tabel map = {}
   for (i, label) in enumerate(label_list.1):
       label_map[label] = i
                                                                                   - [PAD]: 0
  #label_map: {'B': 1, 'l': 2, '0': 3, 'X': 4, '[CLS]': 5, '[SEP]': 6}
  # print("lab<mark>el map: ". label map)</mark>
  # with open(os.path.join(FLAGS.output dir. 'label2id.pkl').'wb') as w:
     pickle.dump(label map.w)
   -textlist = example[1].split(' ') #example.text.split(' ')을 index형태로 바꿈
   labellist = example[0].split(' ') #example.label.split(' ')을 index형태로 바꿈
   tokens = []
   Tabels = \Pi
   for i. word in enumerate(textlist):
                                                                                    - "케라스" → B/I/O 中 B
       token = tokenizer.tokenize(word) #사용!!
       tokens.extend(token)
       label 1 = labellist[i]
                                                                                     → "케라스"를 tokenization하면 [케, ##라, ##스]
       for m in range(len(token)):
          if m==∩:
              labels.append(label_1)
           else:
              Tabels.append("X")
   tokens = tokens[0:(max seg length - 2)]
       labels = labels[0:(max_seq_length - 2)]
   ntokens = []
   segment_ids = []
   label_ids = []
   ntokens.append("[CLS]")
   segment_ids.append(0)
   # append("0") or append("[CLS]") not sure!
   label_ids.append(label_map["[CLS]"])
   for i, token in enumerate(tokens):
       ntokens.append(token)
```

- 객체 사용 X, 인덱스 사용

 \rightarrow labels=[B, X, X]

```
vocab=path+'/vocab.txt' → Biobert의 vocab
label_list=['B'.'l'.'0'.'X'.'[CLS]'.'[SEP]'] #101: [CLS]. 102: [SEP]
max_seq_length= 128
                        → max seq length
tokenizer= FullTokenizer(vocab)
                                   → tokenizer
result input ids=[]
result_seg_ids=[]
y=[]
for i, example in enumerate(ex):
  input_ids, segment_ids, label_ids*convert_single_example(i , example, label_list,max_seq_length, tokenizer)
  result_input_ids.append(input_ids)
  result_seg_ids.append(segment_ids)
  y.append(label_ids)
train_x=[]
train v=[]
train_x.append(np.array(result_input_ids))
train_x.append(np.array(result_seg_ids))
train_y.append(np.array(y))
```

train_x data

```
[arrav([[ 101, 14516, 27412, ....
                                                     01.
                                              0.
                                                     0],
          101, 11350, 131, ...,
          101. 1106, 3531, ...,
                                                     0],
          101. 175, 1377, ...,
                                                     0],
          101. 1103. 3469. ....
                                                     01.
                                                     011).
       [ 101, 1292, 2686, ...,
array([[0, 0, 0, ..., 0, 0, 0].
       [0, 0, 0, \ldots, 0, 0, 0],
       [0, 0, 0, \dots, 0, 0, 0]
       [0, 0, 0, \dots, 0, 0, 0]
       [0, 0, 0, \ldots, 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0]])]
```

train_y data

```
[array([[5, 3, 4, ..., 0, 0, 0], [5, 3, 3, ..., 0, 0, 0], [5, 3, 3, ..., 0, 0, 0], ..., [5, 3, 4, ..., 0, 0, 0], [5, 3, 3, ..., 0, 0, 0], [5, 3, 3, ..., 0, 0, 0]])]
```

{'B': 1, 'I': 2, 'O': 3, 'X': 4, '[CLS]': 5, '[SEP]': 6}



Bert 모형의 자체 Masking 된 텐서들을 풀어주기 위한 NonMasking

```
from keras, lavers import Laver
class NonMasking(Layer):
   def __init__(self, **kwargs):
        self.supports_masking = True
        super(NonMasking, self).__init__(**kwargs)
   def build(self, input_shape):
        input shape = input shape
   def compute_mask(self, input, input_mask=None):
        return None
   def call(self, x, mask=None):
        return x
   def get_output_shape_for(self, input_shape):
        return input_shape
```

```
class MyLaver(Laver):
 def init (self.seg len. **kwargs):
   self.sea len = sea len
   self.supports_masking = True
   super(MyLaver, self), init (**kwargs)
 def build(self, input_shape):각 token별로 outptu값이 layer 개수만큼 나오게 하기 위해
   self.W = self.add_weight(name='kernel', shape=(768,7),initializer='uniform',trainable=True)
   super(MyLayer, self).build(input shape)
 def call(self. x):
   x = K.reshape(x, shape*(-1,128,768))
                                       Shape 를 max seg에 맞게
   x = K.dot(x. self.W)
       # Ner은 각 token별로 softmax를 취해서 Tabel을 유추하므로 이 부분은 제외
       \# x = K.permute dimensions(x. (2.0.1))
       # self.start_logits, self.end_logits = x[ 0], x[1]
       \#x = K.reshape(x. shape=(-1, 128, 7))
                                         Token 별로 softmax
   self.logits = K.softmax(x, axis= -1)
   #self.logits = K.reshape(logit.shape=(-1, 128*7))
   return self.logits
```

def compute_output_shape(self, input_shape):
 return (input_shape[0], self.seq_len)

```
def get_bert_finetuning_model(model):
    inputs = model.inputs[:2] #segment, token 두개의 input
   dense = model.output
   x = NonMasking()(dense)
   output_layer = MyLayer(128)(x)
   model = keras.models.Model(inputs, output_layer)
   model.compile(
      optimizer=RAdam(learning_rate=LR, weight_decay=0.001),
      loss="sparse_categorical_crossentropy",
      metrics=["accuracy"])
     #loss='categorical crossentropy'.
     #metrics=['categorical_accuracy'])
    return model
```

우리의 layer를 쓰기 위한 함수

```
sess = K.get_session()
uninitialized_variables = set([i.decode('ascii') for i in sess.run(tf.report_uninitialized_variables())])
init = tf.variables_initializer([v for v in tf.global_variables() if v.name.split(':')[0] in uninitialized_variables])
sess.run(init)

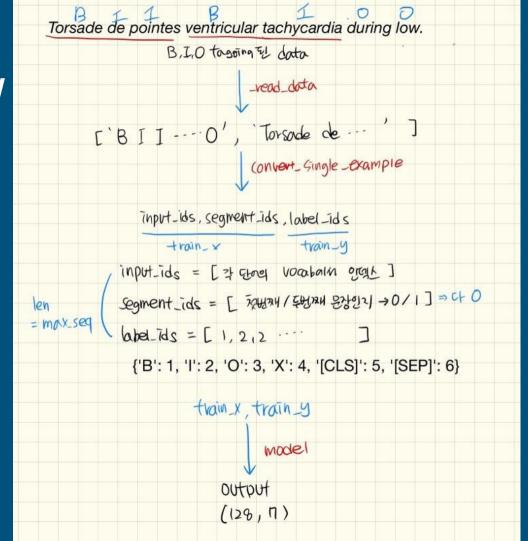
bert_model = get_bert_finetuning_model(model)
bert_model.summary()
# bert_model.fit 함수에 들어가긴 하는데, error가 남
train_y_hot_squ = np.reshape(train_y_hot, (-1,128+7))
history = bert_model.fit(train_x, train_y_hot_squ, batch_size=10, validation_split=0.05, shuffle=False, verbose=1)
```

Model summary

Trainable params: 107,430,144
Non-trainable params: 0

ti (None,	128,	768)	2362368	Encoder-11-FeedForward-Norm[0][0]
ti (None,	128,	768)	0	Encoder-12-MultiHeadSelfAttention
ti (None,	128,	768)	0	Encoder-11-FeedForward-Norm[0][0] Encoder-12-MultiHeadSelfAttention
ti (None,	128,	768)	1536	Encoder-12-MultiHeadSelfAttention
or (None,	128,	768)	4722432	Encoder-12-MultiHeadSelfAttention
t (None,	128,	768)	0	Encoder-12-FeedForward[0][0]
dd (None,	128,	768)	0	Encoder-12-MultiHeadSelfAttention Encoder-12-FeedForward-Dropout[0]
_a (None,	128,	768)	1536	Encoder-12-FeedForward-Add[0][0]
(None,	128,	768)	0	Encoder-12-FeedForward-Norm[0][0]
(None,	128,	768)	5376	non_masking_16[0][0]
1	ti (None, ti (None, or (None, t (None, dd (None, La (None,	ti (None, 128, ti (None, 128, or (None, 128, t (None, 128, dd (None, 128, (None, 128,	ti (None, 128, 768) ti (None, 128, 768) or (None, 128, 768) t (None, 128, 768) dd (None, 128, 768) La (None, 128, 768) (None, 128, 768)	ti (None, 128, 768)

Data flow



학습 결과

Model.fit 과정에서 문제가 있어, 아직 확인 못함!

1. biobert pre-trained weight를 colab에 올리는 과정

```
config_path = os.path.join(pretrained_path, 'bert_config.json')
checkpoint path = os.path.join(pretrained path. 'model.ckpt-1000000')
vocab_path = os.path.join(pretrained_path, 'vocab.txt')
laver_num = 12
model = load trained model from checkpoint( #사전학습된 모델 불러오기
    config_path.
                                                                         Traceback (most recent call last)
                                    DataLossError
    checkpoint_path,
                                   <ipvthon-input-24-5a585d21b790> in <module>()
    training= True.
                                              training= True.
                                              trainable=True,
    trainable=True.
                                    ---> 14
                                              seq_len=SEQ_LEN.)
    seg Ten=SEQ LEN.)
                                    /usr/local/lib/python3.6/dist-packages/tensorflow_core/python/pywrap_tensorflow_internal.py in __init__(self, file
                                       883
                                              def __init__(self, filename):
                                                  this = _pywrap_tensorflow_internal.new_CheckpointReader(filename)
                                    --> 885
                                       888
                                                  trv:
                                       887
                                                     self.this.append(this)
                                    DataLossError: block checksum mismatch
```

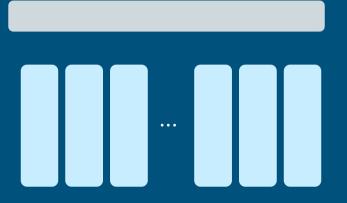
- 1. biobert pre-trained weight를 colab에 올리는 과정
 - * 원인
 - : 이미 이전에 pre-trained weight로 fine-tuning을 한번 진행했었는데,
 - 그 weight를 다시 사용하니 checkpoint에 충돌 발생
 - * 해결
 - : bioBert의 weight를 다시 새롭게 다운받아서 사용

2. Model.fit을 하는 과정에서

```
/usr/local/lib/pvthon3.6/dist-packages/tensorflow_core/pvthon/client/session.pv in _call__(self,
                ret = tf_session.TF_SessionRunCallable(self._session._session,
   1470
   1471
                                                          self._handle, args,
-> 1472
                                                          run_metadata_ptr)
   1473
                if run metadata:
   1474
                  proto_data = tf_session.TF_GetBuffer(run_metadata_ptr)
InvalidArgumentError: 2 root error(s) found.
 (0) Invalid argument: Incompatible shapes: [10] vs. [10.128]
         [[{{node metrics_5/acc/Equal}}]]
         [[loss_5/mul/_43325]]
  (1) Invalid argument: Incompatible shapes: [10] ValueError
                                                                                             Traceback (most recent call last)
         [[{{node metrics_5/acc/Equal}}]]
                                                      <ipvthon-input-133-4d28eecb8ef6> in <module>()
                                                           7 bert_model.summary()
successful operations.
                                                           8 # bert model.fit 함수에 들어가진 하는데, error가 남
0 derived errors ignored.
                                                      ----> 9 history = bert model.fit(train x, train y, batch size=10, validation split=0.05, shuffle=False, verbose=1)
                                                      /usr/local/lib/python3.6/dist-packages/keras/engine/training.utils.py in standardize_input_data(data, names, shapes, check_batch_axis, exception_prefix)
                                                                                       ': expected ' + names[i] + ' to have shape ' +
                                                         140
                                                                                       str(shape) + ' but got array with shape ' +
                                                      --> 141
                                                                                       str(data_shape))
                                                         142
                                                                return data
                                                         143
                                                      ValueError: Error when checking target: expected my_layer_9 to have shape (1,) but got array with shape (128,)
```

- 2. Model.fit을 하는 과정에서
 - 원인

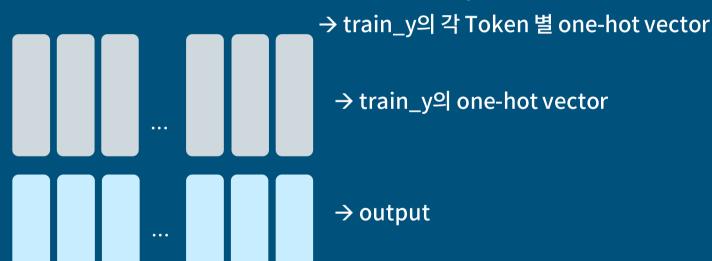
: cross-entropy할때



→ train_y

→ output

2. Model.fit을 하는 과정에서 → 시도 1 train_y의 모든 vector를 one-hot vector로



2. Model.fit을 하는 과정에서

실패한 원인

: 이전 - loss를 token 별로 X, input 문장 별로 구했음

: NER- 문장을 구성하는 token별로 구해야 해서

→ Cross entropy를 사용하지 않고, loss/ matrix에 대한 class를 따로 지정해서 한번 더 해볼 예정

최종 결과 알고리즘 코드

https://drive.google.com/open?id=1dC5csukxyu33w3SvQulHpADI3YflMlDT