

Web Lab2

PB1711585 张永停

一、实验要求

- 对于文档进行适当的预处理
- 选取合适的模型对文本进行建模，并在训练集上进行关系抽取模型训练。
- 在在线平台提交结果验证关系抽取模型的准确率。

二、算法

(一)数据处理

- 从数据集中提出句子以及关系，并将关系转为数字编码

```
1 pattern_text = re.compile(r"\".*\"")
2 pattern_rel = re.compile(r".*\(")
3 pattern_ent1 = re.compile(r"\"(.*\","")
4 pattern_ent1 = re.compile(r"\"(.*\","")
5 info = []
6 with open(train_data_path, "r") as load_f:
7     info = []
8     single_data = {}
9     i = 1
10    for line in load_f.readlines():
11        if i % 2 == 1:
12            single_data = {}
13            single_data['text'] = pattern_text.findall(line)[0]
14            [1:-1]
15            else:
16                single_data['rel'] = pattern_rel.findall(line)[0][:-1]
17                single_data['ent1'] = pattern_ent1.findall(line)[0]
18                [1:-1]
19                single_data['ent2'] = pattern_ent1.findall(line)[0]
20                [1:-1]
21            i += 1
22            info.append(single_data)
```

- 数据划分，从 train.txt 中按7: 2.4: 0.6划分train, test, develop
- 训练时数据准备以及mask选择
 - 使用bert的 tokenizer 进行 encode
 - 同时判断长度进行 padding

```
1 def load_train(data_path):
2     rel2id, id2rel = map_id_rel()
3     max_length=128
4     tokenizer =
5     BertTokenizer.from_pretrained(r'../../data/Sqrti/bert-uncase/bert-
6     base-uncased')
7     train_data = {}
8     train_data['label'] = []
```

```

7     train_data['mask'] = []
8     train_data['text'] = []
9
10    with open(data_path, 'r') as load_f:
11        for line in load_f.readlines():
12            dic = json.loads(line)
13            if dic['rel'] not in rel2id:
14                train_data['label'].append(0)
15            else:
16                train_data['label'].append(rel2id[dic['rel']])
17            #sent=dic['ent1']+dic['ent2']+dic['text']
18            sent = dic['text']
19            indexed_tokens = tokenizer.encode(sent,
20 add_special_tokens=True)
21            avai_len = len(indexed_tokens)
22            while len(indexed_tokens) < max_length:
23                indexed_tokens.append(0) # 0 is id for [PAD]
24            indexed_tokens = indexed_tokens[: max_length]
25            indexed_tokens =
26 torch.tensor(indexed_tokens).long().unsqueeze(0) # (1, L)
27
28            # Attention mask
29            att_mask = torch.zeros(indexed_tokens.size()).long() # (1,
30 L)
31            att_mask[0, :avai_len] = 1
32            train_data['text'].append(indexed_tokens)
33            train_data['mask'].append(att_mask)
34
35    return train_data

```

(二)bert

- bert模型是一种模型迁移，将预训练模型和下游任务模型结合在一起，核心目的就是：是把下游具体NLP任务的活逐渐移到预训练产生词向量上。BERT真真正正意义上同时考虑了上下文。
- 用bert搭建一个实体抽取的模型, 使用Bert用于序列分类的(BertEncoder+Fc+CrossEntropy)模型
bert-base-uncased, 加上MASK-Attention一起送进模型

```

1  def train(net,dataset,num_epochs, learning_rate, batch_size):
2      net.train()
3      optimizer = optim.SGD(net.parameters(),
4 lr=learning_rate,weight_decay=0)
5      train_iter = torch.utils.data.DataLoader(dataset, batch_size,
6 shuffle=True)
7      pre=0
8
9      for epoch in range(num_epochs):
10         correct = 0
11         total=0
12         iter = 0
13         for text,mask, y in train_iter:
14             iter += 1
15             optimizer.zero_grad()
16             if text.size(0)!=batch_size:
17                 break
18             text=text.reshape(batch_size,-1)
19             mask = mask.reshape(batch_size, -1)
20             if USE_CUDA:

```

```

19         text=text.cuda()
20         mask=mask.cuda()
21         y = y.cuda()
22         output = net(text, attention_mask=mask, labels=y)
23         loss = output['loss']
24         logits = output['logits']
25         loss.backward()
26         optimizer.step()
27         _, predicted = torch.max(logits.data, 1)
28
29         total += text.size(0)
30         correct += predicted.data.eq(y.data).cpu().sum()
31         loss=loss.detach().cpu()
32         print("epoch ", str(epoch), " loss: ",
33               loss.mean().numpy().tolist(), "right", correct.cpu().numpy().tolist(),
34               "total", total, "Acc:", correct.cpu().numpy().tolist()/total)
35         acc = eval(model, dev_dataset, 32)
36         if acc > pre:
37             pre = acc
38             torch.save(model, "../model/" + str(acc)+'.pth')
39     return

```

- 约20个epoch收敛

```

epoch 17 loss: 0.04723089933395386 right 11494 total 11648 Acc: 0.9867788461538461
Eval Result: right 707 total 768 Acc: 0.9205729166666666
epoch 18 loss: 0.052388548851013184 right 11547 total 11648 Acc: 0.9913289835164835
Eval Result: right 711 total 768 Acc: 0.92578125
epoch 19 loss: 0.02501794695854187 right 11571 total 11648 Acc: 0.9933894230769231
Eval Result: right 711 total 768 Acc: 0.92578125
epoch 20 loss: 0.028522953391075134 right 11599 total 11648 Acc: 0.9957932692307693
Eval Result: right 710 total 768 Acc: 0.9244791666666666
epoch 21 loss: 0.16468557715415955 right 11602 total 11648 Acc: 0.9960508241758241
Eval Result: right 713 total 768 Acc: 0.9283854166666666
epoch 22 loss: 0.02653355896472931 right 11607 total 11648 Acc: 0.9964800824175825
Eval Result: right 712 total 768 Acc: 0.9270833333333334

```

(三)RNN注意力模型

- 模型定义

```

1 def RNN_model(input_layer, num_class): # RNN model
2     def smoothing_attention(x):
3         e = K.sigmoid(x)
4         s = K.sum(e, axis=-1, keepdims=True)
5         return e / s
6     reg = 0.0001
7     dropout = 0.5
8     hidden_dim = 1024
9     vector = Bidirectional(CuDNNGRU(hidden_dim, return_sequences=False,
10 kernel_regularizer=keras.regularizers.l2(reg)))(input_layer)
11     lstm = Bidirectional(CuDNNGRU(hidden_dim, return_sequences=True,
12 kernel_regularizer=keras.regularizers.l2(reg)))(input_layer)
13     ee = Dot(axis=-1, normalize=True)([vector, lstm])
14     weights = Lambda(smoothing_attention)(ee)
15     weights = RepeatVector(2*hidden_dim)(weights)
16     weights = Permute((2, 1))(weights)
17     output = Multiply()(weights, lstm)

```

```

16     output = Lambda(lambda x: K.sum(x, axis=1))(output)
17     output = Dense(512)(output)
18     output = BatchNormalization()(output)
19     output = Activation("relu")(output)
20     output = Dense(256)(output)
21     output = BatchNormalization()(output)
22     output = Activation("relu")(output)
23     output = Dropout(dropout)(output)
24     output = Dense(num_class, activation='softmax')(output)
25     model = Model(sequence_input, output)
26     print(model.summary())
27     return model

```

Layer (type)	Output Shape	Param #	Connected to
input_1 (InputLayer)	(None, 100)	0	
embedding_1 (Embedding)	(None, 100, 300)	5868900	input_1[0][0]
bidirectional_1 (Bidirectional)	(None, 2048)	8146944	embedding_1[0][0]
bidirectional_2 (Bidirectional)	(None, 100, 2048)	8146944	embedding_1[0][0]
dot_1 (Dot)	(None, 100)	0	bidirectional_1[0][0] bidirectional_2[0][0]
lambda_1 (Lambda)	(None, 100)	0	dot_1[0][0]
repeat_vector_1 (RepeatVector)	(None, 2048, 100)	0	lambda_1[0][0]
permute_1 (Permute)	(None, 100, 2048)	0	repeat_vector_1[0][0]
multiply_1 (Multiply)	(None, 100, 2048)	0	permute_1[0][0] bidirectional_2[0][0]
lambda_2 (Lambda)	(None, 2048)	0	multiply_1[0][0]
dense_1 (Dense)	(None, 512)	1049088	lambda_2[0][0]
batch_normalization_1 (BatchNor	(None, 512)	2048	dense_1[0][0]
activation_1 (Activation)	(None, 512)	0	batch_normalization_1[0][0]
dense_2 (Dense)	(None, 256)	131328	activation_1[0][0]
batch_normalization_2 (BatchNor	(None, 256)	1024	dense_2[0][0]

- 训练

- 使用斯坦福Common Crawl (42B tokens, 1.9M vocab, uncased, 300d vectors, 1.75 GB)作为词向量预训练模型
- 仿照bert的形式，对实体前后做标注 xxxxxxxxxe1xxxxxxxxx 与 ssssssssse1ssssssss

```

1     embeddings_index = {}
2     with open(os.path.join(GLOVE_DIR, 'glove.42B.300d.txt'),
encoding='utf-8') as f: # read pre-trained word embedding
3         for line in f:
4             values = line.split()
5             word = values[0]
6             coefs = np.asarray(values[1:], dtype='float32')
7             embeddings_index[word] = coefs
8     x_train = []
9     y_train = []
10    label_to_y = dict()
11    with open(os.path.join(GLOVE_DIR, "train.txt")) as f:
12        for idx, l in enumerate(f):
13            l = l.strip()
14            if idx % 4 == 0:

```

```

15         ID, sentence = l.split("\t")
16         sentence = sentence[1:-1]
17         sentence = sentence.replace('<e1>',
'xxxxxxxxxe1xxxxxxxxx ')
18         sentence = sentence.replace('<e2>',
'xxxxxxxxxe2xxxxxxxxx ')
19         sentence = sentence.replace('</e1>', '
ssssssssse1sssssssss')
20         sentence = sentence.replace('</e2>', '
ssssssssse2sssssssss')
21         X_train.append(sentence)
22         elif idx % 4 == 1:
23             label = 1
24             if label not in label_to_y:
25                 label_to_y[label] = len(label_to_y)
26                 Y_train.append(label_to_y[label])
27             else:
28                 pass
29         y_to_label = {j:i for i, j in label_to_y.items()}
30         Y_train = np.array(Y_train, dtype=int)
31         num_class = max(Y_train) + 1
32         Y_train = to_categorical(Y_train)
33         tokenizer = Tokenizer(oov_token="UNK")
34         tokenizer.fit_on_texts(X_train)
35         sequences = tokenizer.texts_to_sequences(X_train)
36         X_train = pad_sequences(sequences, maxlen=MAX_SEQUENCE_LENGTH)
37         if VALIDATION_SPLIT > 0: # generate validation set
38             indices = np.arange(len(Y_train))
39             np.random.shuffle(indices)
40             val_index = int(VALIDATION_SPLIT * len(Y_train))
41             X_val = X_train[indices[:val_index]]
42             Y_val = Y_train[indices[:val_index]]
43             X_train = X_train[indices[val_index:]]
44             Y_train = Y_train[indices[val_index:]]
45         X_test = []
46         ID_test = []
47         with open(os.path.join(GLOVE_DIR, "test.txt")) as f:
48             for l in f:
49                 ID, sentence = l.strip().split("\t")
50                 sentence = sentence[1:-1]
51                 sentence = sentence.replace('<e1>', 'xxxxxxxxxe1xxxxxxxxx
')
52                 sentence = sentence.replace('<e2>', 'xxxxxxxxxe2xxxxxxxxx
')
53                 sentence = sentence.replace('</e1>', '
ssssssssse1sssssssss')
54                 sentence = sentence.replace('</e2>', '
ssssssssse2sssssssss')
55                 ID_test.append(ID)
56                 X_test.append(sentence)
57                 sequences = tokenizer.texts_to_sequences(X_test)
58                 X_test = pad_sequences(sequences, maxlen=MAX_SEQUENCE_LENGTH)
59                 word_index = tokenizer.word_index # word dictionary <word, index>
60                 num_words = len(word_index) + 1
61                 embedding_matrix = np.zeros((num_words, EMBEDDING_DIM)) # create
word embedding matrix
62                 for word, i in word_index.items():
63                     embedding_vector = embeddings_index.get(word)

```

```

64         if embedding_vector is not None:
65             embedding_matrix[i] = embedding_vector
66         embedding_layer = Embedding(num_words, # initial word embedding
weights
67                                     EMBEDDING_DIM,
68                                     weights=[embedding_matrix],
69                                     input_length=MAX_SEQUENCE_LENGTH,
70                                     trainable=True)
71         sequence_input = Input(shape=(MAX_SEQUENCE_LENGTH, ),
dtype='int32') # input layer
72         embedded_sequences = embedding_layer(sequence_input) # word
embedding
73         model = RNN_model(embedded_sequences, num_class)
74         model.compile(loss='categorical_crossentropy',
75                       optimizer=keras.optimizers.adam(lr= 0.001,
amsgrad=True, clipvalue=15),
76                       metrics=['accuracy'])
77         early_stop = EarlyStopping(monitor='val_loss' if VALIDATION_SPLIT >
0 else "loss", patience=15, mode='min')
78         model.fit(X_train, Y_train,
79                 batch_size=128,
80                 epochs=50,
81                 callbacks=[early_stop],
82                 validation_data=(X_val, Y_val) if VALIDATION_SPLIT > 0 else
None)
83         Y_pre = model.predict(X_test)
84         Y_pre = np.argmax(Y_pre, axis=1)
85         Y_pre = [y_to_label[i] for i in Y_pre]
86         with open("predict.txt", 'w') as f:
87             for ID, label in zip(ID_test, Y_pre):
88                 f.write(ID + "\t" + label + "\n")
89

```

- 该方法收敛极快，约10个epoch就收敛了

```

Epoch 3/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.8462 - acc: 0.7564 - val_loss: 1.1482 - val_acc: 0.6630
Epoch 4/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.5611 - acc: 0.8450 - val_loss: 1.9637 - val_acc: 0.5104
Epoch 5/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.3579 - acc: 0.9193 - val_loss: 1.0919 - val_acc: 0.6861
Epoch 6/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.2266 - acc: 0.9576 - val_loss: 1.2955 - val_acc: 0.6740
Epoch 7/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.1909 - acc: 0.9687 - val_loss: 1.1721 - val_acc: 0.7135
Epoch 8/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.1329 - acc: 0.9842 - val_loss: 1.1935 - val_acc: 0.7179
Epoch 9/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.1004 - acc: 0.9918 - val_loss: 1.2665 - val_acc: 0.7256
Epoch 10/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.0794 - acc: 0.9959 - val_loss: 1.3510 - val_acc: 0.7212
Epoch 11/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.1161 - acc: 0.9851 - val_loss: 1.3609 - val_acc: 0.6959
Epoch 12/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.0759 - acc: 0.9965 - val_loss: 1.2854 - val_acc: 0.7420
Epoch 13/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.0670 - acc: 0.9974 - val_loss: 1.3495 - val_acc: 0.7014
Epoch 14/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.0571 - acc: 0.9989 - val_loss: 1.2409 - val_acc: 0.7398
Epoch 15/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.0609 - acc: 0.9973 - val_loss: 1.3256 - val_acc: 0.7387
Epoch 16/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.1222 - acc: 0.9816 - val_loss: 1.8178 - val_acc: 0.6103
Epoch 17/50
8206/8206 [=====] - 15s 2ms/step - loss: 0.0939 - acc: 0.9896 - val_loss: 1.6436 - val_acc: 0.6926
Epoch 18/50

```

三、优化

(一)bert

- 优化器如 `sgd`, `adam`, `lr_scheduler` 等的使用

- 发现使用adam无法收敛，初步怀疑是自己的参数设置的不好
- `lr_scheduler` 可以有有效的减少 `sgd` 的反向跳动
- 最终模型准确率在训练集可以达到99%，在验证集可以达到90%，但在在线系统提交后(已经使用优化器)却只有58%，初步猜测数据分布不均衡造成过拟合。

经过统计，发现训练集类别分布不均匀，`other` 类别数目过多

```
Distribution: [1890. 1355. 1415. 1344. 984. 1057. 1088. 1050. 798. 696.]
```

这同时可以从模型的预测发现，模型更倾向于预测0类别

```
tensor([1, 6, 0, 3, 1, 6, 8, 2, 1, 8, 4, 0, 0, 0, 2, 2, 2, 3, 7, 4, 4, 1, 4, 6,
        3, 3, 7, 4, 4, 9, 7, 1], device='cuda:0')
tensor([2, 0, 0, 0, 0, 6, 2, 2, 0, 2, 2, 0, 0, 0, 6, 2, 0, 0, 0, 2, 0, 2, 0, 2,
        0, 0, 0, 0, 0, 2, 2, 2], device='cuda:0')
```

上一行显示为label, 下一行为预测值

- 对交叉熵按分布加权

```
1 ce_loss = F.cross_entropy(output, target, weight=self.weight,
    reduction='none')
```

其中 `weight` 第一次取值为 `[1, 1.39824732, 1.3019039, 1.36501901, 1.90703851, 1.79724656, 1.67953216, 1.82002535, 2.36963696, 2.65925926]`

经过这样处理后模型的预测就比较均匀，最后正确率为 52%(我试了好几种分布，均没有原来的高QWQ)

- 使用**Focal Loss**

Focal Loss 是用于解决 `hard sample` 的一种loss, 它也是对 `loss` 加权, 按照 `pred` 之间的比例来进行加权

```
1 class FocalLoss(nn.Module):
2     def __init__(self, gamma=0, weight=None, reduction='mean'):
3         super(FocalLoss, self).__init__()
4         self.gamma = gamma
5         self.weight = weight
6         self.reduction = reduction
7
8     def forward(self, output, target):
9         # convert output to pseudo probability
10        out_target = torch.stack([output[i, t] for i, t in
    enumerate(target)])
11        probs = torch.sigmoid(out_target)
12        focal_weight = torch.pow(1 - probs, self.gamma)
13
14        # add focal weight to cross entropy
15        ce_loss = F.cross_entropy(output, target,
    weight=self.weight, reduction='none')
16        focal_loss = focal_weight * ce_loss
17
18        if self.reduction == 'mean':
19            focal_loss = (focal_loss / focal_weight.sum()).sum()
20        elif self.reduction == 'sum':
21            focal_loss = focal_loss.sum()
22
23        return focal_loss
```

使用Focal Loss的训练部分仅需要更改loss获得形式即可

```
1 loss_fn = FocalLoss()
2 loss_fn = FocalLoss()
3 loss = loss_fn(log, y)
4 loss.backward()
```

(二)RNN模型

该部分见算法Part 3

四、结果

ACC

(一)bert

- bert+att : 0.58
- bert+att+weighted loss: 0.52
- bert+att+focal loss: 0.59
- bert+att+focal: 0.61

Filename	ACC-Relation	ACC-NER
张永停-PB17111585-7-1.txt	0.613125	0.0

(二)RNN+ATT

- 0.62

Filename	ACC-Relation	ACC-NER
张永停-PB17111585-10.txt	0.6225	0.0

无论哪种方法，均在测试集表现不佳，模型过拟合严重。

当使用bert时，验证集准确率可以达到0.9, 但RNN+ATT只有70%附近。私以为，bert的学习能力要比RNN强(当然这个几乎公认了)。

其中bert在验证集的表现最令我疑惑，为什么验证集这么高但测试集却只有0.6?