

服务器

ssh -p 22 user@10.201.156.73

password: 321kgCODE

路径

/home/user/hotpot/hotpot-master

环境

```
source activate hotpot
```

下载数据

```
./download.sh
```

<input type="checkbox"/>	hotpot_train_v1.1.json	566,426,227	json-文件
<input type="checkbox"/>	hotpot_dev_fullwiki_v1.json	47,454,698	json-文件
<input type="checkbox"/>	hotpot_dev_distractor_v1.json	46,320,117	json-文件

glove

	glove.840B.300d.zip	2,176,768,927
---	---------------------	---------------

```
unzip glove.840B.300d.zip
```

预处理

```
python main.py --mode prepro --data_file hotpot_train_v1.1.json --para_limit  
2250 --data_split train
```

生成

word2idx.json	6,662,434	json-文件
idx2word.json	7,453,340	json-文件
idx2char.json	110,590	json-文件
char2idx.json	96,662	json-文件
char_emb.json	1,218,947	json-文件
word_emb.json	1,018,428,687	json-文件
train_eval.json	1,872,387,536	json-文件
train_record.pkl	28,621,801,707	pkl-文件

```
python main.py --mode prepro --data_file hotpot_dev_distractor_v1.json --
para_limit 2250 --data_split dev
```

生成

dev_eval.json	154,878,531	json-文件
dev_record.pkl	2,350,918,447	pkl-文件

```
python main.py --mode prepro --data_file hotpot_dev_fullwiki_v1.json --
data_split dev --fullwiki --para_limit 2250
```

生成

fullwiki.dev_eval.json	158,904,839	json-文件
fullwiki.dev_record.pkl	2,343,996,887	pkl-文件

预处理流程

1. 通过 `process_file` 得到 `example` 和 `eval_example`

```
process_file(config.data_file, config) #处理一个json文件
_process_article() #处理文件中的一个item
_process(sent, is_sup_fact, is_title=False) # 对每句话进行处理
```

对一条数据，解析成以下格式

```
example = {'context_tokens': context_tokens,
           'context_chars': context_chars,
           'ques_tokens': ques_tokens,
```

```

        'ques_chars': ques_chars,
        'y1s': [best_indices[0]],
        'y2s': [best_indices[1]],
        'id': article['_id'],
        'start_end_facts': start_end_facts # 由多个(start_token_id,
end_token_id, is_sup_fact=True/False)组成, 每句话的起始token和终止token在
context_tokens中的位置, 以及这句话是否是支撑事实。
    }

eval_example = {'context': text_context, # 合并成一个str的context
                'spans': flat_offsets, # context中每个token的span
                'answer': [answer],
                'id': article['_id'],
                'sent2title_ids': sent2title_ids # list, 由多个[para_title,
idx]组成, 表示context中每句话对应的标题以及在标题下的idx, title本身对应的idx为-1;
                } # 汇总后写入{}_eval.json文件中

```

best_indices:

1. Yes, [-1, -1]
2. No, [-2, -2]
3. Yes or no 之外的答案, 且不在上下文中 (0, 1)
4. Yes or no 之外的答案, 且在上下文中, 对应的位置

```

answer = article['answer'].strip()
if answer.lower() == 'yes':
    best_indices = [-1, -1]
elif answer.lower() == 'no':
    best_indices = [-2, -2]
else:
    if article['answer'].strip() not in ''.join(text_context):
        # in the fullwiki setting, the answer might not have been
retrieved

        # use (0, 1) so that we can proceed
        best_indices = (0, 1)
    else:
        _, best_indices, _ = fix_span(text_context, offsets,
article['answer'])
        answer_span = []
        for idx, span in enumerate(flat_offsets):
            if not (best_indices[1] <= span[0] or best_indices[0] >=
span[1]):

                answer_span.append(idx)
        best_indices = (answer_span[0], answer_span[-1])

```

2. 得到字词嵌入矩阵

若word2idx_file文件存在，载入此文件作为word2idx_dict，否则通过get_embedding得到word_emb_mat, word2idx_dict, idx2word_dict；

字嵌入同理；

3. build_features构建特征

Para_limit用来限制context_tokens长度，在训练阶段只选择len(context_tokens) > Para_limit的数据进行训练

```
datapoints.append({'context_idx': torch.from_numpy(context_idx), # 一维张量，句子中每个token的idx，长度为para_limit
                  'context_char_idx': torch.from_numpy(context_char_idx), # 二维张量，[para_limit, charlimit]，字id矩阵，[i, j]第i个token中第j个字符的id
                  'ques_idx': torch.from_numpy(ques_idx), # 类似context_idx
                  'ques_char_idx': torch.from_numpy(ques_char_idx), # 类似context_char_idx
                  'y1': y1, # 见`best_indices`
                  'y2': y2,
                  'id': example['id'],
                  'start_end_facts': example['start_end_facts']}) # 由多个(start_token_id, end_token_id, is_sup_fact=True/False)组成，每句话的起始token和终止token在context_tokens中的位置，以及这句话是否是支撑事实。
```

```
torch.save(datapoints, out_file)
```

存入{}.record.pkl文件中