

tags: ADL

HW2

```
python3 run_qa.py --lr=3e-5 --doc_stride=200 --  
batch_size=2 --num_epoch=1 --  
model_name=hfl/chinese-macbert-large --device=cuda  
python3 qatest.py --lr=3e-5 --doc_stride=200 --  
qa_path=./ckpt/macbert.ckpt --  
context_file=./data/context.json --  
test_file=./data/test.json --pred_file=./pred.csv  
--model_name=hfl/chinese-macbert-large --  
device=cuda
```

Q1.

1.

I use a bert tokenizer, the bert tokenizer is pretrained on massive dataset. It separates words into wordpieces, and transform them into tokens. My macbert tokenizer uses the traditional Chinese Word Segmentation (CWS) tool to split the text into several words. In this way, we could adopt whole word masking(WWM) in Chinese to mask the word instead of individual Chinese characters

2.

- a. I use the `char_to_token` function from tokenizer to convert answer position to token position
- b. I first check top-4 probability of start and end position, to see if they are legal indexes and select the combination with largest prob, then I use the `offset_mapping` in tokenizer to know the correspondence of token position and char position.

Q2.

1.

- a. I use macbert_large for bert, the config is

```

1  {
2    "architectures": [
3      "BertForMaskedLM"
4    ],
5    "attention_probs_dropout_prob": 0.1,
6    "directionality": "bidi",
7    "gradient_checkpointing": false,
8    "hidden_act": "gelu",
9    "hidden_dropout_prob": 0.1,
10   "hidden_size": 1024,
11   "initializer_range": 0.02,
12   "intermediate_size": 4096,
13   "layer_norm_eps": 1e-12,
14   "max_position_embeddings": 512,
15   "model_type": "bert",
16   "num_attention_heads": 16,
17   "num_hidden_layers": 24,
18   "pad_token_id": 0,
19   "pooler_fc_size": 768,
20   "pooler_num_attention_heads": 12,
21   "pooler_num_fc_layers": 3,
22   "pooler_size_per_head": 128,
23   "pooler_type": "first_token_transform",
24   "type_vocab_size": 2,
25   "vocab_size": 21128
26 }

```

b. EM

- public: 0.78571
- private: 0.80126

c. cross entropy loss

d. AdamW, lr=3e-5, batch_size=2

2.

a. I used Langboat/mengzi-bert-base

b. EM

- public: 0.78119
- private: 0.77777

c, d. This model has smaller

hidden_size, intermediate_size, num_hidden_layers

, and the loss drops slower than macbert_large. And the

pretrain masking are different, macbert uses similar words

for the masking purpose. A similar word is obtained by using

Synonymstoolkit, which is based on word2vec. Mengzi uses

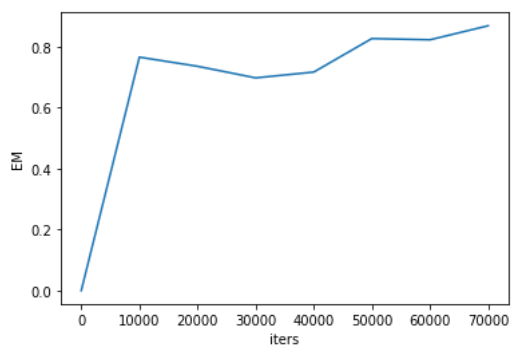
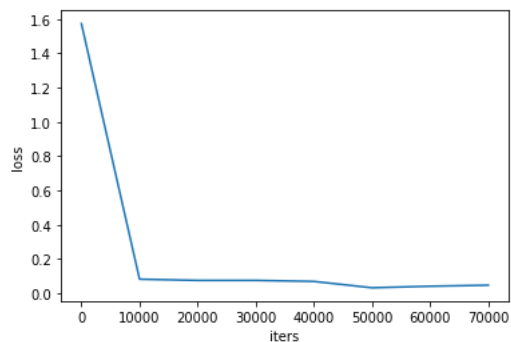
LAMB optimizer and mixed-batch training in pretrain stage.

```

1  {
2    "architectures": [
3      "BertForMaskedLM"
4    ],
5    "attention_probs_dropout_prob": 0.1,
6    "bos_token_id": 0,
7    "directionality": "bidi",
8    "eos_token_id": 2,
9    "gradient_checkpointing": false,
10   "hidden_act": "gelu",
11   "hidden_dropout_prob": 0.1,
12   "hidden_size": 768,
13   "initializer_range": 0.02,
14   "intermediate_size": 3072,
15   "layer_norm_eps": 1e-12,
16   "max_position_embeddings": 512,
17   "model_type": "bert",
18   "num_attention_heads": 12,
19   "num_hidden_layers": 12,
20   "output_past": true,
21   "pad_token_id": 1,
22   "pooler_fc_size": 768,
23   "pooler_num_attention_heads": 12,
24   "pooler_num_fc_layers": 3,
25   "pooler_size_per_head": 128,
26   "pooler_type": "first_token_transform",
27   "position_embedding_type": "absolute",
28   "torch_dtype": "float16",
29   "transformers_version": "4.9.2",
30   "type_vocab_size": 2,
31   "use_cache": true,
32   "vocab_size": 21128
33 }

```

Q3.



Q4

- config:

```

1  {
2      "architectures": [
3          "BertForMaskedLM"
4      ],
5      "attention_probs_dropout_prob": 0.1,
6      "directionality": "bidi",
7      "gradient_checkpointing": false,
8      "hidden_act": "gelu",
9      "hidden_dropout_prob": 0.1,
10     "hidden_size": 256,
11     "initializer_range": 0.02,
12     "intermediate_size": 1024,
13     "layer_norm_eps": 1e-12,
14     "max_position_embeddings": 512,
15     "model_type": "bert",
16     "num_attention_heads": 4,
17     "num_hidden_layers": 6,
18     "pad_token_id": 0,
19     "pooler_fc_size": 256,
20     "pooler_num_attention_heads": 4,
21     "pooler_num_fc_layers": 3,
22     "pooler_size_per_head": 128,
23     "pooler_type": "first_token_transform",
24     "type_vocab_size": 2,
25     "vocab_size": 21128
26 }

```

I train from a config with smaller layers and sizes, and used the same learning rate as the pretrained scene.

```

1  configuration = BertConfig.from_pretrained('./macconf.json')
2  model = BertForQuestionAnswering(configuration).to(device)

```

- EM:0.03

Q5

a.

- Intent: BertForSequenceClassification pretrained on bert-base-uncased
- Slot: BertForTokenClassification pretrained on bert-base-uncased

The model was pretrained with two objectives:

Masked language modeling (MLM): taking a sentence, the model randomly masks 15% of the words in the input then run the entire masked sentence through the model and has to predict the masked words. This is different from traditional recurrent neural networks (RNNs) that usually see the words one after the other, or from autoregressive models like GPT which internally mask the future tokens. It allows the model to learn a bidirectional representation of the sentence.

Next sentence prediction (NSP): the models concatenates two masked sentences as inputs during pretraining.

Sometimes they correspond to sentences that were next to

each other in the original text, sometimes not. The model then has to predict if the two sentences were following each other or not.

b.

- intent: 0.95288

predict_0.csv 11 days ago by b08902047	0.95288	0.95644	
---	---------	---------	---

- slot: 0.82636

spred.csv 10 days ago by b08902047	0.82636	0.83592	
---	---------	---------	---

c. cross entropy loss

d. optimizer = adam, lr = 3e-5, batch_size = 32