

## ADL HW2 Report

### Q1:Data processing

#### 1. Tokenizer:

- The algorithm used in BertTokenizer is “WordPiece”. It is the subword tokenization algorithm and is very similar to BPE. WordPiece first initializes the vocabulary to include every character present in the training data and progressively learns a given number of merge rules. In contrast to BPE, WordPiece does not choose the most frequent symbol pair, but the one that maximizes the likelihood of the training data once added to the vocabulary.

#### 2. Answer span:

- a. When tokenize the question and context, I set “return\_offsets\_mapping” to TRUE which returns (char\_start, char\_end) for each token. Then, find start token index and end token index of the current span in the context by “sequence\_ids” (indicate what is the context, question, and padding). Finally, use while-loop and offsets to find the start and end positions.
  - While token\_start\_index < len(offsets) and offsets[token\_start\_index][0] <= start\_char:  
token\_start\_index += 1. (Final start= token\_start\_index-1)
  - while offsets[token\_end\_index][1] >= end\_char:  
token\_end\_index -= 1. (Final end= token\_end\_index+1)
- b. First, use argsort to find start index and end index from output.start\_logits and output.end\_logits. Then, determine the final start/end position by using offsets mapping.  
(start\_char = offset\_mapping[start\_index][0],  
end\_char = offset\_mapping[end\_index][1])

### Q2:Modeling with BERTs and their variants

#### 1. Describe

- a. model: “bert-base-chinese”
- b. performance: EM:0.796, F1:0.896 (For QA)

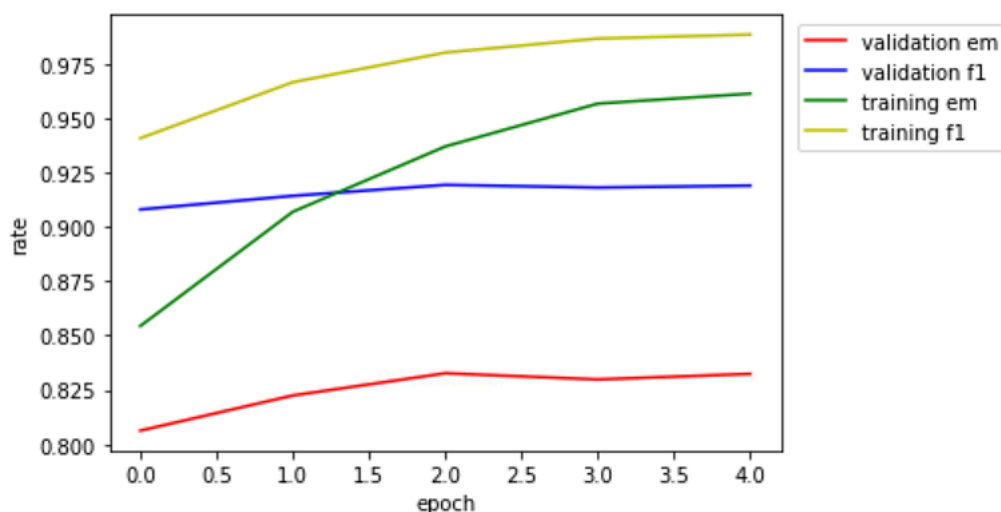
- c. loss function: CrossEntropyLoss
- d. optimization algorithm: AdamW, learning rate:  $1e-5$ , batch size: 16.

## 2. Try another type of pretrained model and describe

- a. model: "hfl/chinese-macbert-base"
- b. performance: EM:0.833, F1:0.919 (For QA)
- c. difference between pretrained model: MacBERT is an improved BERT with novel MLM as correction pre-training task, which mitigates the discrepancy of pre-training and fine-tuning. Instead of masking with [MASK] token, they propose to use similar words which is obtained by using Synonyms toolkit based on word2vec similarity calculations for the masking purpose. If an N-gram is selected to mask, it will find similar words individually. In rare cases, when there is no similar word, they will degrade to use random word replacement. Except for the new pre-training task, they also incorporate the Whole Word Masking (WWM), N-gram masking, and Sentence-Order Prediction (SOP). Note that there is no differences in the main neural architecture from original BERT.

## Q3:Curves

### 1. Plot (learning curve of EM and F1)



#### Q4:Pretrained vs Not Pretrained

- a. The configuration of the model and how do I train this model: I load the configuration from "bert-base-chinese" but change the hidden\_size to 512, the number of attention head to 8, and the number of hidden layers to 6. I use the same training script as pretrained BERT model for this problem.
- b. the performance of this model v.s. BERT:
- Train from scratch: The optimizer I used is AdamW, learning rate is  $1e-4$ , and the number of epoch is 30. The picture below is the final result. The result shows that we can get good performance on training set but poor performance on validation set. (The first line is training loss and accuracy, the second line is validation loss and accuracy.) (The upper picture below). However, we can get good performance on both training set and validation set when using BERT. (The lower picture below)

```
loss: 136.17594998884306 , acc: 0.908725931995084
loss: 2097.9780416488647 , acc: 0.048496880317640385
```

```
■ 6103/? [12:31<00:00, 8.41it/s]
```

```
epoch: 10 , loss: 49.65869615044676 , acc: 0.9612044244162229
```

```
■ 882/? [00:38<00:00, 24.26it/s]
```

```
epoch: 10 , loss: 930.0942116413862 , acc: 0.7702779353374929
```

#### Q5:HW1 with BERTs

- a. model:
- Intent classification: I use the pretrained model "bert-base-uncased". And I utilize "BertForSequenceClassification" in huggingface to do this case.
  - Slot tagging: I use the pretrained model "bert-base-uncased". And I utilize "BertForTokenClassification" in huggingface to do this case.
- b. performance:

- Intent classification: In HW1, I used GRU to train the model. The public score on Kaggle is 0.912, the private score is 0.910. The picture below is the model trained on bert. The time consuming is lower than using RNN model and we can get accuracy of 0.96 on evaluation which is higher than using RNN model.(The second line of "epoch:" is validation loss and accuracy.)

```

100% ██████████ 5/5 [02:39<00:00, 31.93s/it]
100% ██████████ 469/469 [00:30<00:00, 15.58it/s]
epoch: 1 , loss: 1521.5530570745468 , acc: 0.5656666666666667
epoch: 1 , loss: 154.48129081726074 , acc: 0.891
100% ██████████ 469/469 [00:30<00:00, 15.58it/s]
epoch: 2 , loss: 457.7452366948128 , acc: 0.9444
epoch: 2 , loss: 46.163687855005264 , acc: 0.9483333333333334
100% ██████████ 469/469 [00:30<00:00, 15.49it/s]
epoch: 3 , loss: 135.2215090841055 , acc: 0.9877333333333334
epoch: 3 , loss: 23.895966589450836 , acc: 0.9613333333333334
100% ██████████ 469/469 [00:30<00:00, 15.56it/s]
epoch: 4 , loss: 58.983032681047916 , acc: 0.9934
epoch: 4 , loss: 18.864845853298903 , acc: 0.9646666666666667
100% ██████████ 469/469 [00:30<00:00, 15.45it/s]
epoch: 5 , loss: 30.92214596271515 , acc: 0.9969333333333333
epoch: 5 , loss: 17.481978215277195 , acc: 0.9606666666666667

```

- Slot tagging: In HW1, I used GRU to train the model. The public score on Kaggle is 0.783, the private score is 0.786. The picture below is the model trained on bert. The time consuming is lower than using RNN model and we can get accuracy of 0.837 on evaluation which is higher than using RNN model.(The second line of "epoch:" is validation loss and accuracy.)

```

100% ██████████ 15/15 [05:31<00:00, 22.13s/it]
100% ██████████ 227/227 [00:20<00:00, 11.16it/s]
epoch: 1 , loss: 47.79889163002372 , acc: 0.697266703478741
epoch: 1 , loss: 3.2492468804121017 , acc: 0.773
100% ██████████ 227/227 [00:20<00:00, 10.67it/s]
epoch: 2 , loss: 18.310740925371647 , acc: 0.8128106018774158
epoch: 2 , loss: 3.029066652059555 , acc: 0.785
100% ██████████ 227/227 [00:20<00:00, 11.15it/s]
epoch: 3 , loss: 13.64531574305147 , acc: 0.8503589177250138
epoch: 3 , loss: 2.6460428535938263 , acc: 0.824
100% ██████████ 227/227 [00:20<00:00, 11.05it/s]
epoch: 4 , loss: 10.497511544264853 , acc: 0.8825234676974048
epoch: 4 , loss: 2.846926588565111 , acc: 0.809
100% ██████████ 227/227 [00:20<00:00, 10.36it/s]
epoch: 5 , loss: 8.140312465373427 , acc: 0.9046107123136389
epoch: 5 , loss: 3.2746324334293604 , acc: 0.822
100% ██████████ 227/227 [00:20<00:00, 11.87it/s]
epoch: 6 , loss: 7.060944152995944 , acc: 0.9182771949199338
epoch: 6 , loss: 3.485572222620249 , acc: 0.807
100% ██████████ 227/227 [00:20<00:00, 11.40it/s]
epoch: 7 , loss: 5.894104221370071 , acc: 0.9302871341800111
epoch: 7 , loss: 3.3981530955061316 , acc: 0.837

```

c. loss function:

- Intent classification: The loss function used in "BertForSequenceClassification" is CrossEntropyLoss.
- Slot tagging: The loss function used in "BertForTokenClassification" is CrossEntropyLoss.

d. Optimization algorithm, learning rate and batch size.

- Intent classification: The optimizer I used is AdamW from huggingface, learning rate is  $3e-5$ , and batch size is 32.
- Slot tagging: The optimizer I used is AdamW from huggingface, learning rate is  $3e-5$ , and batch size is 32.