ECE 601: Homework 10

Wei Xu

Email: xu1639@purdue.edu

Due date: 11:59 PM, Apr. 25, 2024

(Spring 2024)

1 Programming Tasks

1.1 Dataset

Note: Some of the code for this subsection is borrowed from the DLStudio package (https://engineering.purdue.edu/kak/distDLS/).

The code for loading data sets is shown below. The data sets are provided with pkl files, which are all dictionaries, and can be loaded using pickle.load method in the load_dataset function. There are 10,000 samples in total, and the proportion is 7:2:1 for training, evaluation, and testing sets respectively. Specifically, the *_dict.pkl files contain 'id', 'title', 'context', 'question', 'answers', while the *_processed.pkl files contain 'input_ids', 'attention_mask', 'start_positions', and 'end_positions'.

```
# load datasets
1
   def load_dataset(path='./dataset'):
       with open('{}/train_dict.pkl'.format(path), 'rb') as f:
3
4
           train_dict = pickle.load(f)
5
           f.close()
       with open('{}/eval_dict.pkl'.format(path), 'rb') as f:
6
7
           eval_dict = pickle.load(f)
8
           f.close()
       with open('{}/test_dict.pkl'.format(path), 'rb') as f:
9
           test_dict = pickle.load(f)
10
           f.close()
11
       with open('{}/train_data_processed.pkl'.format(path), 'rb') as f:
12
           train_processed = pickle.load(f)
13
           f.close()
14
       with open('{}/eval_data_processed.pkl'.format(path), 'rb') as f:
15
           eval_processed = pickle.load(f)
16
           f.close()
17
       with open('{}/test_data_processed.pkl'.format(path), 'rb') as f:
18
19
           test_processed = pickle.load(f)
20
           f.close()
21
       # print keys
22
       print(train_dict.keys())
23
       print(eval_dict.keys())
       print(test_dict.keys())
24
25
       print(train_processed.keys())
26
       print(eval_processed.keys())
27
       print(test_processed.keys())
       return train_dict, eval_dict, test_dict, train_processed,
28
          \hookrightarrow eval_processed, test_processed
```

1.2 BERT for Q&A

Note: Some of the code for this subsection is borrowed from the DLStudio package (https://engineering.purdue.edu/kak/distDLS/).

The code for fine-tuning a model is shown below. In the train_model function, the model ('bert-base-uncased') is initialized through BertForQuestionAnswering.from_pretrained. The model structure can be inspected by printing model.modules. The training arguments are set through TrainingArguments. In this experiment, the number of epochs is 10, the batch size for training is 8, and the batch size for evaluation is 8. The log is printed every single epoch. The training and evaluation dataset loaders are created through Dataset.from_pandas with the *_processed sets. Then train the model using Trainer. And save the model after finishing training. The training output of the first 5 epochs can be found below the code snippet.

```
# fine tune model
   def train_model(train_processed, eval_processed):
3
       # initialize model
       model_name = 'bert-base-uncased'
4
       model= BertForQuestionAnswering.from_pretrained(model_name)
5
       print(model._modules)
6
7
       # set training arguments
       training_args = TrainingArguments(output_dir='./results', # output
8
          \hookrightarrow directory
9
                                          use_mps_device=False,
10
                                          num_train_epochs=10, # total
                                             → number of training epochs
11
                                          per_device_train_batch_size=8, #
                                             → batch size per device
                                             → during training
12
                                          per_device_eval_batch_size=8, #
                                             → batch size for evaluation
13
                                          weight_decay=0.01, # strength of
                                             14
                                          logging_strategy='epoch',
15
                                          logging_steps=1
16
17
       # create dataset instance
18
       train_dataset = Dataset.from_pandas(pd.DataFrame(train_processed))
19
       eval_dataset = Dataset.from_pandas(pd.DataFrame(eval_processed))
20
       # train model
       trainer = Trainer(model=model, # the instantiated Transformers
21
          → model to be fine-tuned
22
                          args=training_args, # training arguments, defined
                             → above
                         train_dataset=train_dataset, # training dataset
23
24
                          eval_dataset=eval_dataset # evaluation dataset
25
26
       trainer.train()
27
       # save trainer instance
       trainer.save_model('./results/last_model')
```

Outputs:

The training loss keeps decreasing and the learning rate is decaying for more precise results.

1.3 Testing and evaluation metrics

Note: Some of the code for this subsection is borrowed from the DLStudio package (https://engineering.purdue.edu/kak/distDLS/).

The code for testing the model and calculating metrics is shown below. Call test_model to execute this process. The testing set is loaded through Dataset.from_pandas. The fine-tuned model is loaded through BertForQuestionAnswering.from_pretrained from the saved model. The to-kenizer is created through BertTokenizer.from_pretrained. The prediction can be obtained by trainer.predict.

It is noticed that the predictions are all lowercase and in the subword-level tokenization. The preprocessing is needed. First, convert the prediction from subword-level tokenization to word-level tokenization by subword2word. This function simply removes ##. Secondly, normalize the ground truth by normalize_string. This function converts all letters to lowercase ASCII characters and splits the words and symbols to match the prediction format. The degree symbol of in normalize_string can not display well in the code snippet, so ^{\circ} is used to represent it.

The Exact Match and F1-score metrics are calculated by compute_exact_match and f1_score respectively. The inputs are the prediction and ground truth after preprocessing.

15 samples and the (individual and overall) metrics can be found below the code snippet.

```
# convert Unicode to ASCII
1
   def unicode2ascii(s):
2
3
       return ''.join(
4
           c for c in unicodedata.normalize('NFD', s)
5
           if unicodedata.category(c) != 'Mn'
6
       )
7
8
   # normalize sentence
9
   def normalize_string(s):
       s = unicode2ascii(s.lower().strip())
10
       s = re.sub(r"([,.!?%$\-\',\"\"\(\)^{\circ}])", r" \1 ", s)
11
       s = re.sub(r"[^a-zA-Z0-9,.!?%$\-\','\"\"\(\)^{\circ}]+", r" ", s)
12
       s = re.sub(r'^\s+|\s+\$', '', s)
13
14
       s = re.sub(r" ^{\circ} (circ) ", r"^{\circ} (circ) ", s)
15
       return s
16
   # convert subword-level token to word-level token
17
18 def subword2word(s):
     s = re.sub(r" ##", r"", s)
19
```

```
20
      return s
21
  # compute Exact Match
22
23 def compute_exact_match(prediction, truth):
24
       return int(prediction == truth)
25
26
   # compute F1-score
  def f1_score(prediction, truth):
27
       pred_tokens = prediction.split()
28
29
       truth_tokens = truth.split()
       # if either the prediction or the truth is no-answer then F1 = 1 if
30
          \hookrightarrow they agree, 0 otherwise
       if len(pred_tokens) == 0 or len(truth_tokens) == 0:
31
32
           return int(pred_tokens == truth_tokens)
33
       else:
34
           common_tokens = set(pred_tokens) & set(truth_tokens)
35
           # if there are no common tokens then F1 = 0
           if len(common tokens) == 0:
36
               return 0
37
38
           else:
               prec = len(common_tokens) / len(pred_tokens)
39
               rec = len(common_tokens) / len(truth_tokens)
40
               return 2 * (prec * rec) / (prec + rec)
41
42
43
   # test model
   def test_model(test_dict, test_processed):
44
45
       # create dataset instance
46
       test_dataset = Dataset.from_pandas(pd.DataFrame(test_processed))
47
       # load trained model
       model = BertForQuestionAnswering.from_pretrained('./results/
48
          → last_model')
       trainer = Trainer(model=model)
49
50
       # Initialize the tokenizer
       tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
51
       # test trained model
52
       x = trainer.predict(test_dataset)
53
54
       # retrieve offset mapping of prediction
55
       start_pos, end_pos = x.predictions
       start_pos = np.argmax(start_pos, axis=1)
56
       end_pos = np.argmax(end_pos, axis=1)
57
       # create lists to save metrics
58
       EM_list = []
59
60
       F1_list = []
61
       # print predictions and calculate metrics
       for k, (i, j) in enumerate(zip(start_pos, end_pos)):
62
63
           # convert indices to tokens
           tokens = tokenizer.convert_ids_to_tokens(test_processed[')
64
               → input_ids'][k])
65
           # preprocess prediction and ground truth
           prediction = subword2word(', '.join(tokens[i:j+1]))
66
           truth = normalize_string(test_dict['answers'][k]['text'][0])
67
           print('Question:', test_dict['question'][k])
68
           print('Answer:', prediction)
69
           print('Correct Answer:', truth)
70
```

```
71
          em = compute_exact_match(prediction, truth)
72
          f1 = f1_score(prediction, truth)
          print('Exact Match:', em)
73
          print('F1 Score:', f1)
74
75
          EM_list.append(em)
76
          F1_list.append(f1)
          print('---')
77
78
      # print overall metrics
79
      print('\n')
      print('Exact Match: average: {}; median: {}'.format(statistics.mean
         print('F1 Score: average: {}; median: {}'.format(statistics.mean())
         → F1_list), statistics.median(F1_list)))
```

Outputs:

```
Question: What type of imaging was used to study the relationship between
   → humans and dogs?
Answer: mri
Correct Answer: magnetic resonance imaging
Normalized Correct Answer: magnetic resonance imaging
Exact Match: 0
F1 Score: 0
Question: How much money were possible changes to the Mexico City section
   \hookrightarrow of the film rumored to have saved the production?
Answer: $ 20 million
Correct Answer: $20 million
Normalized Correct Answer: $ 20 million
Exact Match: 1
F1 Score: 1.0
Question: Where did the Chinese government decide that parents who had
   → lost children could go for free treatment?
Answer: fertility clinics
Correct Answer: fertility clinics
Normalized Correct Answer: fertility clinics
Exact Match: 1
F1 Score: 1.0
Question: What is the ultimate goal for Theravadins?
Answer: nibbana
Correct Answer: Nibb na
Normalized Correct Answer: nibbana
Exact Match: 1
F1 Score: 1.0
Question: What university was Lok-Ham Chan a professor at?
Answer: university of washington
Correct Answer: the University of Washington
Normalized Correct Answer: the university of washington
Exact Match: 0
F1 Score: 0.8571428571428571
```

```
Question: Who was the Mongol prince?
Answer: sakya pandita
Correct Answer: Godan
Normalized Correct Answer: godan
Exact Match: 0
F1 Score: 0
Question: How many people were buried in the collapsed schools?
Answer: 1 , 700
Correct Answer: 1,700
Normalized Correct Answer: 1, 700
Exact Match: 1
F1 Score: 1.0
Question: Who mentored contestants in the fourteenth and fifteenth seasons
   \hookrightarrow of American Idol?
Answer: randy jackson
Correct Answer: Scott Borchetta
Normalized Correct Answer: scott borchetta
Exact Match: 0
F1 Score: 0
Question: Thomas Stritch was an editor of which publican from Notre Dame?
Answer: matthew fitzsimons , frederick crosson
Correct Answer: Review of Politics
Normalized Correct Answer: review of politics
Exact Match: 0
F1 Score: 0
Question: Paperback of the Year award from Bestsellers magazine was
   → awarded when?
Answer: 1962
Correct Answer: 1962
Normalized Correct Answer: 1962
Exact Match: 1
F1 Score: 1.0
Question: Gray color is often called what when referring to dogs?
Answer: blue
Correct Answer: blue
Normalized Correct Answer: blue
Exact Match: 1
F1 Score: 1.0
Question: How many persons were still unaccounted for in Yingxiu?
Answer: around 3, 000
Correct Answer: around 9,000
Normalized Correct Answer: around 9, 000
Exact Match: 0
F1 Score: 0.75
Question: How were the semi-finalists divided in season four?
Answer: gender
Correct Answer: by gender
```

```
Normalized Correct Answer: by gender
Exact Match: 0
F1 Score: 0.666666666666666
Question: Who was the host of American Idol in its fourteenth season?
Answer: ryan seacrest
Correct Answer: Ryan Seacrest
Normalized Correct Answer: ryan seacrest
Exact Match: 1
F1 Score: 1.0
Question: In what borough is Godiva based?
Answer: manhattan
Correct Answer: Manhattan
Normalized Correct Answer: manhattan
Exact Match: 1
F1 Score: 1.0
Exact Match: average: 0.572; median: 1.0
F1 Score: average: 0.7179283755599364; median: 1.0
```

Generally, the predictions are precise. There are three cases. (1) Some answers are accurate. For example, the prediction for the question ("How much money were possible changes to the Mexico City section of the film rumored to have saved the production?") is "\$ 20 million", which is the same with the (normalized) correct answer; the prediction for the question ("Where did the Chinese government decide that parents who had lost children could go for free treatment?") is "fertility clinics", which is the same with the (normalized) correct answer. (2) Some other answers are good enough but have bad metrics. For example, the prediction for the question ("What type of imaging was used to study the relationship between humans and dogs?") is "mri", which the abbreviation of the (normalized) correct answer "magnetic resonance imaging" (so in my opinion, this prediction can also be considered as a correct answer); the prediction for the question ("What university was Lok-Ham Chan a professor at?") is "university of washington", which only has one less "the" than the (normalized) correct answer (this can also be considered as a correct answer). (3) Some answers are incorrect. For example, the prediction for the question ("Who was the Mongol prince?") is "sakya pandita", while the correct answer should be "godan"; the prediction for the question ("Who was the Mongol prince?") is "sakya pandita", while the correct answer should be "godan".

The average Exact Match is 0.572, meaning 57.2% predictions are exactly correct. The average F1-score is 0.718, which is a not-bad level. The medians for the metrics are both 1.0, meaning at least a half answers are perfect.

However, considering the case (2), in which the answers are good enough but have bad metrics, the real performance of the model should be better than what the metrics tell.

1.4 Comparison

Note: Some of the code for this subsection is borrowed from the DLStudio package (https://engineering.purdue.edu/kak/distDLS/).

The code for comparison with another fine-tuned model is shown below. The compare_model function loads 'distilbert-base-cased-distilled-squad' model from Hugging Face. Then run it over the testing set and calculate the Exact Match and F1-score metrics.

15 samples and the (individual and overall) metrics can be found below the code snippet.

```
# compare with another open-source fine-tuned model
   def compare_model(test_dict):
3
       # load model
       question_answerer = pipeline('question-answering', model='
4

    distilbert-base-cased-distilled-squad')
       # create lists to save metrics
5
       EM_list = []
6
7
       F1_list = []
       # print predictions and calculate metrics
8
       for i in range(len(test_dict['question'])):
9
10
           result = question_answerer(question=test_dict['question'][i],
              context=test_dict['context'][i])
11
           print('Question:', test_dict['question'][i])
           print('Answer:', result['answer'])
12
           print('Correct:', test_dict['answers'][i]['text'][0])
13
           em = compute_exact_match(result['answer'], test_dict['answers']
14
              → ][i]['text'][0])
15
           f1 = f1_score(result['answer'], test_dict['answers'][i]['text'
              \hookrightarrow ][0])
           print('Exact Match:', em)
16
           print('F1 Score:', f1)
17
           EM_list.append(em)
18
19
           F1_list.append(f1)
20
           print('---')
       # print overall metrics
21
       print('\n')
22
       print('Exact Match: average: {}; median: {}'.format(statistics.mean
23
          24
       print('F1 Score: average: {}; median: {}'.format(statistics.mean(
          → F1_list), statistics.median(F1_list)))
```

Outputs:

```
Question: What type of imaging was used to study the relationship between
   \hookrightarrow humans and dogs?
Answer: magnetic resonance imaging
Correct: magnetic resonance imaging
Exact Match: 1
F1 Score: 1.0
Question: How much money were possible changes to the Mexico City section

→ of the film rumored to have saved the production?

Answer: $20 million
Correct: $20 million
Exact Match: 1
F1 Score: 1.0
Question: Where did the Chinese government decide that parents who had
   → lost children could go for free treatment?
Answer: fertility clinics
Correct: fertility clinics
Exact Match: 1
F1 Score: 1.0
```

```
Question: What is the ultimate goal for Theravadins?
Answer: Nibb na
Correct: Nibb na
Exact Match: 1
F1 Score: 1.0
Question: What university was Lok-Ham Chan a professor at?
Answer: University of Washington
Correct: the University of Washington
Exact Match: 0
F1 Score: 0.8571428571428571
Question: Who was the Mongol prince?
Answer: Godan
Correct: Godan
Exact Match: 1
F1 Score: 1.0
Question: How many people were buried in the collapsed schools?
Answer: 1,700
Correct: 1,700
Exact Match: 1
F1 Score: 1.0
Question: Who mentored contestants in the fourteenth and fifteenth seasons

→ of American Idol?

Answer: Scott Borchetta
Correct: Scott Borchetta
Exact Match: 1
F1 Score: 1.0
Question: Thomas Stritch was an editor of which publican from Notre Dame?
Answer: The Review of Politics
Correct: Review of Politics
Exact Match: 0
F1 Score: 0.8571428571428571
Question: Paperback of the Year award from Bestsellers magazine was
   → awarded when?
Answer: 1962
Correct: 1962
Exact Match: 1
F1 Score: 1.0
Question: Gray color is often called what when referring to dogs?
Answer: blue
Correct: blue
Exact Match: 1
F1 Score: 1.0
Question: How many persons were still unaccounted for in Yingxiu?
Answer: 9,000
Correct: around 9,000
```

```
Exact Match: 0
F1 Score: 0.666666666666666
Question: How were the semi-finalists divided in season four?
Answer: by gender
Correct: by gender
Exact Match: 1
F1 Score: 1.0
Question: Who was the host of American Idol in its fourteenth season?
Answer: Ryan Seacrest
Correct: Ryan Seacrest
Exact Match: 1
F1 Score: 1.0
Question: In what borough is Godiva based?
Answer: Manhattan
Correct: Manhattan
Exact Match: 1
F1 Score: 1.0
Exact Match: average: 0.769; median: 1.0
F1 Score: average: 0.8906788305864256; median: 1.0
```

The same three cases also occur, but the overall performance is better. It pays attention to the case of letters and does not need to normalize the ground truth before calculating the metrics. In terms of the metrics, the average of Exact Match is 0.769, and the average of F1-score is 0.891. They are both higher than the model fine-tuned by myself. And no surprise that the medians for the metrics are both 1.0.

2 Complete Source Code

The degree symbol $^{\circ}$ in normalize_string can not display well in the code snippet, so $\{\hat{\}$ is used to represent it.

```
1 import numpy as np
2 | import random
3 import torch
4 import os
5 import argparse
6 | import pickle
7 import pandas as pd
8 | import unicodedata
9 | import re
10 | import statistics
11 | from transformers import BertForQuestionAnswering, TrainingArguments,
      → Trainer, BertTokenizer, pipeline
12 from datasets import Dataset
13
14 # set ramdom seed
15
   def random_seed_setting(seed):
       random.seed(seed)
16
     torch.manual_seed(seed)
```

```
18
       torch.cuda.manual_seed(seed)
19
       np.random.seed(seed)
       torch.backends.cudnn.deterministic = True
20
       torch.backends.cudnn.benchmarks = False
21
       os.environ['PYTHONHASHSEED'] = str(seed)
22
23
24
   # load datasets
25
   def load_dataset(path='./dataset'):
       with open('{}/train_dict.pkl'.format(path), 'rb') as f:
26
27
           train_dict = pickle.load(f)
           f.close()
28
29
       with open('{}/eval_dict.pkl'.format(path), 'rb') as f:
30
            eval_dict = pickle.load(f)
31
           f.close()
       with open('{}/test_dict.pkl'.format(path), 'rb') as f:
32
33
           test_dict = pickle.load(f)
34
           f.close()
       with open('{}/train_data_processed.pkl'.format(path), 'rb') as f:
35
36
           train_processed = pickle.load(f)
37
           f.close()
       with open('{}/eval_data_processed.pkl'.format(path), 'rb') as f:
38
39
           eval_processed = pickle.load(f)
40
           f.close()
41
       with open('{}/test_data_processed.pkl'.format(path), 'rb') as f:
42
           test_processed = pickle.load(f)
           f.close()
43
44
       # print keys
45
       print(train_dict.keys())
46
       print(eval_dict.keys())
47
       print(test_dict.keys())
       print(train_processed.keys())
48
49
       print(eval_processed.keys())
50
       print(test_processed.keys())
       return train_dict, eval_dict, test_dict, train_processed,
51
          → eval_processed, test_processed
52
   # fine tune model
53
54
   def train_model(train_processed, eval_processed):
55
       # initialize model
       model_name = 'bert-base-uncased'
56
57
       model= BertForQuestionAnswering.from_pretrained(model_name)
58
       print(model._modules)
59
       # set training arguments
60
       training_args = TrainingArguments(output_dir='./results', # output
          \hookrightarrow directory
61
                                           use_mps_device=False,
62
                                           num_train_epochs=10, # total
                                               → number of training epochs
63
                                           per_device_train_batch_size=8, #
                                               → batch size per device
                                               → during training
                                           per_device_eval_batch_size=8, #
64
                                               → batch size for evaluation
```

```
65
                                             weight_decay=0.01, # strength of
                                                66
                                             logging_strategy='epoch',
67
                                             logging_steps=1
68
69
        # create dataset instance
70
        train_dataset = Dataset.from_pandas(pd.DataFrame(train_processed))
71
        eval_dataset = Dataset.from_pandas(pd.DataFrame(eval_processed))
72
        # train model
73
        trainer = Trainer(model=model, # the instantiated Transformers
           → model to be fine-tuned
                           args=training_args, # training arguments, defined
74
                               → above
                           train_dataset=train_dataset, # training dataset
75
76
                           eval_dataset=eval_dataset # evaluation dataset
77
78
        trainer.train()
        # save trainer instance
79
        trainer.save_model('./results/last_model')
80
81
    # convert Unicode to ASCII
82
    def unicode2ascii(s):
83
        return ''.join(
84
            c for c in unicodedata.normalize('NFD', s)
85
86
            if unicodedata.category(c) != 'Mn'
        )
87
88
    # normalize sentence
89
   def normalize_string(s):
90
        s = unicode2ascii(s.lower().strip())
91
92
        s = re.sub(r"([,.!?%$\-\'\'\"\"(\)^{\circ}])", r" \1 ", s)
        s = re.sub(r"[^a-zA-Z0-9,.!?%$\-\'\'\"\"\(\)^{\circ}]+", r" ", s)
93
94
        s = re.sub(r'^\s+|\s+\$', '', s)
        s = re.sub(r" ^{\langle circ \rangle} ", r"^{\langle circ \rangle}", s)
95
        return s
96
97
    # convert subword-level token to word-level token
98
99
    def subword2word(s):
100
        s = re.sub(r" ##", r"", s)
101
        return s
102
    # compute Exact Match
103
    def compute_exact_match(prediction, truth):
104
105
        return int(prediction == truth)
106
107
    # compute F1-score
108
   def f1_score(prediction, truth):
109
        pred_tokens = prediction.split()
110
        truth_tokens = truth.split()
111
        # if either the prediction or the truth is no-answer then F1 = 1 if
           \hookrightarrow they agree, 0 otherwise
        if len(pred_tokens) == 0 or len(truth_tokens) == 0:
112
            return int(pred_tokens == truth_tokens)
113
114
        else:
```

```
115
            common_tokens = set(pred_tokens) & set(truth_tokens)
116
            # if there are no common tokens then F1 = 0
            if len(common_tokens) == 0:
117
                return 0
118
119
            else:
120
                prec = len(common_tokens) / len(pred_tokens)
121
                rec = len(common_tokens) / len(truth_tokens)
122
                return 2 * (prec * rec) / (prec + rec)
123
124
   # test model
   def test_model(test_dict, test_processed):
125
        # create dataset instance
126
127
        test_dataset = Dataset.from_pandas(pd.DataFrame(test_processed))
128
        # load trained model
129
        model = BertForQuestionAnswering.from_pretrained('./results/
           → last_model')
        trainer = Trainer(model=model)
130
        # Initialize the tokenizer
131
        tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
132
133
        # test trained model
        x = trainer.predict(test_dataset)
134
135
        # retrieve offset mapping of prediction
136
        start_pos, end_pos = x.predictions
137
        start_pos = np.argmax(start_pos, axis=1)
138
        end_pos = np.argmax(end_pos, axis=1)
139
        # create lists to save metrics
        EM_list = []
140
141
        F1_list = []
142
        # print predictions and calculate metrics
143
        for k, (i, j) in enumerate(zip(start_pos, end_pos)):
            # convert indices to tokens
144
145
            tokens = tokenizer.convert_ids_to_tokens(test_processed[')

    input_ids'][k])

            # preprocess prediction and ground truth
146
            prediction = subword2word(' '.join(tokens[i:j+1]))
147
            truth = normalize_string(test_dict['answers'][k]['text'][0])
148
149
            print('Question:', test_dict['question'][k])
            print('Answer:', prediction)
150
151
            print('Correct Answer:', test_dict['answers'][k]['text'][0])
            print('Normalized Correct Answer:', truth)
152
153
            em = compute_exact_match(prediction, truth)
154
            f1 = f1_score(prediction, truth)
155
            print('Exact Match:', em)
156
            print('F1 Score:', f1)
            EM_list.append(em)
157
158
            F1_list.append(f1)
            print('---')
159
160
        # print overall metrics
161
        print('\n')
        print('Exact Match: average: {}; median: {}'.format(statistics.mean
162
           print('F1 Score: average: {}; median: {}'.format(statistics.mean())
163
           → F1_list), statistics.median(F1_list)))
164
```

```
165
    # compare with another open-source fine-tuned model
166
    def compare_model(test_dict):
167
        # load model
        question_answerer = pipeline('question-answering', model='
168

    distilbert-base-cased-distilled-squad')
169
        # create lists to save metrics
170
        EM_list = []
171
        F1_list = []
172
        # print predictions and calculate metrics
173
        for i in range(len(test_dict['question'])):
            result = question_answerer(question=test_dict['question'][i],
174

    context=test_dict['context'][i])

            print('Question:', test_dict['question'][i])
175
            print('Answer:', result['answer'])
176
            print('Correct:', test_dict['answers'][i]['text'][0])
177
            em = compute_exact_match(result['answer'], test_dict['answers'
178
               → ][i]['text'][0])
            f1 = f1_score(result['answer'], test_dict['answers'][i]['text'
179
               \hookrightarrow ][0])
180
            print('Exact Match:', em)
            print('F1 Score:', f1)
181
182
            EM_list.append(em)
183
            F1_list.append(f1)
184
            print('---')
185
        # print overall metrics
186
        print('\n')
        print('Exact Match: average: {}; median: {}'.format(statistics.mean
187
           print('F1 Score: average: {}; median: {}'.format(statistics.mean(
188
           → F1_list), statistics.median(F1_list)))
189
    if __name__ == '__main__':
190
191
        # '1' -- fine tune model
192
        # '2' -- test fine-tuned model
        # '3' -- compare with another fine-tuned model
193
194
        parser = argparse.ArgumentParser()
195
        parser.add_argument('-t', '--task', type=str, default='1',\
            help='choose a task', choices=['1','2','3'])
196
197
        args = parser.parse_args()
198
199
        # set random seed
        seed = 60146
200
201
        random_seed_setting(seed)
202
203
        train_dict, eval_dict, test_dict, train_processed, eval_processed,
           → test_processed = load_dataset()
204
        if args.task == '1':
            train_model(train_processed, eval_processed)
205
206
        if args.task == '2':
207
            test_model(test_dict, test_processed)
208
        if args.task == '3':
209
            compare_model(test_dict)
```