

How Multilingual is Multilingual BERT?

He HUANG 22107447

November 23, 2023

1 Introduction

2 Fine-Tuning mBERT

```
1 import conllu
2 # Extracting PoS information from a conllu file
3 def load_conllu(filename):
4     for sentence in conllu.parse(open(filename, "rt", encoding="utf-8").read()):
5         tokenized_words = [token["form"] for token in sentence]
6         gold_tags = [token["upos"] for token in sentence]
7         yield tokenized_words, gold_tags
8 corpus = list(load_conllu("fr_sequoia-ud-test.conllu"))
```

2.1 The Universal Dependencies Project

1. Why do we need “consistent annotation” for our experiment ?

Our experiment aims to test a model across multiple languages. Consistent annotation ensures that similar linguistic structures or features are annotated in a uniform way between them. Therefore, it allows us to compare more effectively these languages that differ in grammar, syntax, and morphology.

2. What do we use the `yield` keyword rather than a simple `return` (line 7) ? Why do we have to call the `list` constructor (line 9) ?

The keyword `yield` creates a generator object, it will process one sentence at a time instead of loading the entire file, while `return` executes the function immediately and returns the results, thus it would be more time-consuming and take up a lot of memory.

The `list` will iterate over the entire generator and convert the generator into a list.

3. What is the distribution of labels in the test set of the Sequoia treebank ?

```
1 # count the number of occurrences of each POS tag
2 from collections import Counter
3 tag_counts = Counter(tag for _, tags in corpus for tag in tags)
4
5 # plot the distribution of POS tags
6 import matplotlib.pyplot as plt
```

```

7 plt.figure(figsize=(10, 8))
8 plt.bar(tag_counts.keys(), tag_counts.values())
9 plt.xlabel("Label", fontsize=10)
10 plt.ylabel("Occurrences", fontsize=10)
11 plt.title("Distribution of POS tags", fontsize=14)
12 plt.show()

```

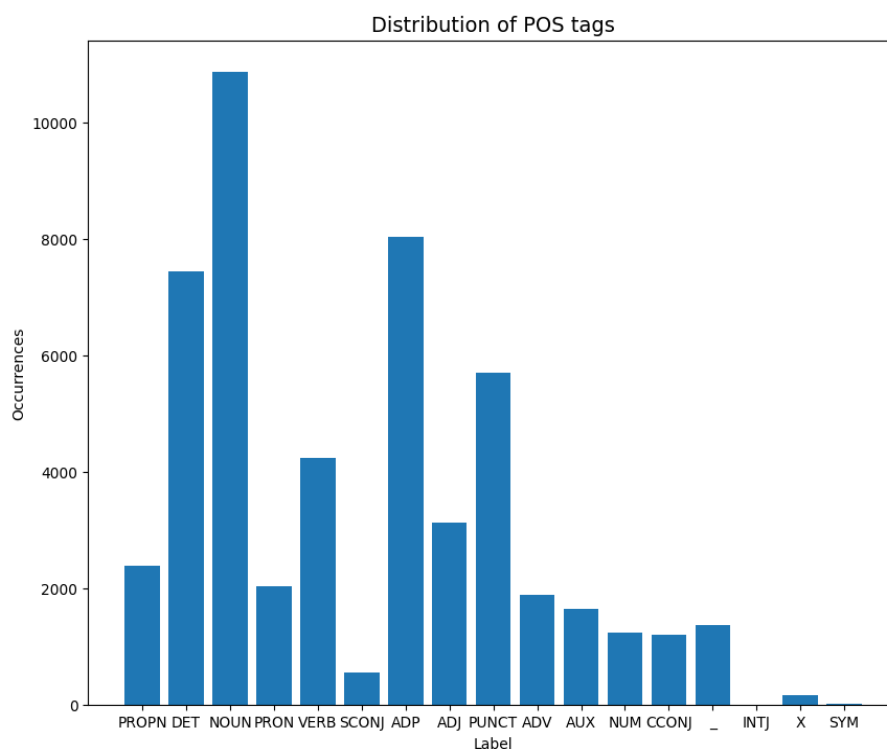


Figure 1: Distribution of POS tags in Sequoia train corpus

4. What are the multiword tokens in the test set of Sequoia? How are they annotated ?

According to the document UD French Sequoia, there are 8 types of multi-word tokens. Examples: des, du, au, aux, auxquels, auxquelles, duquel, desdites. They are annotated by label "_" or "ADP", "DET".

5. Why did the UD project make the decision of splitting multiword tokens into several (grammatical) words ? What is your opinion on this decision ?

The basic units of annotation are syntactic words (not phonological or orthographic words)¹. A single word in French "au", can be a combination of the preposition "à" and the definite article "le", which contains morphological and syntactic characteristics. Therefore, splitting these words is necessary for exploring morphosyntactic dependencies.

6. According to UD tokenization guidelines, a token can contain spaces. Are there any token in the Sequoia corpora that contain spaces ? If so these spaces should be removed.

```

1 # find the tokens that contains spaces then move the spaces

```

¹<https://universaldependencies.org/u/overview/tokenization.html>

```

2 num_space = 0
3 for sents, _ in corpus:
4     for token in sents:
5         if " " in token:
6             num_space += 1
7             token = token.replace(" ", "")
8 print(f"Number of spaces: {num_space} in Sequoia text corpus")
9
10 # Number of spaces: 13 in Sequoia text corpus

```

2.2 mBERT tokenization

7. Why does mBERT use a subword tokenization ?

Using subword tokenization allows us to have a smaller vocabulary size that can reduce the amount of memory. If we regard every possible word and wordform as a distinct token, the vocabulary size would become very large.

8. How is the sentence “Pouvez-vous donner les mêmes garanties au sein de l’Union Européene” tokenized according to the UD convention ? How is it tokenized by mBERT tokenizer ?

The result of mBERT tokenizer is:

```

1 !pip install transformers
2
3 sentence = "Pouvez-vous donner les mêmes garanties au sein de l'Union Européene"
4 # Tokenizing a single sentence with mBERT tokenizer.
5 from transformers import AutoTokenizer
6
7 tokenizer = AutoTokenizer.from_pretrained("bert-base-multilingual-cased")
8 # show the resulting tokenization
9 mbert_tokenized = tokenizer.tokenize(sentence)
10 print(mbert_tokenized)
11
12 Output:
13 ['Po', '##uve', '##z', '-', 'vous', 'donner', 'les', 'mêmes', 'gara', '##nties',
    'au', 'sein', 'de', 'l', '', 'Union', 'Euro', '##pée', '##ne']

```

According to the UD convention:

```

1 ['Pouvez', '-vous', 'donner', 'les', 'mêmes', 'garanties', 'à', 'le', 'sein', 'de',
    'l', '', 'Union', 'Européenne', '?']

```

9. Why is this difference in tokenization a problem for training a mBERT-based PoS tagger ?

A word with UD convention that’s tokenized as a single unit might be tokenized into multiple subword in mBERT model. This difference directly affects the size of the vocabulary, and the learning and prediction ability.

2.3 Reconciling the Two Tokenizations

10. Write a function that takes as parameter a sentence tokenized according to UD rules (i.e. a variable of type List[str] and its gold labels (also a variable of type List[str]) and implements the first principle explained above.

```

1 # first principle: concatenation of the label for multiword tokens
2 def concat_labels(tokens, labels):
3     updated_ud_labels, updated_ud_tokens = [], []

```

```

4
5     i = 0
6     while i < len(tokens):
7         updated_ud_tokens.append(tokens[i])
8
9         # this is a maker of multiword tokens
10        if labels[i] == "_":
11            updated_ud_labels.append(labels[i+1] + "+" + labels[i+2])
12            i += 2
13        else:
14            updated_ud_labels.append(labels[i])
15            i += 1
16    return updated_ud_tokens, updated_ud_labels
17
18 # Apply the principle to the corpus corresponding to the UD dataset
19 def up_corpus(corpus):
20     new_ud_corpus = []
21
22     for tok, tag in corpus:
23         new_ud_tokens, new_ud_labels = concat_labels(tok, tag)
24         new_ud_corpus.append((new_ud_tokens, new_ud_labels))
25     return new_ud_corpus

```

11. Write a function that takes as input the sentences and labels created by “normalizing” UD sentences, apply the mBERT tokenizer and compute the corresponding label. You must ensure each subtoken (including padding symbols) has a label.

```

1 # second principle: normalization the corresponding labels
2 from copy import deepcopy
3
4 def normalize_labels(corpus, tokenizer, MAX_LENGTH=50):
5     """
6     corpus is a list of tuples updated (ud_tokens, ud_labels)
7     """
8     updated_mbert_labels = []
9     updated_mbert_corpus = []
10
11    for data in corpus:
12        texts = data[0]
13        texts_tags = data[1]
14
15        model = tokenizer(texts,
16                           is_split_into_words=True,
17                           return_offsets_mapping=True,
18                           max_length=MAX_LENGTH,
19                           padding='max_length',
20                           truncation=True)
21
22        tokens = tokenizer.tokenize(texts,
23                                    is_split_into_words=True,
24                                    padding=True,
25                                    truncation=True)
26        tokens.insert(0, '<pad>')
27        tokens.append('<pad>')
28
29        input_ids = model["input_ids"]
30        attention_masks = model["attention_mask"]
31

```

```

32     updated_mbert_labels = [None] * len(model["offset_mapping"])
33     original_labels = deepcopy(texts_tags)
34
35     for i in range(len(model["offset_mapping"])):
36         if model["offset_mapping"][i][0] != 0 or model["offset_mapping"][i
37 ] [1] == 0:
38             updated_mbert_labels[i] = '<pad>'
39         else:
40             updated_mbert_labels[i] = original_labels[0]
41             original_labels.pop(0)
42     updated_mbert_corpus.append((tokens, updated_mbert_labels, input_ids,
43 attention_masks))
44     return updated_mbert_corpus

```

12. Write a function that encodes the labels into integers.

```

1 # Encodes the mbert labels into integers: <pad> label must be mapped to -100.
2 def create_l2i(mbert_corpus_updated):
3
4     all_labels = set()
5     labels_to_idx = {}
6
7     for pair in mbert_corpus_updated:
8         for labels in pair[1]:
9             all_labels.add(labels)
10    labels_to_idx = {label: i for i, label in enumerate(all_labels)}
11    labels_to_idx["<pad>"] = -100
12    return labels_to_idx

```

2.4 Creating a Dataset

13. Using the `Dataset.from_list` method, write a method that creates a `Dataset` that encapsulates a corpus in the conllu.

```

1 # create the dataset
2 !pip install datasets
3
4 from datasets import Dataset
5 def create_dataset(corpus):
6     """
7     corpus is a list of tuples (mbert_tokens, mbert_labels, input_ids,
8     attention_masks)
9     """
10    examples = []
11
12    for data in corpus:
13        # Encode the labels into integers
14        ids_labels = [l2i[label] for label in data[1]]
15        examples.append({"input_ids": data[2], "attention_mask": data[3], "
16 labels": ids_labels})
17    ds = Dataset.from_list(examples)
18    return ds

```

14. Create three instances of `Dataset` : one for the train set, one for the dev set and the last one for the test set.

```

1 # create three instances of the dataset
2
3 # train set
4 train_corpus = list(load_conllu(trainfile))
5
6 train_ud_corpus = up_corpus(train_corpus)
7 train_mbert_corpus = normalize_labels(train_ud_corpus,tokenizer)
8 l2i = create_l2i(train_mbert_corpus)
9 ds_train = create_dataset(train_mbert_corpus)
10
11 # dev set
12 dev_corpus = list(load_conllu(devfile))
13
14 dev_ud_corpus = up_corpus(dev_corpus)
15 dev_mbert_corpus = normalize_labels(dev_ud_corpus,tokenizer)
16 ds_dev = create_dataset(dev_mbert_corpus)
17
18 # test set
19 test_corpus = list(load_conllu(testfile))
20
21 test_ud_corpus = up_corpus(test_corpus)
22 test_mbert_corpus = normalize_labels(test_ud_corpus,tokenizer)
23 ds_test = create_dataset(test_mbert_corpus)

```

2.5 Fine-Tuning mBERT

15. How can you modify the code of Figure 4 to report the PoS tagging accuracy during optimization. Why is this information important ?

I add a function `compute_metrics` to use with loss as the model evaluation, and pass it to the trainer :

```

1 # add the evaluation metric
2 !pip install evaluate
3
4 import numpy as np
5 import evaluate
6
7 from transformers import AutoModelForTokenClassification, TrainingArguments,
8     Trainer
9 model_checkpoint = "bert-base-multilingual-cased"
10 model = AutoModelForTokenClassification.from_pretrained(model_checkpoint,
11     num_labels=len(l2i))
12
13 # evaluate the model
14 def compute_metrics(eval_pred):
15     metric = evaluate.load("accuracy")
16     logits, labels = eval_pred
17     predictions = np.argmax(logits, axis=-1)
18     predictions = predictions[:,1]
19     labels = labels[:,1]
20     return metric.compute(predictions=predictions, references=labels)
21
22 def train_process(ds_train, ds_dev, batch_size=16):
23
24     training_args = TrainingArguments(
25         output_dir="./results",
26         evaluation_strategy = "epoch",
27         learning_rate=2e-5,
28         per_device_train_batch_size=batch_size,

```

```

27     per_device_eval_batch_size=batch_size,
28     num_train_epochs=3,
29     weight_decay=0.01,
30     logging_dir='./logs',
31     logging_steps=10,
32     optim = "adamw_torch") # avoid deprecation warning
33
34     trainer = Trainer(
35         model=model,
36         args=training_args,
37         train_dataset=ds_train, # the train set
38         eval_dataset=ds_dev, # the dev set
39         compute_metrics=compute_metrics
40     )
41     trainer.train()
42     trainer.evaluate()
43
44 train_process(ds_train, ds_dev)

```

| Epoch | Training Loss | Validation Loss | Accuracy |
|-------|---------------|-----------------|----------|
| 1 | 0.135300 | 0.107672 | 0.961165 |
| 2 | 0.066700 | 0.074755 | 0.983010 |
| 3 | 0.044800 | 0.069465 | 0.983010 |

Figure 2: Loss and accuracy on the french train and dev set

3 Evaluating the multilingual capacity of mBERT

16. Choose 5 languages from the UD project. For each language train a PoS tagger using the code of the previous section and test it on the 5 languages you have chosen.

For the experiment, I chose five languages, they are Vietnamese, Korean, Dutch, Italian and Turkish. Here are the number of training sets, validation sets for each experimental language :

| Language | Train set | Dev set |
|------------|-----------|---------|
| Vietnamese | 1400 | 1123 |
| Korean | 4400 | 950 |
| Dutch | 5789 | 676 |
| Italian | 1781 | 156 |
| Turkish | 3435 | 1100 |

17. Can you use the same hyperparameters for the different languages ?

For each language, I didn't change the hyperparameters during the training.

18. What can you conclude ?

The findings indicate that among the five languages evaluated, Vietnamese and Turkish yielded accuracies below 80%. In contrast, the remaining three languages—Korean, Italian, and Dutch—surpassed the 90% accuracy threshold. Notably, Dutch achieved the pinnacle of precision with a remarkable 97% accuracy rate.

In terms of corpus size, both Vietnamese and Italian have considerably smaller datasets, each with fewer than 2,000 training examples. Despite this, Italian demonstrates a significantly higher precision rate of 91%, which is markedly greater than that of Vietnamese. In addition, Turkish has more than 3000 training examples, but the accuracy is not much different from that of Vietnamese.

Vietnamese is an isolating language, where word meanings are typically modified through the use of unique words rather than changing word forms directly. On the other hand, Turkish is an agglutinative language characterized by complex words composed of multiple affixes, each denoting distinct grammatical features. mBERT employs subword tokenization, it struggles to capture meanings in uninflected words as in Vietnamese, and similarly faces difficulty in managing the extended word structures inherent to Turkish.

Here are all the results:

| Epoch | Training Loss | Validation Loss | Accuracy |
|-------|---------------|-----------------|----------|
| 1 | 0.573100 | 0.547699 | 0.745325 |
| 2 | 0.394600 | 0.470056 | 0.782725 |
| 3 | 0.297800 | 0.465663 | 0.795191 |

Figure 3: Loss and accuracy on the Vietnamese train and dev set

| Epoch | Training Loss | Validation Loss | Accuracy |
|-------|---------------|-----------------|----------|
| 1 | 0.270400 | 0.223100 | 0.934737 |
| 2 | 0.145400 | 0.192093 | 0.949474 |
| 3 | 0.135200 | 0.178400 | 0.955789 |

Figure 4: Loss and accuracy on the Korean train and dev set

| Epoch | Training Loss | Validation Loss | Accuracy |
|-------|---------------|-----------------|----------|
| 1 | 0.139600 | 0.124084 | 0.970414 |
| 2 | 0.081400 | 0.107161 | 0.977811 |
| 3 | 0.049300 | 0.108571 | 0.976331 |

Figure 5: Loss and accuracy on the Dutch train and dev set

| Epoch | Training Loss | Validation Loss | Accuracy |
|-------|---------------|-----------------|----------|
| 1 | 0.180000 | 0.196991 | 0.903846 |
| 2 | 0.110600 | 0.159628 | 0.923077 |
| 3 | 0.099000 | 0.152014 | 0.916667 |

Figure 6: Loss and accuracy on the Italian train and dev set

| Epoch | Training Loss | Validation Loss | Accuracy |
|-------|---------------|-----------------|----------|
| 1 | 0.520900 | 0.541689 | 0.738182 |
| 2 | 0.366900 | 0.454626 | 0.787273 |
| 3 | 0.283200 | 0.430282 | 0.795455 |

Figure 7: Loss and accuracy on the Turkish train and dev set