

Multilingual NLP

Lab 2 – How Multilingual is Multilingual BERT?

AN Ji

November 19, 2024

1 Introduction

Installation:

- HuggingFace: transformer, datasets

2 Fine-Tuning mBERT

Goals:

- Evaluate the performance of a PoS tagger fine-tuned on a given language and compare it to the performance of a PoS tagger trained only on this language
- Evaluate on a language B the performance of a PoS tagger fine-tuned on a language A

2.1 The Universal Dependencies Project

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

For this experiment, we need to train a PoS tagger, evaluate and compare its performance on different languages. Applying a common set of PoS tags cross-linguistically allows us to do the job correctly and more easily, as well as other comparative morphosyntactic researches using UD projects. It's therefore crucial to being “uniform”, although this method could imply some kind of compromise between languages that differs a lot in some subtle morphosyntactic aspects.

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)?

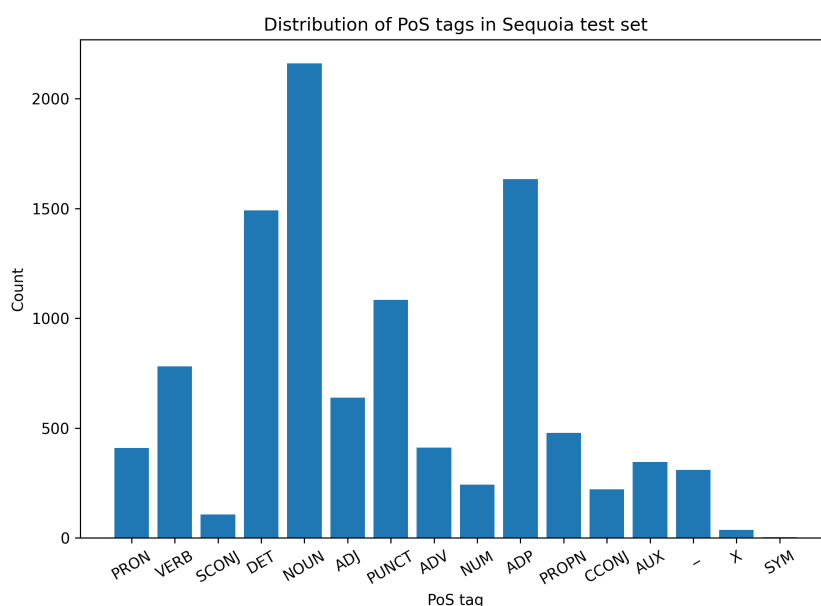
```
1 import conllu
2
3 def load_conllu(filename):
4     for sentence in conllu.parse(open(filename, "rt",
5         encoding="utf-8").read()):
6         tokenized_words = [token['form'] for token in sentence]
7         gold_tags = [token["upos"] for token in sentence]
8         yield tokenized_words,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

The return keyword requires to load all data into memory at once before returning a list, which may often be intensive to memory use, especially when we should deal with extremely large datasets and haven't decided how many data to process with that result later. Using yield to

create a generator first instead of returning directly a huge list of sentences is a better and versatile choice on both memory management and data flexibility. So when we decide to use all data, we can simply combine `list()` with `yield` to transform the generator into a full list of sentences, like what we do in line 9.

3. Distribution of labels in the test set of the Sequoia treebank

```
1 from collections import Counter
2 import matplotlib.pyplot as plt
3
4 def plot_pos_distribution(corpus):
5     """plot the distribution of PoS tags in the corpus"""
6     pos_counts = Counter(tag for _, tags in corpus for tag in tags)
7     plt.figure(figsize=(9, 6))
8     plt.bar(pos_counts.keys(), pos_counts.values())
9     plt.xlabel("PoS tag")
10    plt.ylabel("Count")
11    plt.xticks(rotation=30)
12    plt.title("Distribution of PoS tags in Sequoia test set")
13    plt.show()
14
15 plot_pos_distribution(corpus)
```



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

```
1 multi_tokens = [token.lower() for tokens, tags in corpus for token, tag in
2                 zip(tokens, tags) if tag == '_']
3 dict(Counter(multi_tokens))
```

```
{'des': 114, 'du': 115, 'aux': 22, 'au': 58, 'desdites': 1}
```

From the conllu file of Sequoia test set, we find the following annotation rule for the multiword tokens above: the token itself is marked with an underscore as `upos`, followed by its 2 constituents labeled as `ADP` (adposition) and `DET` (determiner) respectively. For example, `desdites` in the sentence of line 195 and `des` in the sentence of line 451:

En tant que remarque isolée, je voudrais donc remercier Mme Cederschiöld du fait qu'elle ait clarifié par exemple la notion desdites copies temporaires dans son avis.

ID	FORM	LEMMA	UPOS
24-25	desdites		
24	de	de	ADP
25	lesdites	ledit	DET

Comme vous tous, j'entends les appels à l'aide des mères, des épouses et des soeurs des Kosovars albanais qui croupissent dans les geôles serbes.

ID	FORM	LEMMA	UPOS
12-13	des		
12	de	de	ADP
13	les	le	DET

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

Since it is emphasized that the basic units of annotation are **syntactic words** (not phonological or orthographic words)¹, French contractions are actually not compatible with any of the PoS tags defined in UD tokenization guidelines. They are mostly composed of a preposition (marked as ADP, like *à*, *de*) and a determiner (marked as DET, like *le*, *les*). If we choose only one, it would become unavoidable to lose important syntactic information of the ignored component; if we introduce a new tag for them, we will be forced to decrease consistency and comparability from a cross-linguistic perspective. Hence splitting these forms seems to be rather a reasonable solution (a step towards de-specialization?).

6. Tokens with spaces inside are numbers:

```
1 tokens_spaces = {token for tokens, _ in corpus for token in tokens if ' ' in token}
2 tokens_spaces
3
4 [Out]:
5
6 {'100 000', '15 000', '190 500', '2 000', '3 092', '3 852', '3 862', '50 000', '500 000', '67 025', '7 736', '80 000', '800 000'}
```

Remove spaces:

```
1 corpus = [[token.replace(' ', '') for token in tokens], tags] for tokens, tags in corpus]
```

2.2 mBERT Tokenization

7. Why does mBERT use a subword tokenization?

Subword tokenization algorithms, here **WordPiece** used in mBERT, are able to keep frequently used words as tokens on the one hand, and break down relatively less used and unseen words into frequently occurring subword units that the model has seen and learned somewhere in the pretraining data. This approach reduces notably vocabulary size and memory use while maximizing the model's semantic representation of data, making it more efficient in multilingual text processing (e.g. for morphologically rich languages).

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

¹UD Guidelines – Tokenization and word segmentation

```

1 # Tokenization by UD convention:
2
3 ['Pouvez', '-vous', 'donner', 'les', 'mêmes', 'garanties', 'au', 'à', 'le',
  'sein', 'de', 'l', 'Union', 'européenne', '?']
4
5 # Tokenization by mBERT tokenizer:
6
7 !pip install transformers
8 from transformers import AutoTokenizer
9
10 tokenizer = AutoTokenizer.from_pretrained("bert-base-multilingual-cased")
11 sentence = "Pouvez-vous donner les mêmes garanties au sein de l'Union
  Européene ?"
12 print(tokenizer.tokenize(sentence))
13
14 [Out]:
15 ['Po', '##uve', '##z', '-', 'vous', 'donner', 'les', 'mêmes', 'gara',
  '##nties', 'au', 'sein', 'de', 'l', '', 'Union', 'Euro', '##pée', '##ne',
  '?']

```

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

Because if we want to train an mBERT-based PoS tagger, we had better ensure that all our training tokens conform with mBERT's inherent vocabulary forms so that knowledge about them already mastered by mBERT can be optimally transferred and facilitate our training. If we don't solve this difference, then the model's pre-training vocabulary advantages would not be maximally used and it will produce a significant impact upon final performance, since it means increasing vocabulary size as well as built-in mappings unnecessarily.

2.3 Reconciling 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 def process_ud_multiwords(tokens, tags):
2     """
3     multiword tokens: keep contractions unchanged.
4     E.g. 'aux','à','les' -> 'aux'
5     multiword tags: replace '_' with the following 2 tags.
6     E.g. ' ','ADP','DET' -> 'ADP+DET'
7     """
8     updated_tokens, updated_tags = list(), list()
9     i = 0
10    while i < len(tokens):
11        updated_tokens.append(tokens[i])
12        if tags[i] == '_':
13            updated_tags.append(f'{tags[i+1]}+{tags[i+2]}')
14            i += 3
15        else:
16            updated_tags.append(tags[i])
17            i += 1
18    return updated_tokens, updated_tags

```

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 # implement 2nd principle

```

```

2
3 def tokenize_align_tags(corpus, tokenizer=tokenizer, MAX_LENGTH=64):
4     """
5     Let mBERT re-tokenize each token, align sub-tokens with tags, pad
        sentences to MAX_LENGTH, save model inputs (& raw tokens/tags) of each
        sentence in a dict, then all dicts in a list
6     """
7     all_sentences = list()
8     for tokens, tags in corpus:
9         sentence = dict()
10        model_inputs = tokenizer(
11            tokens,
12            is_split_into_words=True,
13            # return_offsets_mapping=True, # no need if using word_ids
14            # truncation=True, # no need if setting max_length
15            padding='max_length',
16            max_length=MAX_LENGTH,
17        )
18        # save model inputs
19        sentence['subtokens'] = model_inputs.tokens()
20        sentence['input_ids'] = model_inputs['input_ids']
21        sentence['attention_mask'] = model_inputs['attention_mask']
22        # use for tag alignment e.g. [None, 0, 0, 0, 1, 1, 2, 3, 4, 4, 5, None]
23        word_ids = model_inputs.word_ids()
24
25        aligned_tags = []
26        previous_word_id = None
27
28        for word_id in word_ids:
29            if word_id is None: # [CLS], [SEP] & [PAD]
30                aligned_tags.append('<pad>')
31            elif word_id != previous_word_id: # word head
32                aligned_tags.append(tags[word_id])
33            else: # rest of word
34                aligned_tags.append('<pad>')
35
36            previous_word_id = word_id
37
38        sentence['tags'] = aligned_tags
39        all_sentences.append(sentence)
40    return all_sentences

```

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

See 14. For convenience I handle this step after creating the multilingual dataset.

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.

Here I combined functions defined above to transform a conllu corpus to a Dataset:

```

1 def conllu2ds(lang, split):
2     """
3     Encapsulate one conllu corpus into HuggingFace datasets
4     """
5     filename = pick_conllu(lang, split)
6     corpus = list(load_conllu(filename))
7     updated_corpus = [process_ud_multiwords(tokens, tags) for tokens, tags in
8                       corpus]
9     all_sentences = tokenize_align_tags(updated_corpus, tokenizer)
10    dataset = Dataset.from_list(all_sentences)

```

```
10 return dataset
```

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 code2lang = {
2     'zh': 'Chinese',
3     'yue': 'Cantonese',
4     'ja': 'Japanese',
5     'ug': 'Uyghur',
6     'sa': 'Sanskrit',
7     'th': 'Thai',
8 }
9
10 def encode_ds_tags():
11     """
12     First combine all language datasets, then encode all possible tags into
13     integers
14     """
15     ds_dict = DatasetDict()
16     for code in code2lang.keys():
17         ds = conllu2ds(code, 'train') if code == 'ug' else conllu2ds(code,
18             'test')
19         ds_dict.update({code: ds})
20
21     unique_tags = set([tag for _, ds in ds_dict.items() for sent in ds for tag
22         in sent['tags'] if tag != '<pad>'])
23     unique_tags = sorted(list(unique_tags))
24     tag2int = {tag: i for i, tag in enumerate(unique_tags)}
25     tag2int['<pad>'] = -100
26     # print(f"Overall 'meaningful' unique tags: {unique_tags}")
27     # print(f"Final mapping of tags to integers (including '<pad>'):"
28         f"{tag2int}")
29     # print(f"Final total of unique tags (including '<pad>'): {len(tag2int)}")
30
31     for lang, ds in ds_dict.items():
32         ds = ds.map(lambda x: {'labels': [tag2int[tag] for tag in x['tags']]})
33         ds_dict[lang] = ds
34
35     return ds_dict
36
37 # Generate an integer-labeled dataset for all languages
38 ds_dict = encode_ds_tags()
39 ds_dict
40
41 [Out]:
42 Overall 'meaningful' unique tags:
43 ['ADJ', 'ADP', 'ADV', 'AUX', 'CCONJ', 'DET', 'INTJ', 'NOUN', 'NUM', 'PART',
44 'PRON', 'PROPN', 'PUNCT', 'SCONJ', 'SYM', 'VERB', 'X']
45 Final mapping of tags to integers (including '<pad>'):
46 {'ADJ': 0, 'ADP': 1, 'ADV': 2, 'AUX': 3, 'CCONJ': 4, 'DET': 5, 'INTJ': 6,
47 'NOUN': 7, 'NUM': 8, 'PART': 9, 'PRON': 10, 'PROPN': 11, 'PUNCT': 12,
48 'SCONJ': 13, 'SYM': 14, 'VERB': 15, 'X': 16, '<pad>': -100}
49 Final total of unique tags (including '<pad>'):
50 18
51
52 DatasetDict({
53     zh: Dataset({
54         features: ['subtokens', 'input_ids', 'attention_mask', 'tags',
55             'labels'],
56         num_rows: 1004
57     })
58 })
```

```

51     yue: Dataset({
52         features: ['subtokens', 'input_ids', 'attention_mask', 'tags',
53                 'labels'],
54         num_rows: 1004
55     })
56     ja: Dataset({
57         features: ['subtokens', 'input_ids', 'attention_mask', 'tags',
58                 'labels'],
59         num_rows: 1000
60     })
61     ug: Dataset({
62         features: ['subtokens', 'input_ids', 'attention_mask', 'tags',
63                 'labels'],
64         num_rows: 1656
65     })
66     sa: Dataset({
67         features: ['subtokens', 'input_ids', 'attention_mask', 'tags',
68                 'labels'],
69         num_rows: 2709
70     })
71     th: Dataset({
72         features: ['subtokens', 'input_ids', 'attention_mask', 'tags',
73                 'labels'],
74         num_rows: 1000
75     })
76 })
77
78 def split_dataset(ds, train_size=0.8, seed=42):
79     splits = ds.train_test_split(train_size=train_size, seed=seed,
80                                 shuffle=True)
81     train_ds, temp_ds = splits['train'], splits['test']
82     temp = temp_ds.train_test_split(train_size=0.5, seed=seed)
83     dev_ds, test_ds = temp['train'], temp['test']
84
85     return DatasetDict({
86         'train': train_ds,
87         'dev': dev_ds,
88         'test': test_ds
89     })
90
91 for lang, ds in ds_dict.items():
92     ds_dict[lang] = split_dataset(ds)

```

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?**

We can add a `compute_metrics` function to design the evaluation steps and pass it to the `Trainer` function as an argument. It's crucial because we need to check the reported accuracy and loss on both train set and dev set to make sure the model converges well instead of ending up with overfitting. We also need to use accuracy as metric to report the performance on all test sets.

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.**

The following languages are chosen.

- Cantonese (yue, Sino-Tibetan, Chinese): HK

- Chinese (zh, Sino-Tibetan, Chinese): also HK, **parallel with Cantonese**
- Japanese (ja, Japanese): PUD
- Thai (th, Tai-Kadai): also PUD, **parallel with Japanese**
- Uyghur (ug, Turkic, Southeastern): UDT
- Sanskrit (sa, Indo-European, Indic): Vedic

```

1 def compute_metrics(pred):
2     """ metrics computing setup """
3     metrics = evaluate.load("accuracy")
4     logits, labels = pred
5     predictions = np.argmax(logits, axis=-1)
6     predictions = predictions.flatten()
7     labels = labels.flatten()
8     mask = labels != -100
9
10    return metrics.compute(predictions=predictions[mask],
11                           references=labels[mask])
12
13 # training settings
14 batch_size = 16
15 model_checkpoint = "bert-base-multilingual-cased"
16
17 results = {lang: {test_lang: None for test_lang in ds_dict.keys()} for lang in
18               ds_dict.keys()}
19
20 for lang, ds in ds_dict.items():
21     ds_train, ds_dev, ds_test = ds['train'], ds['dev'], ds['test']
22     model = AutoModelForTokenClassification.from_pretrained(
23         model_checkpoint,
24         num_labels=18,
25     )
26
27     training_args = TrainingArguments(
28         output_dir=f"./results/{lang}",
29         eval_strategy="epoch",
30         learning_rate=2e-5,
31         per_device_train_batch_size=batch_size,
32         per_device_eval_batch_size=batch_size,
33         num_train_epochs=3,
34         weight_decay=0.01,
35         logging_dir=f"./results/logs",
36         logging_steps=10,
37         report_to="none",
38         push_to_hub=True,
39         hub_model_id=f"ankei/{lang}_pos_tagger",
40         hub_strategy="end",
41     )
42
43     trainer = Trainer(
44         model=model,
45         args=training_args,
46         train_dataset=ds_train,
47         eval_dataset=ds_dev,
48         compute_metrics=compute_metrics,
49     )
50
51     print(f"Training {code2lang[lang]} ...")
52     trainer.train()
53     trainer.evaluate()
54
55     for test_lang, test_ds in ds_dict.items():

```



```

55     test_results = trainer.predict(test_ds["test"])
56     test_accuracy = test_results.metrics["test_accuracy"]
57     results[lang][test_lang] = test_accuracy
58     print(f"Accuracy of {code2lang[lang]} model on {code2lang[test_lang]}
        test set: {test_accuracy:.4f}")

```

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

Theoretically yes, and I didn't change the hyperparameters for the final version of fine-tuning (i.e. using `batch_size=16`, `epoch=3`, `max_length=64` for padding and truncation). However, I tested a session out of curiosity using `epoch=10` to see how the models would perform, and I found the overall loss and accuracy became the most stable around 5th epoch and started to fluctuate heavily after. So based on this, I thought maybe for each trained language we could increase epoch accordingly to make each model better, but after all since the datasets are limited to be small on purpose (from 1000 to 2700), it is hardly possible to boost performance only by modifying epoch, which might not help much in terms of models' cross-linguistic generalization ability.

18. What can you conclude?

Language-specific performance

For each language's sake, the model fine-tuned on Japanese yields the best predictive ability (95.8%), followed by Chinese (91.0%), Cantonese (85.7%) and Thai (85.2%). Japanese, Chinese and Thai are well supported and represented during mBERT's pre-training process because of their high-quality Wikipedia contents, so their performance is less surprising. Interestingly, Cantonese was not officially chosen², but the Cantonese model is still not bad compared to Thai model (at least their performance gap is usually slight after comparing several trainings). Reasons for that are various, but the most important one may be that Cantonese benefits a lot from Chinese since there is considerable vocabulary overlap. In addition, we used parallel corpora for these two languages, which are mainly composed of Hong Kong student film subtitles and Hong Kong Legislative Council's meeting records. Given that the Chinese subtitles still contain many Cantonese-specific words, meetings tend to remain lexically-formal even though they were held in Cantonese before being translated into Chinese records, and mBERT's Chinese pre-training also relies largely on Traditional Chinese contents that are de facto an amalgam of different non-mainland Chinese varieties (at least so lexically), it's no wonder that Cantonese model did well.

While the fine-tuned model excels on the above languages, we should also notice that for the rest – Uyghur and Sanskrit, that mBERT has not been explicitly pre-trained on, it still yields good results with an accuracy of 74.5% and 77.0% respectively after learning from only 1490 and 2438 sentences. On the other hand, they are also the most difficult to predict. By now we could not explain granularly enough on these observations, perhaps numerous factors can account for them more or less: lack of pre-training resources, quality and size of corpus used for fine-tuning, morphological features, alphabetical idiosyncrasies, etc. It would be reasonable to expect a much better performance using more data. We could say that these results benefit a lot from the WordPiece subword tokenization strategy that mBERT applied to set up and refine its vocabulary, which enables it to handle low-resource synthetic languages (Uyghur and Sanskrit here) more efficiently.

Cross-linguistic generalization

– Chinese model produced the most spread out performances on other languages (from Sanskrit 10.8% to Cantonese 78.1%), followed by Cantonese model (from Sanskrit 25.4% to Chinese 82.0%).

²See <https://github.com/google-research/bert/blob/master/multilingual.md#list-of-languages>

Trained on Tested on	zh	yue	ja	ug	sa	th
zh	0.910364	0.819795	0.612512	0.401494	0.324930	0.428571
yue	0.780992	0.857438	0.547521	0.393251	0.265840	0.363636
ja	0.560086	0.457797	0.957797	0.424177	0.293991	0.404149
ug	0.250252	0.278449	0.442598	0.744713	0.389225	0.286506
sa	0.108095	0.254195	0.229023	0.346989	0.770484	0.201876
th	0.418750	0.382955	0.373864	0.293182	0.336932	0.851705

Table 1: Performance of fine-tuned pos-taggers

This may indicate that models fine-tuned on morphologically-simple (isolating) languages would predict less accurately on morphologically-complex ones.

– Models trained on Chinese and Cantonese that are two members of sinitic languages show high transferability between them. This is expected given their close linguistic relatedness.

– If we ignore the reciprocal high transferability between Chinese and Cantonese, there is an obvious transfer gap of about 33%–43% between each language model’s performances on its dominant language and the best-performed other language. Additionally, we can also observe that models trained on languages with relatively rich morphological features (like Uyghur and Sanskrit) yield narrowest generalization performance range (12%–13%), while models trained on more isolating languages spread out much more heavily (Chinese 45.0%, Japanese 38.4%, Thai 22.7% and Cantonese 20.4%). Is it a hint that when predicting other languages, models based on morphologically complex languages tend to produce smaller deviation than those trained on simple ones?

– Horizontally speaking, Sanskrit, an ancient language with a rich inflectional morphology, is the least predictable among all, while Japanese (agglutinating, and a bit isolating?) and Chinese (isolating) are the opposite (although the performances on them are still limited).

– Thai model seems to be struggling to generalize since it shows largest transfer gap of at least 42% when predicting other languages. More interestingly, it reports the worst score on Sanskrit from which Thai historically borrowed and derived a lot of words.

The above results seem to reflect some kind of mapping (or correlation) between morphosyntactic distinctions and language-specific fine-tuned model’s cross-linguistic generalization ability. We can say that mBERT is multilingual enough in terms of its compatibility and performance on languages it is fine-tuned on, but mBERT’s transferability proves to be relatively limited and might correlate with typological relatedness between languages.