

Open in app ↗



Search Medium



Published in Towards Data Science



Ruben Winastwan

Follow

May 3, 2022 · 11 min read · ✨ · Listen



Save



Named Entity Recognition with BERT in PyTorch

How to leverage a pre-trained BERT model for custom data to predict the entity of each word in a text



459



7





Photo by [Aaron Burden](#) on [Unsplash](#)

When it comes to dealing with NLP problems, BERT oftentimes comes up as a machine learning model that we can count on in terms of its performance. The fact that it's been pre-trained on more than 2,500M words and its bidirectional nature to learn information from a sequence of words makes it a powerful model to use.

I wrote about how we can leverage BERT for text classification before, and in this article, we're going to focus more on how to use BERT for named entity recognition (NER) tasks.



What is NER?

NER is a task in NLP to identify and extract meaningful information (or we can call it entities) in a sentence or text. An entity can be a single word or even a group of words that refer to the same category.

As an example, let's say we the following sentence and we want to extract information about a person's name from this sentence.

The name is Bond, James Bond

The first step of a NER task is to detect an entity. This can be a word or a group of words that refer to the same category. As an example:

- ‘*Bond*’  an entity that consists of a single word
- ‘*James Bond*’  an entity that consists of two words, but they are referring to the same category.

To make sure that our BERT model knows that an entity can be a single word or a group of words, then we need to provide information about the beginning and the ending of an entity on our training data via the so-called Inside-Outside-Beginning (IOB) tagging. We will see more about this on our dataset later in this article.

After detecting an entity, the next step in a NER task is to categorize the detected entity. The categories of an entity can be anything depending on our use case. Below is an example of categories of entities:

- **Person:** Bond, James Bond, Sam, Anna, Frank, Leonardo DiCaprio
- **Location:** New York, Vienna, Munich, London
- **Organization:** Google, Apple, Stanford University, Deutsche Bank
- **Location:** Central Park, Brandenburger Tor, Times Square

These entities are basically the label of our data during the training process of our BERT model, which we will look at in detail later in the following section.

BERT for NER

As previously mentioned, BERT is a transformers-based machine learning model that will come in pretty handy if we want to solve NLP-related tasks.

If you're not yet familiar with BERT, I recommend you to read my previous article about text classification with BERT before reading this article. There you'll find information about what BERT actually is, what kind of input data the model expects, and the output that you'll get from the model.

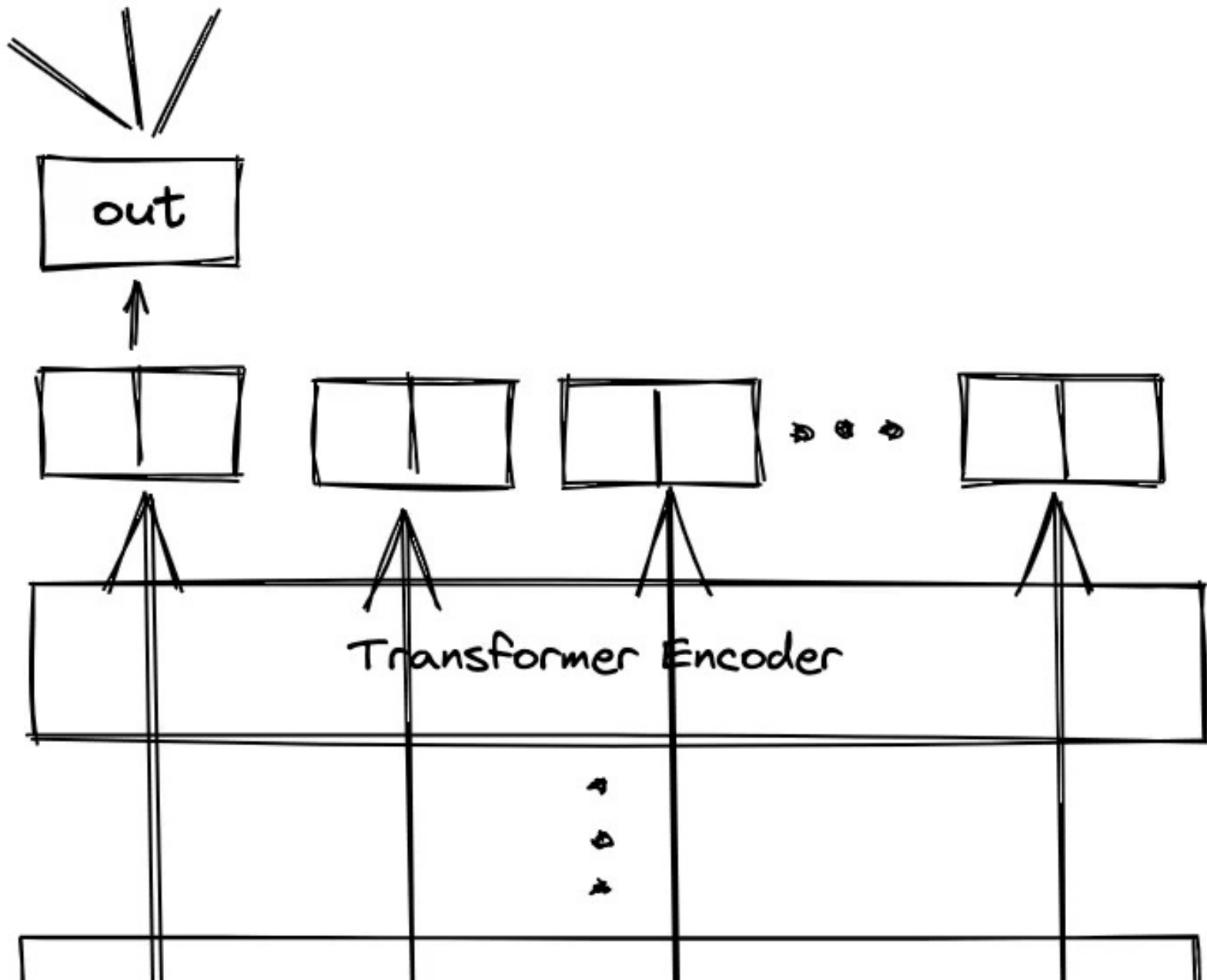
Text Classification with BERT in PyTorch

How to leverage a pre-trained BERT model from Hugging Face to classify text of news articles

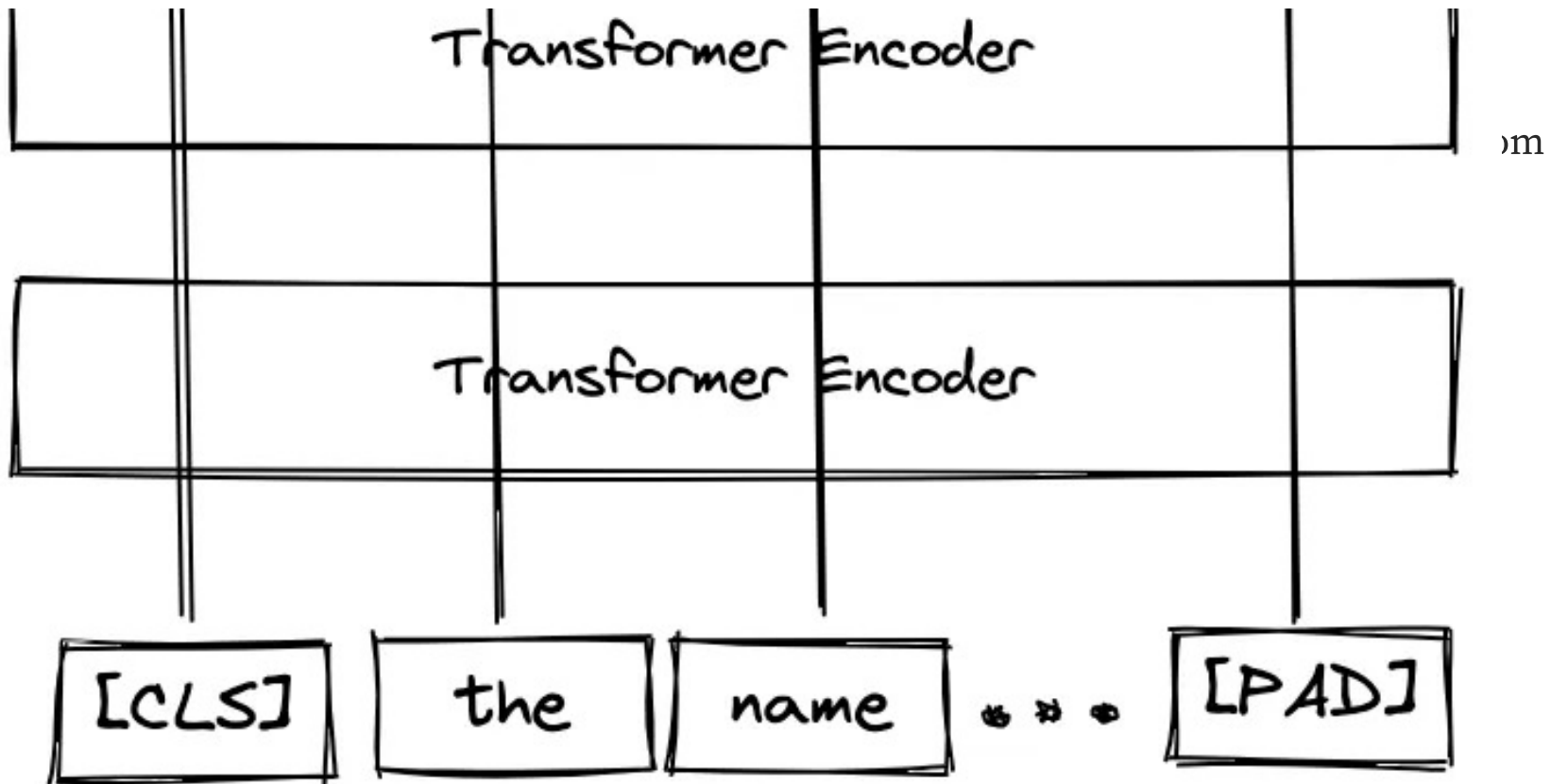
towardsdatascience.com

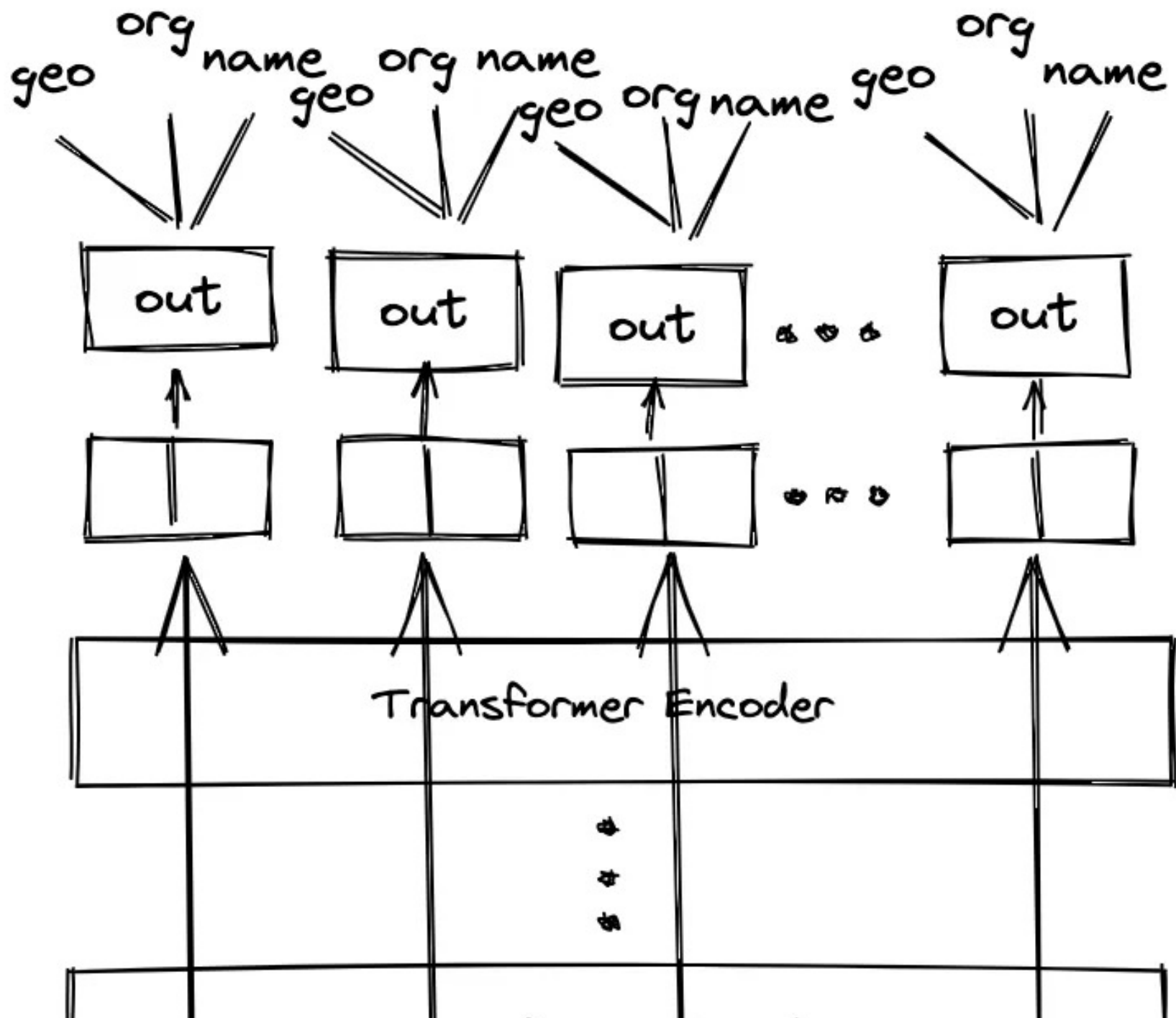
What differentiates between BERT for text classification and the NER problem is how we set the output of the model. For a text classification problem, we only use the embedding vector output from the special [CLS] token, as you can see in the visualization below:

geo org name



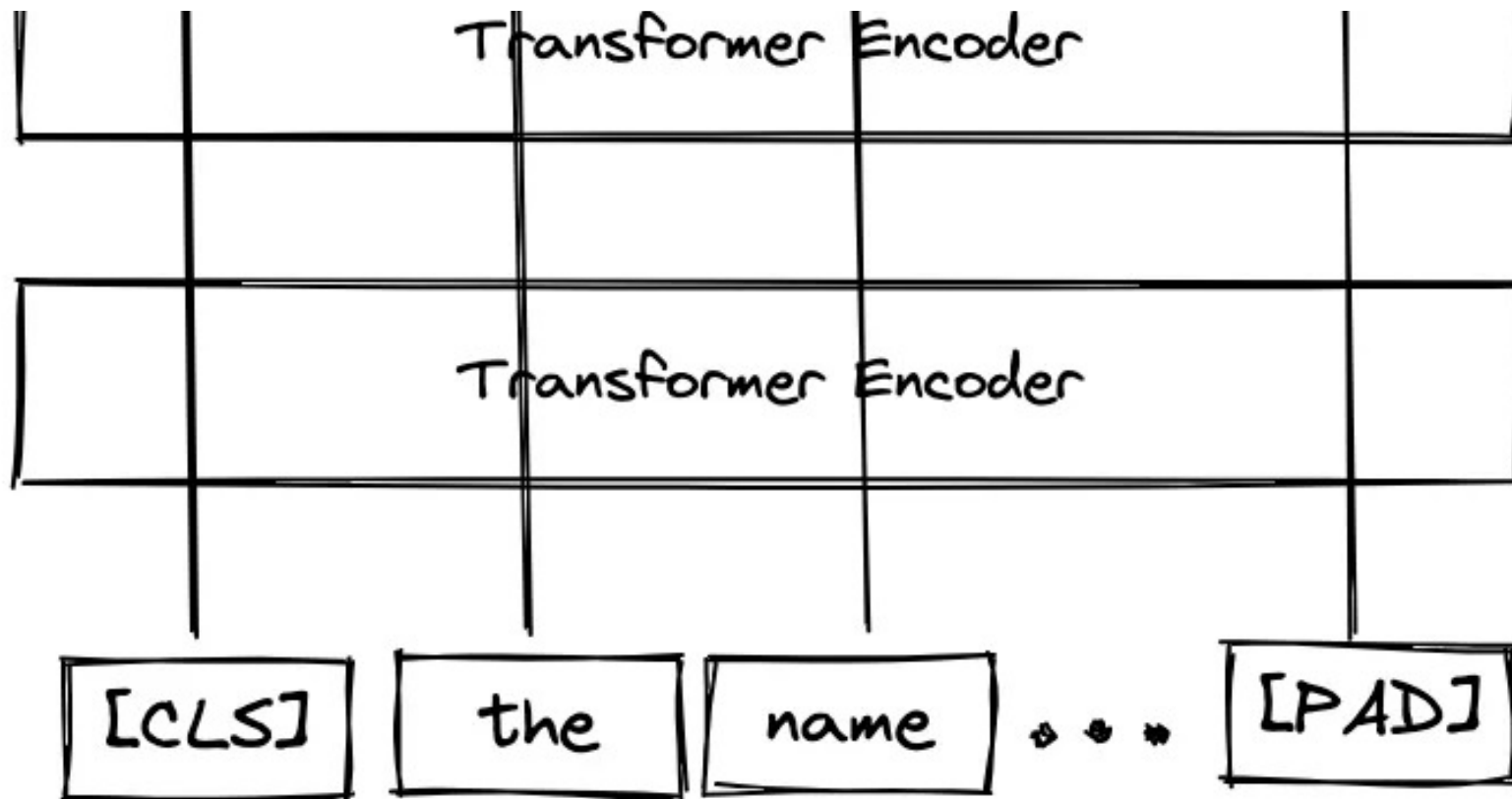
Mea
all c





By u
leve
Nov

n
en.



Ab
The
spe

NER Data

Named Entity Recognition Data

www.kaggle.com

This dataset is distributed under Open Database v1.0 license, so we are free to share and use this dataset for our own purpose. Now let's take a look at what the dataset looks like.

```
In [1]: import pandas as pd
```

```
In [2]: df = pd.read_csv('ner.csv')  
df.head()
```

```
Out[2]:
```

| | text | labels |
|---|---|---|
| 0 | Thousands of demonstrators have marched throug... | O O O O O O B-geo O O O O O B-geo O O O O O B-... |
| 1 | Iranian officials say they expect to get acces... | B-gpe O O O O O O O O O O O O O O O O B-tim O O O ... |
| 2 | Helicopter gunships Saturday pounded militant ... | O O B-tim O O O O O B-geo O O O O O B-org O O ... |
| 3 | They left after a tense hour-long standoff wit... | O O O O O O O O O O O O |
| 4 | U.N. relief coordinator Jan Egeland said Sunda... | B-geo O O B-per I-per O B-tim O B-geo O B-gpe ... |

```
In [ ]:
```

As we can see above, we have a dataframe which consists of the text and the label. The label corresponds to entity category of each word in a text.

In total, there are 9 entity categories, which are:

- `geo` for geographical entity
- `org` for organization entity
- `per` for person entity
- `gpe` for geopolitical entity
- `tim` for time indicator entity
- `art` for artifact entity
- `eve` for event entity
- `nat` for natural phenomenon entity
- `o` is assigned if a word doesn't belong to any entity.

Let's take a look at the unique labels available on our dataset:

```

1 # Split labels based on whitespace and turn them into a list
2 labels = [i.split() for i in df['labels'].values.tolist()]
3
4 # Check how many labels are there in the dataset
5 unique_labels = set()
6
7 for lb in labels:
8     [unique_labels.add(i) for i in lb if i not in unique_labels]
9
10 print(unique_labels)
11 >>> {'B-tim', 'B-art', 'I-art', 'O', 'I-gpe', 'I-per', 'I-nat', 'I-geo', 'B-eve', 'B-org', 'B-gpe', 'I-eve', 'B-per'}
12
13 # Map each label into its id representation and vice versa
14 labels_to_ids = {k: v for v, k in enumerate(sorted(unique_labels))}
15 ids_to_labels = {v: k for v, k in enumerate(sorted(unique_labels))}
16 print(labels_to_ids)
17 >>> {'B-art': 0, 'B-eve': 1, 'B-geo': 2, 'B-gpe': 3, 'B-nat': 4, 'B-org': 5, 'B-per': 6, 'B-tim': 7, 'I-art': 8, 'I-geo': 9, 'I-gpe': 10, 'I-nat': 11, 'I-per': 12, 'I-eve': 13, 'O': 14}

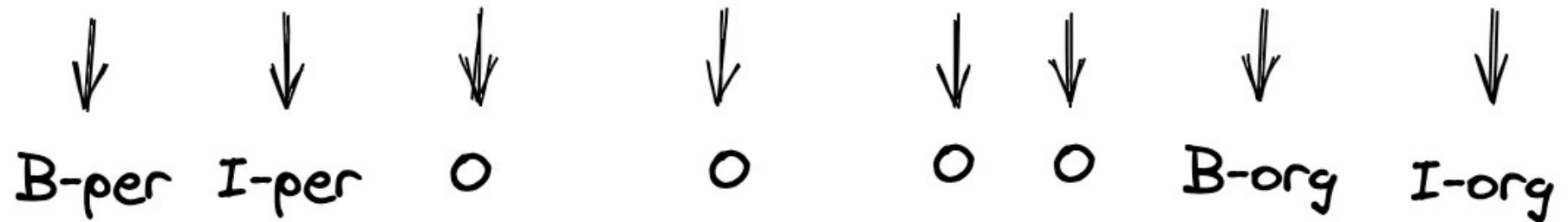
```

ner_1.py hosted with ❤ by GitHub

[view raw](#)

As you might notice, each entity category is preceded with the letter **I** or **B**. This corresponds to what previously mentioned as IOB tagging. **I** means *Intermediate* and **B** means *Beginning*. Let's take a look at the following sentence to understand the concept of IOB tagging a little bit more.

Kevin Durant plays basketball for the Brooklyn Nets



- 'Kevin' has B-pers label since it's the beginning of a person entity
- 'Durant' has I-pers label because it's the continuation of a person entity
- 'Brooklyn' has B-org since it's the beginning of an organization entity
- 'Nets' has I-org label since it's the continuation of an organization entity
- Other words are assigned O label as they don't belong to any entity

Data Preprocessing


Before we are able to use a BERT model to classify the entity of a token, of course, we need to do data preprocessing first, which includes two parts: tokenization and adjusting the label to match the tokenization. Let's start with tokenization first.

Tokenization

Tokenization can be easily implemented with BERT, as we can use `BertTokenizerFast` class from a pretrained BERT base model with HuggingFace.

To give you an example how BERT tokenizer works, let's take a look at one of the texts from our dataset:

```
1 # Let's take a look at how can we preprocess the text - Take first example
2 text = df['text'].values.tolist()
3 example = text[36]
4
5 print(example)
6 >>> 'Prime Minister Geir Haarde has refused to resign or call for early elections.'
```

ner_2.py hosted with  by GitHub

[view raw](#)

Tokenizing the text above with `BertTokenizerFast` is very straightforward:

```
1 from transformers import BertTokenizerFast
2
3 tokenizer = BertTokenizerFast.from_pretrained('bert-base-cased')
4 text_tokenized = tokenizer(example, padding='max_length', max_length=512, truncation=True, return_tensors="pt")
```

ner_3.py hosted with ❤️ by GitHub

[view raw](#)

We provide several arguments when calling `tokenizer` method from `BertTokenizerFast` class above:

- `padding` : to pad the sequence with a special [PAD] token to the maximum length that we specify. The maximum length of a sequence for a BERT model is 512.
- `max_length` : maximum length of a sequence.
- `truncation` : this is a Boolean value. If we set the value to `True`, then tokens that exceed the maximum length will not be used.
- `return_tensors` : the tensor type that is returned, depending on machine learning frameworks that we use. Since we're using PyTorch, then we use `pt` .

And below is the output of the tokenization process:


```
30         0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
31         0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
32         0, 0, 0, 0, 0, 0, 0, 0, 0]]})
```

ner_4.py hosted with ❤️ by GitHub

[view raw](#)

the special **[CLS]** token, the id 102 is reserved for the special **[SEP]** token, and the id 0 is reserved for **[PAD]** token.

- `token_type_ids` : To identify the sequence in which a token belongs to. Since we only have one sequence per text, then all the values of `token_type_ids` will be 0.
- `attention_mask` : To identify whether a token is a real token or padding. The value would be 1 if it's a real token, and 0 if it's a **[PAD]** token.

From the `input_ids` above, we can decode the ids back into the original sequence with `decode` method as follows:

```
1 print(tokenizer.decode(text_tokenized.input_ids[0]))
2
3 >>> '[CLS] Prime Minister Geir Haarde has refused to resign or call for early elections. [SEP] [PAD] [PAD] [PAD] [PA
```

ner_5.py hosted with ❤️ by GitHub

[view raw](#)

We got our original sequence back after implementing `decode` method with the addition of special tokens from BERT such as **[CLS]** token at the beginning of the sequence, **[SEP]** token at the end of the sequence, and a bunch of **[PAD]** tokens to fulfill the required maximum length of 512.

After this tokenization process, we need to proceed to the next step, which is adjusting the label of each token.

Adjusting Label After Tokenization

This is a very important step that we need to do after the tokenization process. This is because the length of the sequence is no longer matching the length of the original label after the tokenization process.

The BERT tokenizer uses the so-called word-piece tokenizer under the hood, which is a sub-word tokenizer. This means that BERT tokenizer will likely to split one word into one or more meaningful sub-words.

As an example, let's say we have the following sequence:

Prime Minister Geir Haarde has refused to resign or call for early elections

↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓

○ ○ B-per I-per ○ ○ ○ ○ ○ ○ ○ ○ ○

The sequence above has in total 13 tokens and thus, it also has 13 labels. However, after BERT tokenization, we get the following result:

```
1 print(tokenizer.convert_ids_to_tokens(text_tokenized["input_ids"][0]))
2
3 >>> ['[CLS]', 'Prime', 'Minister', 'G', '##ei', '##r', 'Ha', '##ard', '##e', 'has', 'refused', 'to', 'resign', 'or',
```

ner_6.py hosted with ❤ by GitHub [view raw](#)

There are two problems that we need to address after tokenization process:

- The addition of special tokens from BERT such as [CLS], [SEP], and [PAD]
- The fact that some tokens are splitted into sub-words.

As sub-word tokenization, word-piece tokenization splits uncommon words into their sub-words, such as 'Geir' and 'Haarde' in the example above. This sub-word tokenization helps the BERT model to learn the semantic meaning of related words.

The consequence of this word piece tokenization and the addition of special tokens from BERT is that the sequence length after tokenization is no longer matching the length of the initial label.

From the example above, now there are in total 512 tokens in the sequence after tokenization, while the length of the label is still the same as before. Also, the first token in a sequence is no longer the word *'Prime'*, but the newly added [CLS] token, so we need to shift our label as well.

To solve this problem, we need to adjust the label such that it has the same length as the sequence after tokenization. To do this, we can utilize the `word_ids` method from the tokenization result as follows:

```
1 word_ids = text_tokenized.word_ids()
2 print(tokenizer.convert_ids_to_tokens(text_tokenized["input_ids"][0]))
3 print(word_ids)
4
5 >>> ['[CLS]', 'Prime', 'Minister', 'G', '##ei', '##r', 'Ha', '##ard', '##e', 'has', 'refused', 'to', 'resign', 'or',
6 >>> [None, 0, 1, 2, 2, 2, 3, 3, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, None, None, None, None, None, None, None, None,
```

ner_7.py hosted with ❤ by GitHub

[view raw](#)

As you can see from the code snippet above, each splitted token shares the same `word_ids`, where special tokens from BERT such as [CLS], [SEP], and [PAD] all do not have specific `word_ids`.

These `word_ids` will be very useful to adjust the length of the label by applying either of these two methods:

1. We only provide a label to the first sub-word of each splitted token. The continuation of the sub-word then will simply have '-100' as a label. All tokens that don't have `word_ids` will also be labeled with '-100'.
2. We provide the same label among all of the sub-words that belong to the same token. All tokens that don't have `word_ids` will be labeled with '-100'.

The function in the code snippet below will do exactly the step defined above.

```
1  def align_label_example(tokenized_input, labels):
2
3      word_ids = tokenized_input.word_ids()
4
5      previous_word_idx = None
6      label_ids = []
7
8      for word_idx in word_ids:
9
10         if word_idx is None:
11             label_ids.append(-100)
12
13         elif word_idx != previous_word_idx:
14             try:
15                 label_ids.append(labels_to_ids[labels[word_idx]])
16             except:
17                 label_ids.append(-100)
18
19         else:
20             label_ids.append(labels_to_ids[labels[word_idx]] if label_all_tokens else -100)
21             previous_word_idx = word_idx
22
23
24     return label_ids
```

If you want to apply the first method, set `label_all_tokens` to `False`. If you want to apply the second method, set `label_all_tokens` to `True`, as you can see in the following code snippet:

```
1 label = labels[36]
2
3 #If we set label_all_tokens to True.....
4 label_all_tokens = True
5
6 new_label = align_label_example(text_tokenized, label)
7 print(new_label)
8 print(tokenizer.convert_ids_to_tokens(text_tokenized["input_ids"][0]))
9
10 >>> [-100, 16, 16, 6, 6, 6, 14, 14, 14, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16, -100, -100, -100, -100, -100, -100,
11 >>> ['[CLS]', 'Prime', 'Minister', 'G', '##ei', '##r', 'Ha', '##ard', '##e', 'has', 'refused', 'to', 'resign', 'or']
12
13
14 #If we set label_all_tokens to False.....
15 label_all_tokens = False
16
17 new_label = align_label_example(text_tokenized, label)
18 print(new_label)
19 print(tokenizer.convert_ids_to_tokens(text_tokenized["input_ids"][0]))
20
21 >>> [-100, 16, 16, 6, -100, -100, 14, -100, -100, 16, 16, 16, 16, 16, 16, 16, 16, 16, 16, -100, -100, -100, -100, -
22 >>> ['[CLS]', 'Prime', 'Minister', 'G', '##ei', '##r', 'Ha', '##ard', '##e', 'has', 'refused', 'to', 'resign', 'or']
```

In the rest of this article, we're going to implement the first method, in which we will only provide a label to the first sub-word in each token and set `label_all_tokens` to `False`.

Dataset Class

Before we train our BERT model for NER task, we need to create a dataset class to generate and fetch data in a batch.

```
1  import torch
2
3  def align_label(texts, labels):
4      tokenized_inputs = tokenizer(texts, padding='max_length', max_length=512, truncation=True)
5
6      word_ids = tokenized_inputs.word_ids()
7
8      previous_word_idx = None
9      label_ids = []
10
11     for word_idx in word_ids:
12
13         if word_idx is None:
14             label_ids.append(-100)
15
16         elif word_idx != previous_word_idx:
17             try:
18                 label_ids.append(labels_to_ids[labels[word_idx]])
19             except:
20                 label_ids.append(-100)
21         else:
22             try:
23                 label_ids.append(labels_to_ids[labels[word_idx]] if label_all_tokens else -100)
24             except:
25                 label_ids.append(-100)
26         previous_word_idx = word_idx
27
28     return label_ids
29
```

```
30 class DataSequence(torch.utils.data.Dataset):
31
32     def __init__(self, df):
33
34         lb = [i.split() for i in df['labels'].values.tolist()]
35         txt = df['text'].values.tolist()
36         self.texts = [tokenizer(str(i),
37                             padding='max_length', max_length = 512, truncation=True, return_tensors="pt") for i
38                       self.labels = [align_label(i,j) for i,j in zip(txt, lb)]
39
```

```
1 import numpy as np
2
3 df = df[0:1000]
4 df_train, df_val, df_test = np.split(df.sample(frac=1, random_state=42),
5                                       [int(.8 * len(df)), int(.9 * len(df))])
```

ner_11.py hosted with ❤ by GitHub

[view raw](#)

```
47
48     def get_batch_labels(self, idx):
49
50         return torch.LongTensor(self.labels[idx])
51
52     def __getitem__(self, idx):
53
54         batch_data = self.get_batch_data(idx)
55         batch_labels = self.get_batch_labels(idx)
56
57         return batch_data, batch_labels
```

BERT model that will act as token-level classifiers.

```
1  from transformers import BertForTokenClassification
2
3  class BertModel(torch.nn.Module):
4
5      def __init__(self):
6
7          super(BertModel, self).__init__()
8
9          self.bert = BertForTokenClassification.from_pretrained('bert-base-cased', num_labels=len(unique_labels))
10
11      def forward(self, input_id, mask, label):
12
13          output = self.bert(input_ids=input_id, attention_mask=mask, labels=label, return_dict=False)
14
15          return output
```

In the code snippet above, first, we instantiate the model and set the output of each token classifier equal to the number of unique entities on our dataset, which in our case is 17.

Next, we will define a function for the training loop.

Training Loop

The training loop for our BERT model is the standard PyTorch training loop with a few additions, as you can see below:

```
1  def train_loop(model, df_train, df_val):
2
3      train_dataset = DataSequence(df_train)
4      val_dataset = DataSequence(df_val)
5
6      train_dataloader = DataLoader(train_dataset, num_workers=4, batch_size=BATCH_SIZE, shuffle=True)
7      val_dataloader = DataLoader(val_dataset, num_workers=4, batch_size=BATCH_SIZE)
8
9      use_cuda = torch.cuda.is_available()
10     device = torch.device("cuda" if use_cuda else "cpu")
11
12     optimizer = SGD(model.parameters(), lr=LEARNING_RATE)
13
14     if use_cuda:
15         model = model.cuda()
16
17     best_acc = 0
18     best_loss = 1000
19
20     for epoch_num in range(EPOCHS):
21
22         total_acc_train = 0
23         total_loss_train = 0
24
25         model.train()
26
27         for train_data, train_label in tqdm(train_dataloader):
28
29             train_label = train_label.to(device)
```

```

30     mask = train_data['attention_mask'].squeeze(1).to(device)
31     input_id = train_data['input_ids'].squeeze(1).to(device)
32
33     optimizer.zero_grad()
34     loss, logits = model(input_id, mask, train_label)
35
36     for i in range(logits.shape[0]):
37
38         logits_clean = logits[i][train_label[i] != -100]
39         label_clean = train_label[i][train_label[i] != -100]
40

```

```

100%|██████████| 800/800 [02:41<00:00, 4.96it/s]
Epochs: 1 | Loss: 0.477 | Accuracy: 0.871 | Val_Loss: 0.362 | Accuracy: 0.895
100%|██████████| 800/800 [02:41<00:00, 4.96it/s]
Epochs: 2 | Loss: 0.343 | Accuracy: 0.903 | Val_Loss: 0.326 | Accuracy: 0.912
100%|██████████| 800/800 [02:41<00:00, 4.96it/s]
Epochs: 3 | Loss: 0.278 | Accuracy: 0.917 | Val_Loss: 0.290 | Accuracy: 0.915
100%|██████████| 800/800 [02:41<00:00, 4.96it/s]
Epochs: 4 | Loss: 0.216 | Accuracy: 0.932 | Val_Loss: 0.308 | Accuracy: 0.919
100%|██████████| 800/800 [02:41<00:00, 4.96it/s]
Epochs: 5 | Loss: 0.172 | Accuracy: 0.945 | Val_Loss: 0.285 | Accuracy: 0.926

```

```

50
51     total_acc_val = 0
52     total_loss_val = 0
53
54     for val_data, val_label in val_dataloader:
55
56         val_label = val_label.to(device)
57         mask = val_data['attention_mask'].squeeze(1).to(device)
58         input_id = val_data['input_ids'].squeeze(1).to(device)
59

```

```
59
60     loss, logits = model(input_id, mask, val_label)
61
62     for i in range(logits.shape[0]):
63
64         logits_clean = logits[i][val_label[i] != -100]
65         label_clean = val_label[i][val_label[i] != -100]
66
67         predictions = logits_clean.argmax(dim=1)
68         acc = (predictions == label_clean).float().mean()
69         total_acc_val += acc
70         total_loss_val += loss.item()
71
72     val_accuracy = total_acc_val / len(df_val)
73     val_loss = total_loss_val / len(df_val)
74
75     print(
76         f'Epochs: {epoch_num + 1} | Loss: {total_loss_train / len(df_train): .3f} | Accuracy: {total_acc_train
77
78     LEARNING_RATE = 5e-3
79     EPOCHS = 5
80     BATCH_SIZE = 2
81
82     model = BertModel()
83     train_loop(model, df_train, df_val)
```

```
1  def evaluate(model, df_test):
2
3      test_dataset = DataSequence(df_test)
4
5      test_dataloader = DataLoader(test_dataset, num_workers=4, batch_size=1)
6
7      use_cuda = torch.cuda.is_available()
8      device = torch.device("cuda" if use_cuda else "cpu")
9
10     if use_cuda:
11         model = model.cuda()
12
13     total_acc_test = 0.0
14
15     for test_data, test_label in test_dataloader:
16
17         test_label = test_label.to(device)
18         mask = test_data['attention_mask'].squeeze(1).to(device)
19
20         input_id = test_data['input_ids'].squeeze(1).to(device)
21
22         loss, logits = model(input_id, mask, test_label)
23
24         for i in range(logits.shape[0]):
25
26             logits_clean = logits[i][test_label[i] != -100]
27             label_clean = test_label[i][test_label[i] != -100]
28
29             predictions = logits_clean.argmax(dim=1)
```

```
30         acc = (predictions == label_clean).float().mean()
31         total_acc_test += acc
32
33     val_accuracy = total_acc_test / len(df_test)
34     print(f'Test Accuracy: {total_acc_test / len(df_test): .3f}')
35
36
37 evaluate(model, df_test)
```

ner_14.py hosted with  by GitHub

[view raw](#)

```
1  def align_word_ids(texts):
2
3      tokenized_inputs = tokenizer(texts, padding='max_length', max_length=512, truncation=True)
4
5      word_ids = tokenized_inputs.word_ids()
6
7      previous_word_idx = None
8      label_ids = []
9
10     for word_idx in word_ids:
11
12         if word_idx is None:
13             label_ids.append(-100)
14
15         elif word_idx != previous_word_idx:
16             try:
17                 label_ids.append(1)
18             except:
19                 label_ids.append(-100)
20         else:
21             try:
22                 label_ids.append(1 if label_all_tokens else -100)
23             except:
24                 label_ids.append(-100)
25         previous_word_idx = word_idx
26
27     return label_ids
28
29
```

```
30 def evaluate_one_text(model, sentence):
31
32
33     use_cuda = torch.cuda.is_available()
34     device = torch.device("cuda" if use_cuda else "cpu")
35
36     if use_cuda:
37         model = model.cuda()
38
39     text = tokenizer(sentence, padding='max_length', max_length = 512, truncation=True, return_tensors="pt")
40
41     mask = text['attention_mask'].to(device)
42     input_id = text['input_ids'].to(device)
43     label_ids = torch.Tensor(align_word_ids(sentence)).unsqueeze(0).to(device)
44
45     logits = model(input_id, mask, None)
46     logits_clean = logits[0][label_ids != -100]
47
48     predictions = logits_clean.argmax(dim=1).tolist()
49     prediction_label = [ids_to_labels[i] for i in predictions]
50     print(sentence)
51     print(prediction_label)
52
53 evaluate_one_text(model, 'Bill Gates is the founder of Microsoft')
54
55 >>> 'Bill Gates is the founder of Microsoft'
56 >>> ['B-per', 'I-per', 'O', 'O', 'O', 'O', 'B-org']
```

