

① Get Dataset & build tokenizer

1.1 : raw dataset from hugging face.

Huggingface

from datasets import load_dataset

raw_dataset = load_dataset({'config': "data_source"}, {'config': "lang_src"} - {'config': "lang_tgt"}, split='train')

ex) load_dataset('opus_basw', "en-fr", split='train')

↳ Dataset ↳ Subset

1.2 build tokenizer (need tokenizers for encoder and decoder)

★ in our case, encoder-tokenizer is tokenizing English sentence

decoder-tokenizer is tokenizing French sentence

★ Tokenizer Pipeline: Normalization → model → Pre-tokenization → train → Post-processing (optional)

- Pre-tokenization is the act of splitting a text into smaller objects. A good way to think of pre-tokenization is that it will split your text into words.

Huggingface

from tokenizers import Tokenizer

from tokenizers.models import WordLevel

from tokenizers.trainers import WordLevelTrainer

from tokenizers.pre_tokenizers import Whitespace

tokenizer_src = Tokenizer(^{model} WordLevel(^{no unknown tokens} unk_token='UNK'))

tokenizer_src.pre_tokenizer = Whitespace() # Split text into words using white spaces.

trainer = WordLevelTrainer(special_tokens=["[UNK]", "[PAD]", "[SOS]", "[EOS]"], min_count=2)

tokenizer_src.train_from_iterator(iterator=get_all_sentences(data, "en"), trainer=trainer)

↳ def get_all_sentences(data, lang):
 for item in data:
 yield item["translation"][lang]

tokenizer_tgt = Tokenizer(^{model} WordLevel(^{unk_token='UNK'})))

tokenizer_tgt.pre_tokenizer = Whitespace() # Split text into words using white spaces.

trainer = WordLevelTrainer(special_tokens=["[UNK]", "[PAD]", "[SOS]", "[EOS]"], min_count=2)

tokenizer_tgt.train_from_iterator(iterator=get_all_sentences(data, "fr"), trainer=trainer)

1.3 Split train & test

train_size = int(0.9 * len(raw_dataset))

test_size = len(raw_dataset) - train_size

train_val = random_split(raw_dataset, [train_size, val_size])

↳ torch.utils.data.random_split

1.4 Create a custom dataset by using Tokenizer (torch.utils.data.Dataset)

A custom Dataset class must implement three functions: `__init__`, `__len__`, and `__getitem__`.

class BilingualDataset(Dataset):

```
def __init__(self, data, tokenizer_src, tokenizer_tgt, src_lang, tgt_lang, seq_len):
    super().__init__()
    self.data = data
    self.tokenizer_src = tokenizer_src
    self.tokenizer_tgt = tokenizer_tgt
```

```
self.src_tokens = torch.tensor([tokenizer_src.tokenize(s) for s in data], dtype=torch.int64)
self.egs_src = self.src_tokens
self.pads_src = torch.zeros((1, seq_len), dtype=torch.int64)
```

```
def __getitem__(self, index):
    return torch.cat([self.src_tokens[index], self.pads_src], dim=1)
```

```
def __len__(self):
    return len(self.src_tokens)

src_text = pairs['translation'][self.src_lang]
tgt_text = pairs['translation'][self.tgt_lang]
```

src_text → src_tokenizer → Add special tokens → encoder_input

tgt_text → tgt_tokenizer → Add special tokens → decoder_input

encoder_mask = (encoder_input != self.pad_token).unsqueeze(0).int() # size: (B, seq_len)

```
[[[1, 1, 1, 1, 1, 1, 0, 0, ..., 0, 0, 0],
  [1, 1, 1, 1, 1, 1, 1, 1, ..., 0, 0, 0],
  ...
  [1, 1, 1, 1, 0, 0, 0, 0, ..., 0, 0, 0]]]
```

Batch size

seq_len = 350

decoder_mask = (decoder_input != self.pad_token).unsqueeze(0).int() & self.causal_mask(decoder_input.size(0)) # (Batch, seq_len, seq_len)

```
[[[1, 0, 0, 0, ..., 0, 0, 0],
  [1, 1, 0, 0, ..., 0, 0, 0],
  [1, 1, 1, 0, ..., 0, 0, 0],
  [1, 1, 1, 1, ..., 0, 0, 0],
  ...
  [1, 1, 1, 1, ..., 0, 0, 0]]]
```

seq_len

Batch size

```
[[[1, 0, 0, 0, ..., 0, 0, 0],
  [1, 1, 0, 0, ..., 0, 0, 0],
  [1, 1, 1, 0, ..., 0, 0, 0],
  [1, 1, 1, 1, ..., 0, 0, 0],
  ...
  [1, 1, 1, 1, ..., 0, 0, 0]]]
```

seq_len

Refer

② Build Transformer.

Pipeline : (let's say $N=2$: # of blocks)

