# 모델 학습을 위한 데이터 전처리

- 데이터 로드

- 토큰화

```python
31    def load_dataset(self, dataset_id, tokenizer):
32        def tokenization(examples):
33            sources = []
34            targets = []
35            prompt = self.PROMPT_TEMPLATE
36            for instruction, input, output in zip(examples['instruction'],examples['input'],examples['output']):
37                if input is not None and input !="":
38                    instruction = instruction+'\n'+input
39                source = prompt.format_map({'instruction':instruction})
40                target = f"{output}{tokenizer.eos_token}"
41
42                sources.append(source)
43                targets.append(target)
44
45            tokenized_sources = tokenizer(sources,return_attention_mask=False)
46            tokenized_targets = tokenizer(targets,return_attention_mask=False,add_special_tokens=False)
47
48            all_input_ids = []
49            all_labels = []
50            for s,t in zip(tokenized_sources['input_ids'],tokenized_targets['input_ids']):
51                input_ids = torch.LongTensor(s + t)[:self.MAX_SEQ_LEN]
52                labels = torch.LongTensor([-100] * len(s) + t)[:self.MAX_SEQ_LEN]
53                assert len(input_ids) == len(labels)
54                all_input_ids.append(input_ids)
55                all_labels.append(labels)
56
57            results = {'input_ids':all_input_ids, 'labels': all_labels}
58            return results
59
60        all_datasets = []
61
62        raw_dataset = load_dataset(dataset_id)        # Hugging Face Hub Model load
63        tokenization_func = tokenization
```

```python
31    def load_dataset(self, dataset_id, tokenizer):
32        def tokenization(examples):
33            sources = []
34            targets = []
35            prompt = self.PROMPT_TEMPLATE
36            for instruction, input, output in zip(examples['instruction'],examples['input'],examples['output']):
37                if input is not None and input !="":
38                    instruction = instruction+'\n'+input
39                source = prompt.format_map({'instruction':instruction})
40                target = f"{output}{tokenizer.eos_token}"
41
42                sources.append(source)
43                targets.append(target)
44
45            tokenized_sources = tokenizer(sources,return_attention_mask=False)
46            tokenized_targets = tokenizer(targets,return_attention_mask=False,add_special_tokens=False)
47                                                          # Tokenizer Text to Number ID
48            all_input_ids = []
49            all_labels = []
50            for s,t in zip(tokenized_sources['input_ids'],tokenized_targets['input_ids']):
51                input_ids = torch.LongTensor(s + t)[:self.MAX_SEQ_LEN]
52                labels = torch.LongTensor([-100] * len(s) + t)[:self.MAX_SEQ_LEN]
53                assert len(input_ids) == len(labels)
54                all_input_ids.append(input_ids)
55                all_labels.append(labels)
56
57            results = {'input_ids':all_input_ids, 'labels': all_labels}
58            return results
59
60        all_datasets = []
61
62        raw_dataset = load_dataset(dataset_id)
63        tokenization_func = tokenization
```

GitHub 주소: https://github.com/EDDI-RobotAcademy/SI-Follow-LLAMA/tree/main/dataset

고용노동부   KOREATECH 직업능력심사평가원

# 모델 학습을 위한 데이터 전처리

- 배치 및 Collator

- 자동화 생성

```python
from dataclasses import dataclass
from typing import Dict, Sequence

import torch
import transformers


@dataclass
class DataCollatorForSupervisedDataset:
    tokenizer: transformers.PreTrainedTokenizer

    def __call__(self, instances: Sequence[Dict]) -> Dict[str, torch.Tensor]:
        input_ids, labels = tuple(
            [instance[key] for instance in instances] for key in ("input_ids", "labels")
        )
        input_ids = torch.nn.utils.rnn.pad_sequence(
            input_ids, batch_first=True, padding_value=self.tokenizer.pad_token_id
        )
        labels = torch.nn.utils.rnn.pad_sequence(
            labels, batch_first=True, padding_value=-100
        )
        return dict(
            input_ids=input_ids,
            labels=labels,
            attention_mask=input_ids.ne(self.tokenizer.pad_token_id),
        )
```

패딩을 통해 배치의 데이터 길이 최적화

```python
raw_dataset = load_dataset(dataset_id)
tokenization_func = tokenization
tokenized_dataset = raw_dataset.map(
    tokenization_func,
    batched=True,
    remove_columns=["instruction","input","output"],
    keep_in_memory=False,
    desc="preprocessing on dataset",
)
processed_dataset = tokenized_dataset
processed_dataset.set_format('torch')
all_datasets.append(processed_dataset['train'])
all_datasets = concatenate_datasets(all_datasets)
return all_datasets

def get_data_collator(self, tokenizer):
    collator = DataCollatorForSupervisedDataset(tokenizer)
    return collator
```

Tokenizing DataSet

GitHub 주소: https://github.com/EDDI-RobotAcademy/SI-Follow-LLAMA/tree/main/dataset

고용노동부    KOREATECH 직업능력심사평가원