

9 주차 과제



| | |
|-------|------------|
| 학번 | 2020136087 |
| 이름 | 윤아현 |
| 수강 과목 | 자연어처리 |
| 담당 교수 | 임상훈 교수님 |
| 제출일 | 2024.05.16 |

문제 및 해결 방법 (구현 코드)

1. IMDB PLM 학습 (45 점)

- * HW 2 에서 진행한 IMDB 데이터셋을 학습 및 분류하는 PLM 기반 분류기를 작성하시오.
- * 실습에 사용하지 않은 PLM 모델을 huggingface hub 에서 불러와 사용하시오.
- * 실습에 사용한 데이터는 한국어 데이터이지만 IMDB 의 경우 영어 데이터임
- * 최소 3 개의 PLM 을 동일 hyperparameter 를 통해 학습 및 평가 해 보고 성능을 비교 분석하시오
- * 각자가 구현 및 제출한 HW2 의 모델과 비교하여 어떠한지
- * 각 PLM 별로 성능이 어떠한고 왜 그런지

GRADING

- * PLM 1 개당 15 점 (총 45 점)

Code

```
1 # plm_name = "distilbert-base-uncased"
2 # plm_name = "albert-base-v2"
3 plm_name = "xlnet-base-cased"

1 # 알아서 Classification 수행해줌
2 # sentiment_classifier 모델을 직접 안짜도 알아서 다 해줌
3 from transformers import AutoTokenizer, AutoModel
4 from transformers import AutoModelForSequenceClassification
5
6 tokenizer = AutoTokenizer.from_pretrained(plm_name)
7 # label, 대상을 알려준다.
8 model = AutoModelForSequenceClassification.from_pretrained(plm_name, num_labels=2)

1 from torch.utils.data import Dataset, DataLoader
2
3 # define dataset class
4 class SentimentDataset(Dataset):
5     def __init__(self, data, tokenizer):
6         self.data = data
7         self.tokenizer= tokenizer
8
9     def __len__(self):
10         return len(self.data)
11
12     def __getitem__(self, index):
13         label = int(self.data[index][2])
14         text = self.data[index][1]
15         # tokenizer가 자동으로 수행함
16         inputs = self.tokenizer(text, return_tensors="pt", padding="max_length",
17                                 truncation=True, max_length=128)
18         inputs = {key: inputs[key].squeeze() for key in inputs}
19
20         # inputs['label'] = torch.tensor(label)
21
22         return inputs|{'label':torch.tensor(label)} # tokenizer에 label만 추가해주면 됨
```

```

1 import lightning as pl
2
3 class SentimentClassifierPL(pl.LightningModule):
4     def __init__(self, sentiment_classifier):
5         super(SentimentClassifierPL, self).__init__()
6         self.model = sentiment_classifier
7
8         self.validation_step_outputs = []
9         self.test_step_outputs = []
10        self.save_hyperparameters()
11
12    def forward(self, inputs):
13        return self.model(inputs)
14
15    def training_step(self, batch, batch_idx):
16        labels = batch.pop('label')
17        outputs = model(**batch, labels=labels)
18        loss = outputs.loss # 자동으로 loss 계산해주기 때문에 바로 가져와서 사용
19        logits = outputs.logits # 자동으로 logits 계산해주기 때문에 바로 가져와서 사용
20        self.log("train_loss", loss)
21        return loss
22
23    def validation_step(self, batch, batch_idx):
24        labels = batch.pop('label')
25        outputs = model(**batch, labels=labels)
26        loss = outputs.loss
27        logits = outputs.logits
28        self.log("val_loss", loss)
29        self.validation_step_outputs.append((loss, logits, labels))
30        return loss, outputs, labels
31
32    def on_validation_epoch_end(self):
33        outputs = self.validation_step_outputs
34        avg_loss = torch.stack([x[0] for x in outputs]).mean()
35
36        self.log("avg_val_loss", avg_loss)
37
38        all_outputs = torch.cat([x[1] for x in outputs])
39        all_labels = torch.cat([x[2] for x in outputs])
40        all_preds = all_outputs.argmax(dim=1)
41        accuracy = (all_preds == all_labels).float().mean()
42        self.log("val_accuracy", accuracy)
43        self.validation_step_outputs.clear()
44
45    def test_step(self, batch, batch_idx):
46        labels = batch.pop('label')
47        outputs = model(**batch, labels=labels)
48        loss = outputs.loss
49        logits = outputs.logits
50        self.log("test_loss", loss)
51        self.test_step_outputs.append((loss, logits, labels))
52        return loss, outputs, labels
53
54    def on_test_epoch_end(self):
55        outputs = self.test_step_outputs
56        avg_loss = torch.stack([x[0] for x in outputs]).mean()
57        self.log("avg_test_loss", avg_loss)
58
59        all_outputs = torch.cat([x[1] for x in outputs])
60        all_labels = torch.cat([x[2] for x in outputs])
61        all_preds = all_outputs.argmax(dim=1)
62        accuracy = (all_preds == all_labels).float().mean()
63        self.log("test_accuracy", accuracy)
64        self.test_step_outputs.clear()
65
66    def configure_optimizers(self):
67        optimizer = torch.optim.AdamW(self.model.parameters(), lr=5e-6)
68        return optimizer

```

```

5 wandb.login()
6
7 def check_performance(model,tokenizer, train_data, val_data, test_data, wandb_log_name):
8     wandb_logger = WandbLogger(project="NLP", name=wandb_log_name, group="Lec07")
9
10    pl_model = SentimentClassifierPL(model)
11
12    train_dataset = SentimentDataset(train_data,tokenizer)
13    train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True)
14    val_dataset = SentimentDataset(val_data,tokenizer)
15    val_loader = DataLoader(val_dataset, batch_size=32, shuffle=False)
16    test_dataset = SentimentDataset(test_data,tokenizer)
17    test_loader = DataLoader(test_dataset, batch_size=32, shuffle=False)
18
19    trainer = pl.Trainer(
20        max_epochs=1,
21        accelerator="gpu",
22        logger=wandb_logger,
23        callbacks=[ModelSummary(max_depth=2)],
24        precision=16
25    )
26
27    trainer.fit(
28        model=pl_model,
29        train_dataloaders=train_loader,
30        val_dataloaders=val_loader
31    )
32
33    trainer.test(dataloaders=test_loader)
34
35    wandb.finish()

```

```

1 pl_model = SentimentClassifierPL(model)

```

```

1 check_performance(pl_model, tokenizer, train_data, val_data, test_data, "distilbert")

```

```

1 check_performance(pl_model, tokenizer, train_data, val_data, test_data, "albert")

```

```

1 check_performance(pl_model, tokenizer, train_data, val_data, test_data, "xlnet")

```

Outputs

1 번. DistilBERT

Google BERT 의 Model 에서 학습 시간과 모데 크기를 줄인 BERT 의 경량화 버전

성능: 64.67%

| Test metric | DataLoader 0 |
|---------------|--------------------|
| avg_test_loss | 0.6585080027580261 |
| test_accuracy | 0.6466667056083679 |
| test_loss | 0.6591495871543884 |

2 번. ALBERT

Cross-Layer 파라미터 공유와 같은 기법을 통해 모델 크기를 줄인 BERT 모델

성능: 52.33%

| Test metric | DataLoader 0 |
|---------------|--------------------|
| avg_test_loss | 0.688616931438446 |
| test_accuracy | 0.5233333706855774 |
| test_loss | 0.680382490158081 |

3 번. XLNet

Transformer 을 기반으로 하여 문맥적 이해가 중요한 task 에 적합한 모델

성능: 57.67%

| Test metric | DataLoader 0 |
|---------------|--------------------|
| avg_test_loss | 0.672248363494873 |
| test_accuracy | 0.5766666531562805 |
| test_loss | 0.6722908616065979 |

※기존의 HW2 모델과 다른 점

기존 설계 모델의 성능: 51.33%

당연하게도, 직접 설계한 모델보다는 3 개의 사전 훈련 모델을 가져와 사용하는 것이 더 높은 성능을 보였다. 기존의 설계 모델은 데이터셋에만 있는 단어와 단순한 선형망을 통해 훈련을 수행하기 때문에 문맥에 대한 정확한 이해가 어렵다.

사전 훈련 모델은 대규모 언어들로 학습을 시켰기 때문에 기존 모델보다는 훨씬 더 많은 단어에 대해 학습을 수행하였으며, Transformer 과 같은 모델을 사용하였기에 문맥 정보 또한 더 잘 파악할 수 있다는 장점이 있다.

※3 개의 사전 학습 모델 비교

XLNet 이 자연어 처리(복잡한 문맥 파악)에 더 적합하여 성능이 가장 높을 것이라고 가정하였다. 하지만, 예상과 달리 DistilBERT 모델이 더 좋은 성능을 낸 것을 확인할 수 있다.

ALBERT 의 경우, 각 계층에서 같은 파라미터를 공유하기에 각 계층간 다양한 특징을 추출하기에는 한계가 있을 것이라고 생각하였다.

DistilBERT 는 기존 BERT 의 경량화 버전이기에, BERT 보다 더 간단해졌지만 유사한 방식으로 데이터를 처리하기에 성능이 높은 것이라고 추측한다. 특히 BERT 는 Transformer 의 Encoder 부분만을

사용하기에 기존 문장들의 문맥 정보를 잘 파악하여 감정 분류에는 더 좋은 성능을 보인다고 생각하였다.

XLNet 또한, 전체적인 문맥을 파악하는 데에 높은 성능을 보이지만, IMDB 리뷰 데이터셋에 정밀한 파라미터 튜닝이 수행되지 않아 성능이 낮게 나왔다고 생각하였다.

WANDB 결과



2. Hyperparameter 조정 (15 점)

* 1 번 과제의 가장 성능이 좋은 모델을 통해 hyperparameter 변화에 따른 성능을 비교하시오.

* Batch size, learning rate 에 대한 성능 비교 필수

1 번 실험

Epoch: 10

Batch_size: 16

Learning_rate: 2e-5

성능: 70.67%

2 번 실험

Epoch: 10

Batch_size: 16

Learning_rate: 1e-4

성능: 71.00%

3 번 실험

Epoch: 10

Batch_size: 32

Learning_rate: 2e-5

성능: 69.67%

4 번 실험

Epoch: 10

Batch_size: 64

Learning_rate: 1e-4

성능: 68.99%

※ weight_decay = 0.01, eps=1e-8 고정

1 번. 학습률 영향

배치 사이즈를 16 으로 고정해놓고 학습률만 수정하여 훈련시켰을 경우, 더 높은 학습률이 조금 더 나은 결과를 가져왔지만 엄청난 영향을 미치지 않은 것을 확인할 수 있다.

2 번. 배치 사이즈 영향

배치 사이즈를 16, 32, 64 로 다양하게 훈련시켰을 경우, 배치 사이즈가 가장 낮은 16 일 때 성능이 제일 좋았던 것을 확인할 수 있다.

배치 사이즈를 늘리면, 일반화된 학습은 가능하지만 각 문장들이 가지는 구체적인 문맥 정보들을 파악하는 것이 어렵다고 추측된다. (local minimum 문제)

3. T5 를 통한 Machine Translation 모델 학습 (40 점)

- * Huggingface hub 의 T5 모듈 및 weight 을 통해 Seq2Seq 실습에 사용한 English-French MT 모델을 구현하시오.

- * Seq2Seq 실습에 사용한 데이터셋 사용 (<http://www.manythings.org/anki/fra-eng.zip>)

- * T5 모델은 "google-t5/t5-small"를 사용

- * 참고:

https://huggingface.co/docs/transformers/model_doc/t5#transformers.T5ForConditionalGeneration

- * 실습에 사용한 모델과의 성능 비교

GRADING

- * 학습 및 평가 (40 점)

Code

```
1 from transformers import AutoTokenizer, T5ForConditionalGeneration
2
3 plm_model = "google-t5/t5-small"
4
5 tokenizer = AutoTokenizer.from_pretrained(plm_model)
6 model = T5ForConditionalGeneration.from_pretrained(plm_model)
```



```

3 class TranslationDataset(Dataset):
4     def __init__(self, data, tokenizer, max_length=30):
5         self.data = data
6         self.tokenizer = tokenizer
7         self.input_texts = [item[0] for item in data] # 영어 문장
8         self.target_texts = [item[1] for item in data] # 프랑스어 문장
9         self.max_length = max_length
10
11     # Pre-tokenize all inputs and targets
12     self.encoder_inputs = self.tokenizer(
13         [f"translate English to French: {text}" for text in self.input_texts],
14         max_length=max_length,
15         padding="max_length",
16         truncation=True,
17         return_tensors="pt"
18     )
19
20     self.decoder_inputs = self.tokenizer(
21         self.target_texts,
22         max_length=max_length,
23         padding="max_length",
24         truncation=True,
25         return_tensors="pt"
26     )
27
28     # Replace padding token id's in labels by -100 so they are ignored by the loss
29     self.decoder_inputs['input_ids'][self.decoder_inputs['input_ids'] == self.tokenizer.pad_token_id] = -100
30
31     def __len__(self):
32         return len(self.input_texts)
33
34     def __getitem__(self, idx):
35         encoder_input = {key: val[idx] for key, val in self.encoder_inputs.items()}
36         labels = self.decoder_inputs['input_ids'][idx]
37         return encoder_input, labels

```

```

1 import lightning as pl
2
3 class TranslateModelPL(pl.LightningModule):
4     def __init__(self, translate_model):
5         super(TranslateModelPL, self).__init__()
6         self.model = translate_model
7         self.loss = nn.CrossEntropyLoss()
8
9         self.validation_step_outputs = []
10        self.test_step_outputs = []
11        self.save_hyperparameters()
12
13    def forward(self, inputs):
14        return self.model(inputs)
15
16    def training_step(self, batch, batch_idx):
17        encoder_input, labels = batch
18        outputs = model(**encoder_input, labels=labels)
19        loss = outputs.loss
20        self.log("train_loss", loss, prog_bar=True)
21        return loss
22
23    def validation_step(self, batch, batch_idx):
24        encoder_input, labels = batch
25        outputs = model(**encoder_input, labels=labels)
26        loss = outputs.loss
27        self.log("val_loss", loss)
28        self.validation_step_outputs.append((loss, outputs, labels))
29        return loss, outputs, labels
30
31    def on_validation_epoch_end(self):
32        outputs = self.validation_step_outputs
33        avg_loss = torch.stack([x[0] for x in outputs]).mean()
34        self.log("avg_val_loss", avg_loss)
35

```

```

36     self.validation_step_outputs.clear()
37
38     def test_step(self, batch, batch_idx):
39         encoder_input, labels = batch
40         outputs = model(**encoder_input, labels=labels)
41         loss = outputs.loss
42         self.log("test_loss", loss)
43         self.test_step_outputs.append((loss, outputs, labels))
44         return loss, outputs, labels
45
46     def on_test_epoch_end(self):
47         outputs = self.test_step_outputs
48         avg_loss = torch.stack([x[0] for x in outputs]).mean()
49         self.log("avg_test_loss", avg_loss)
50
51         self.test_step_outputs.clear()
52
53     def configure_optimizers(self):
54         optimizer = torch.optim.AdamW(self.model.parameters(), lr=5e-6)
55         return optimizer

```

```

1 import wandb
2 from lightning.pytorch.loggers import WandbLogger
3 from lightning.pytorch.callbacks import ModelSummary
4
5 wandb.login()
6
7 def check_performance(model, tokenizer, train_data, val_data, test_data, wandb_log_name):
8     wandb_logger = WandbLogger(project="NLP", name=wandb_log_name, group="Lec07")
9
10    pl_model = TranslateModelPL(model)
11
12    train_dataset = TranslationDataset(train_data, tokenizer)
13    train_loader = DataLoader(train_dataset, batch_size=16, shuffle=True)
14    val_dataset = TranslationDataset(val_data, tokenizer)
15    val_loader = DataLoader(val_dataset, batch_size=16, shuffle=False)
16    test_dataset = TranslationDataset(test_data, tokenizer)
17    test_loader = DataLoader(test_dataset, batch_size=16, shuffle=False)
18
19    trainer = pl.Trainer(
20        max_epochs=1,
21        accelerator="gpu",
22        logger=wandb_logger,
23        callbacks=[ModelSummary(max_depth=2)],
24        precision=16
25    )
26
27    trainer.fit(
28        model=pl_model,
29        train_data_loaders=train_loader,
30        val_data_loaders=val_loader
31    )
32
33    trainer.test(data_loaders=test_loader)
34
35    wandb.finish()

```

```

1 pl_model = TranslateModelPL(model)

```

```

1 check_performance(pl_model, tokenizer, train, vali, test, "t5")

1 input_sentence = "Are we leaving soon?"
2
3 task_prefix = "translate English to French: "
4 input_text = task_prefix + input_sentence
5
6 inputs = tokenizer(input_text, return_tensors="pt")
7
8 outputs = model.generate(**inputs, max_length=30, num_beams=4, early_stopping=True)
9
10 translation = tokenizer.decode(outputs[0], skip_special_tokens=True)
11
12 print(f"Input: {input_sentence}")
13 print("Output: y allons nous bientôt ?")
14 print(f"Translation: {translation}")

```

주의사항

Out Of Memory Issue 로 GPU 가 안돌아가서, Train, Vali, Test 데이터 개수를 대폭 줄여서 학습을 진행하였다.

Outputs

```

Input: Are we leaving soon?
Output: y allons nous bientôt ?
Translation: sortons nous bientôt?

```

위의 결과는 마지막 코드의 출력 값이다.

이전의 Seq2Seq 모델을 훈련시켰을 때 제대로 된 결과가 나오지 않았다. (동일한 단어가 2 번 이상 나오거나 문장이 너무 짧거나 긴 이슈)

하지만, T5 모델로 훈련하였을 때 Epoch 를 1 번만 실행시켜본 것 치고는 엄청난 결과가 나왔다.

앞의 단어 2 개를 제외한 뒤의 단어를 모두 맞힌 것을 확인할 수 있다.

사전 훈련된 T5 모델은 기존의 Seq2Seq 의 모델과 달리 양방향 Encoder 를 사용한다. (Decoder 는 어쩔 수 없이 순차적으로 진행이 필수임) 또한, T5 모델의 특징점에 대해 살펴보니 모델이 다른 문장들(Input 값)에서 획득한 언어 패턴과 문맥 정보를 다른 문장에 사용하여 여러 문장에서의 문맥 정보를 더 풍부하게 획득할 수 있다는 장점이 있다.

이를 통해, 적절한 문맥과 뉘앙스 같은 것을 파악할 수 있다.

WANDB

Charts 2

