



Practice!



실습

- Transformer를 이용하여 1문 1답 챗봇 프로그램을 작성하시오.

- 입력 데이터셋
 - 1문 1답 텍스트: 질문 w/ 답
(음절로 분리, 공백은 <SP>로 변환)
 - 어휘 사전: 음절 집합
- 문제
 - 1문 1답 텍스트를 학습하여 새로운 질문에 대한 적절한 응답을 생성

train.txt - Windows 메모장

파일(F) 편집(E) 서식(O) 보기(V) 도움말(H)

뭐 하고 있 어 ? 이 제 <SP> 수 업 <SP> 끝 냈 어 ! </S>
년 <SP> 뭐 해 대 학 생 이 야 ? </S>
난 <SP> 지 금 <SP> 과 자 <SP> 먹 어 응 <SP> 대 학 생 </S>
응 <SP> 대 학 생 전 공 이 <SP> 뭐 야 ? </S>
전 공 이 <SP> 뭐 야 ? 년 <SP> 집 인 가 보 네 <SP> 부 럽 다 <SP> 나
난 <SP> 국 어 국 문 학 과 오 <SP> 좋 은 <SP> 학 문 <SP> 배 우 네 </S>
집 에 <SP> 가 고 싶 다 응 <SP> 무 슌 <SP> 일 <SP> 하 길 래 ? ? </S>
퇴 근 이 <SP> 언 제 야 평 범 한 <SP> 직 장 인 이 지 <SP> 뭐 </S>
안 념 하 세 요 안 념 하 세 요 ! </S>
밥 <SP> 먹 었 나 요 ! ㅎㅎ 먹 었 습 니 다 </S>
저 는 <SP> 국 밥 ㅋㅋㅋㅋㅋㅋ 가 성 비 가 <SP> 좋 은 <SP> 국 밥 을

vocab.txt - Windows 메모장

파일(F) 편집(E) 서식(O) 보기(V) 도움말(H)

각
값
간
값
값
값

실습

```
from tqdm import tqdm
```

100% |██████████| 547958/547958 [00:09<00:00, 59999.28it/s]

모델 설계: __init__

```
class TransformerChat(nn.Module):
```

```
    def __init__(self, config):
        super().__init__()
```

```
        # 전체 단어(음절) 개수
```

```
        self.vocab_size = config['vocab_size']
```

```
        # 단어(음절) 벡터 크기
```

```
        self.embedding_size = config['embedding_size']
```

```
        # Transformer의 Attention Head 개수
```

```
        self.num_heads = config['num_heads']
```

```
        # Transformer Encoder의 Layer 수
```

```
        self.num_encoder_layers = config['num_encoder_layers']
```

```
        # Transformer Decoder의 Layer 수
```

```
        self.num_decoder_layers = config['num_decoder_layers']
```

```
        # 입력 Sequence의 최대 길이
```

```
        self.max_length = config['max_length']
```

```
        # Transformer 내부 FNN 크기
```

```
        self.hidden_size = config['hidden_size']
```

```
        # Token Embedding Matrix 선언
```

```
        self.embeddings = nn.Embedding(self.vocab_size, self.embedding_size)
```

```
        # Transformer Encoder-Decoder 설계(선언)
```

```
        self.transformer = nn.Transformer(d_model=self.embedding_size, nhead=self.num_heads, num_encoder_layers=self.num_encoder_layers,
                                           num_decoder_layers=self.num_decoder_layers, dim_feedforward=self.hidden_size)
```

```
        # 입력 길이 L에 대한 (L x L) mask 생성: 이전 토큰들의 정보만을 반영하기 위한 mask
```

```
        # [[1, -inf, -inf, -inf],
```

```
        #      [1, 1, -inf, -inf],
```

```
        #      [1, 1, 1, -inf],
```

```
        #      [1, 1, 1, 1]]
```

```
        self.mask = self.transformer.generate_square_subsequent_mask(self.max_length).cuda()
```

```
        # 전체 단어 분포로 변환하기 위한 linear
```

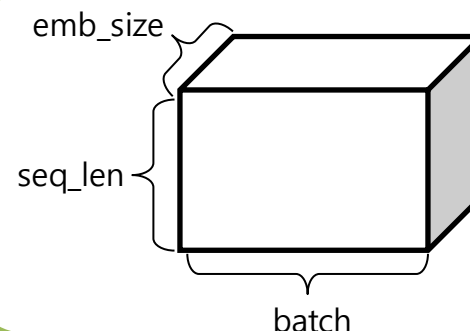
```
        self.projection_layer = nn.Linear(self.embedding_size, self.vocab_size)
```



실습

모델 설계: forward

```
def forward(self, enc_inputs, dec_inputs):  
  
    # enc_inputs: [batch, seq_len], dec_inputs: [batch, seq_len]  
    # enc_input_features: [batch, seq_len, emb_size] -> [seq_len, batch, emb_size]  
    enc_input_features = self.embeddings(enc_inputs).transpose(0, 1)  
  
    # dec_input_features: [batch, seq_len, emb_size] -> [seq_len, batch, emb_size]  
    dec_input_features = self.embeddings(dec_inputs).transpose(0, 1)  
  
    # dec_output_features: [seq_len, batch, emb_size]  
    dec_output_features = self.transformer(src=enc_input_features, tgt=dec_input_features, src_mask = self.mask, tgt_mask = self.mask)  
  
    # hypothesis : [seq_len, batch, vocab_size]  
    hypothesis = self.projection_layer(dec_output_features)  
  
    return hypothesis
```



실습

사전 읽기

```
# 어휘사전(vocabulary) 생성 함수
def load_vocab(file_dir):

    with open(file_dir, 'r', encoding='utf8') as vocab_file:
        char2idx = {}
        idx2char = {}
        index = 0
        for char in vocab_file:
            char = char.strip()
            char2idx[char] = index
            idx2char[index] = char
            index+=1

    return char2idx, idx2char
```

전처리

```
# 문자 입력열을 인덱스로 변환하는 함수
def convert_data2feature(config, input_sequence, char2idx, decoder_input=False):

    # 고정 길이 벡터 생성
    input_features = np.zeros(config["max_length"], dtype=np.int)

    if decoder_input:
        # Decoder Input은 Target Sequence에서 Right Shift
        # Target Sequence : ["안", "녕", "하", "세", "요", "</S>"]
        # Decoder Input Sequence : ["<S>", "안", "녕", "하", "세", "요"]
        input_sequence = " ".join(["<S>"] + input_sequence.split()[1:-1])

    for idx, token in enumerate(input_sequence.split()):
        if token in char2idx.keys():
            input_features[idx] = char2idx[token]
        else:
            input_features[idx] = char2idx['<UNK>']

    return input_features
```

데이터 로드

```
# 데이터 읽기 함수
def load_dataset(config):

    # 어휘사전 읽어오기
    char2idx, idx2char = load_vocab(config['vocab_file'])

    file_dir = config['train_file']
    data_file = open(file_dir, 'r', encoding='utf8').readlines()

    # 데이터를 저장하기 위한 리스트 생성
    enc_inputs, dec_inputs, dec_outputs = [], [], []

    for line in tqdm(data_file):

        line = line.strip().split('##')

        input_sequence = line[0]
        output_sequence = line[1]

        enc_inputs.append(convert_data2feature(config, input_sequence, char2idx))
        dec_inputs.append(convert_data2feature(config, output_sequence, char2idx, True))
        dec_outputs.append(convert_data2feature(config, output_sequence, char2idx))

    # 전체 데이터를 저장하고 있는 리스트를 텐서 형태로 변환
    enc_inputs = torch.tensor(enc_inputs, dtype=torch.long)
    dec_inputs = torch.tensor(dec_inputs, dtype=torch.long)
    dec_outputs = torch.tensor(dec_outputs, dtype=torch.long)

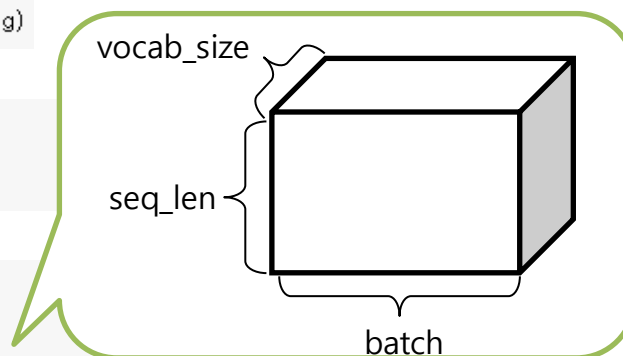
    return enc_inputs, dec_inputs, dec_outputs, char2idx, idx2char
```



실습

학습

```
def train(config):  
  
    # Transformer Seq2Seq 모델 객체 생성  
    model = TransformerChat(config).cuda()  
  
    # 데이터 읽기  
    enc_inputs, dec_inputs, dec_outputs, word2idx, idx2word = load_dataset(config)  
  
    for epoch in range(config["epoch"] + 1):  
        for (step, batch) in enumerate(train_dataloader):  
  
            enc_inputs, dec_inputs, dec_outputs = batch  
  
            # hypothesis: [seq_len, batch, vocab_size] -> [seq_len+batch, vocab_size]  
            hypothesis = model(enc_inputs, dec_inputs).view(-1, config['vocab_size'])  
  
            # labels: [batch, seq_len] -> [seq_len, batch] -> [seq_len(max_length)+batch]  
            labels = dec_outputs.transpose(0, 1)  
            labels = labels.reshape(config["max_length"]*dec_inputs.size(0))  
  
            # 비용 계산 및 역전파 수행: cross_entropy 내부에서 labels를 원핫벡터로 변환 (골드레이블은 항상 1차원으로 입력)  
            loss = loss_func(hypothesis, labels)  
            loss.backward()  
            optimizer.step()
```



실습

평가

```
def do_test(config, model, word2idx, idx2word, input_sequence="오늘 약속있으세요?"):

    # 평가 모드 셋팅
    model.eval()

    # 입력된 문자열의 음절을 공백 단위 토큰으로 변환. 공백은 <SP>로 변환: "오늘 약속" -> "오 늘 <SP> 약 속"
    input_sequence = " ".join([e if e != " " else "<SP>" for e in input_sequence])

    # 텐서 변환: [1, seq_len]
    enc_inputs = torch.tensor([convert_data2feature(config, input_sequence, word2idx)], dtype=torch.long).cuda()

    # input_ids : [1, seq_len] -> 첫번째 디코더 입력 "<S>" 만들기
    dec_inputs = torch.tensor([convert_data2feature(config, "", word2idx, True)], dtype=torch.long).cuda()

    # 시스템 응답 문자열 초기화
    response = ''

    # 최대 입력 길이 만큼 Decoding Loop
    for decoding_step in range(config['max_length']-1):

        # dec_outputs: [vocab_size]
        dec_outputs = model(enc_inputs, dec_inputs)[decoding_step, 0, :]
        # 가장 큰 출력력을 갖는 인덱스 얻어오기
        dec_output_idx = np.argmax(tensor2list(dec_outputs))

        # 생성된 토큰은 dec_inputs에 추가 (첫번째 차원은 배치)
        dec_inputs[0][decoding_step+1] = dec_output_idx

        # </S> 심볼 생성 시, Decoding 종료
        if idx2word[dec_output_idx] == "</S>":
            break

        # 생성 토큰 추가
        response += idx2word[dec_output_idx]

    # <SP>를 공백으로 변환한 후 응답 문자열 출력
    print(response.replace("<SP>", " "))
```

```
def test(config):

    # 어휘사전 읽어오기
    word2idx, idx2word = load_vocab(config['vocab_file'])

    # Transformer Seq2Seq 모델 객체 생성
    model = TransformerChat(config).cuda()

    # 학습한 모델 파일로부터 가중치 불러옴
    model.load_state_dict(torch.load(os.path.join(config["output_dir"], config["trained_model_name"])))

    while(True):
        input_sequence = input("문장을 입력하세요. (종료는 exit을 입력하세요.) : ")
        if input_sequence == 'exit':
            break
        do_test(config, model, word2idx, idx2word, input_sequence)
```



실습

메인 함수

```
if(__name__=="__main__"):

    root_dir = "/gdrive/My Drive/colab/transformer/chatbot/"
    output_dir = os.path.join(root_dir, "output")
    if not os.path.exists(output_dir):
        os.makedirs(output_dir)
    config = {"mode": "train",
              "vocab_file": os.path.join(root_dir, "vocab.txt"),
              "train_file": os.path.join(root_dir, "train.txt"),
              "trained_model_name": "epoch{}.pt".format(10),
              "output_dir": output_dir,
              "epoch": 10,
              "learn_rate": 0.00005,
              "num_encoder_layers": 6,
              "num_decoder_layers": 6,
              "num_heads": 4,
              "max_length": 20,
              "batch_size": 128,
              "embedding_size": 256,
              "hidden_size": 512,
              "vocab_size": 4427
             }

    if(config["mode"] == "train"):
        train(config)
    else:
        test(config)
```

```
100%|██████████| 547958/547958 [00:08<00:00, 63267.86it/s]
Current Step : 200 / 1 Current Loss : 2.773489
Current Step : 400 / 1 Current Loss : 2.226587
Current Step : 600 / 1 Current Loss : 2.164585
Current Step : 800 / 1 Current Loss : 2.179907
Current Step : 1000 / 1 Current Loss : 1.884135
Current Step : 1200 / 1 Current Loss : 1.987994
Current Step : 1400 / 1 Current Loss : 1.859728
Current Step : 1600 / 1 Current Loss : 2.048557
```

```
문장을 입력하세요. (종료는 exit을 입력하세요.) : 안녕하세요~
안녕하세요
문장을 입력하세요. (종료는 exit을 입력하세요.) : 밥은 먹었어?
아직 안먹었어
문장을 입력하세요. (종료는 exit을 입력하세요.) : 부대찌개 먹을까?
그래서 그런가
문장을 입력하세요. (종료는 exit을 입력하세요.) : 우리 놀러가자!
ㅋㅋㅋㅋㅋㅋ
문장을 입력하세요. (종료는 exit을 입력하세요.) : 영화 보러 갈래요?
영화보고 있어요
문장을 입력하세요. (종료는 exit을 입력하세요.) : 학교 같이 가자!
아 그래?
문장을 입력하세요. (종료는 exit을 입력하세요.) : 지금 뭐하고 있어?
나 일하고 있어
문장을 입력하세요. (종료는 exit을 입력하세요.) : exit
```



학습데이터 다운로드

- NIA AI-Hub 개방 데이터
 - 가입: <https://aihub.or.kr>

– 다운로드

- 온라인 구어체 말뭉치

회원유형 : 표시 항목은 필수 입력 항목입니다.

회원유형 ☐ 대기업 ☐ 중소기업 ☒ 대학교 ☐ 정부/공공기관 ☐ 연구기관 ☐ 개인 ☐ 기타

로그인 정보

아이디 : [\[비밀번호 인증 안함\]](#) [이메일 인증](#)

이메일 인증번호 : [인증확인](#)

비밀번호 : [\[사용가능\]](#)

* 9자 이상 20자 이하,
영문(대/소문자)/숫자/특수문자 등 3개의 문자 유형 모두 포함

비밀번호 확인 : [\[비밀번호 일치\]](#)

개인정보

이름 :

데이터 신청 동의 : ☒ 동의 ☐ 미동의

- - [휴대전화인증](#)

* 휴대전화 번호를 입력해주세요. 휴대전화 인증 후 데이터 이용신청을 할 수 있습니다.

기업 기관명 :

이메일 수신 동의 : ☒ 이메일 수신 동의

* 데이터 신청 처리 현황 등 주요 서비스 처리 현황에 대한 알림 메일이 발송됩니다.

#챗봇 #스마트스피커 #AI #인공지능 대규모 한국어 학습용 데이터 #데이터생태계 #크라우드소싱 #문자인식

온라인 구어체 말뭉치

분야 유형

갱신년월: 2022-08 구축년도: 2021 조회수: 962 다운로드: 28 용량: 1.61 GB

[다운로드](#) [샘플 데이터](#) [관심데이터 등록](#) [13](#)



학습데이터 다운로드

- NIA AI-Hub 개방 데이터
 - 데이터 선별: 선호하는 분야 1만 문장 선별

No.	항목명	비율
1	영상 데이터 엔터테인먼트 댓글	26.1%
2	영상 데이터 사회 댓글	23.4%
3	DIATV 게시판	18.5%
4	영상 데이터 게임 댓글	5.1%
5	영상 데이터 생활 댓글	4.9%
6	영상 데이터 동물 댓글	3.7%
7	영상 데이터 제품리뷰 댓글	3.7%
8	영상 데이터 스포츠 댓글	3.6%
9	영상 데이터 취미 댓글	3.1%
10	영상 데이터 연예 댓글	2.2%
11	영상 데이터 음식 댓글	2.0%
12	영상 데이터 뷰티 댓글	1.5%

- 데이터 분할
 - 학습데이터:평가데이터=9:1로 분할



과제

- NIA AI-Hub 데이터를 이용하여 Transformer 기반 챗봇 모델을 만들고 학습데이터 준비 과정, 소스코드, 실행 예제를 제출하세요.



질의응답

Q & A

Homepage: <http://nlp.konkuk.ac.kr>
E-mail: nlpdrkim@konkuk.ac.kr

