

Transformer



본 챗터의 내용은 "Ta-Chun (Bgg/Gene) Su's blog (<https://medium.com/@bgg/seq2seq-pay-attention-to-self-attention-part-2-cf81bf32c73d>)" 자료를 참고하여 편집되었습니다.

Attention Models

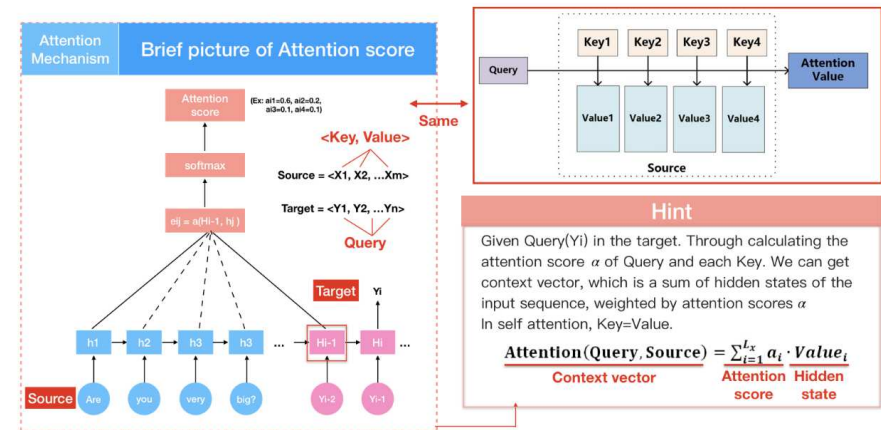


그림 출처: Ta-Chun (Bgg/Gene) Su's blog



Edited by Harksoo Kim

Attention Models in Detail

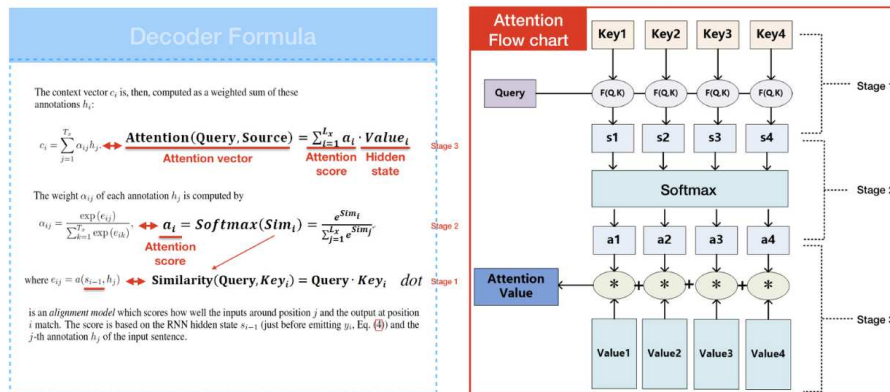


그림 출처: Ta-Chun (Bgg/Gene) Su's blog



Edited by Harksoo Kim

Problems of Attention-Based Models

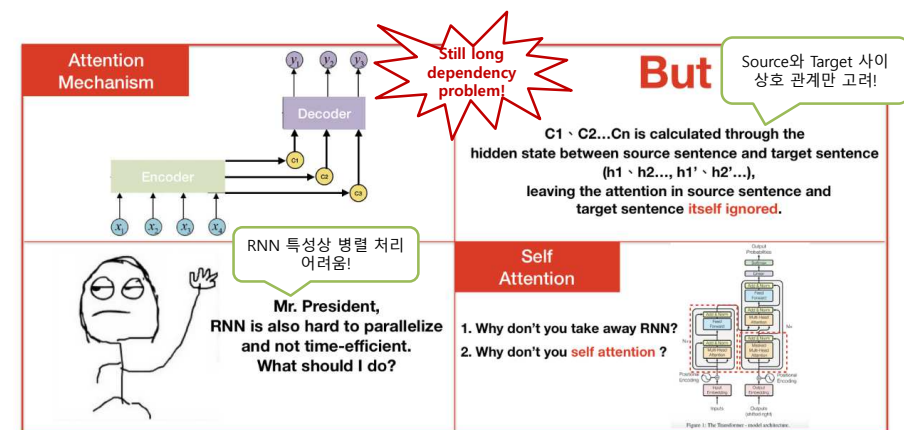
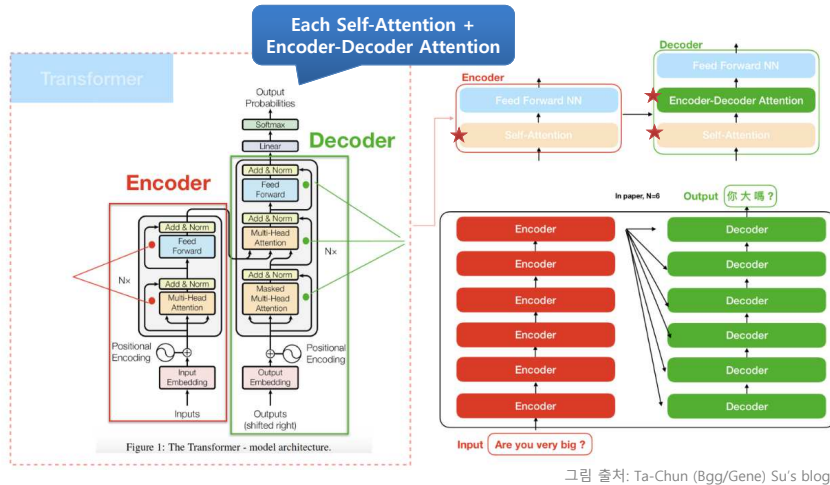


그림 출처: Ta-Chun (Bgg/Gene) Su's blog

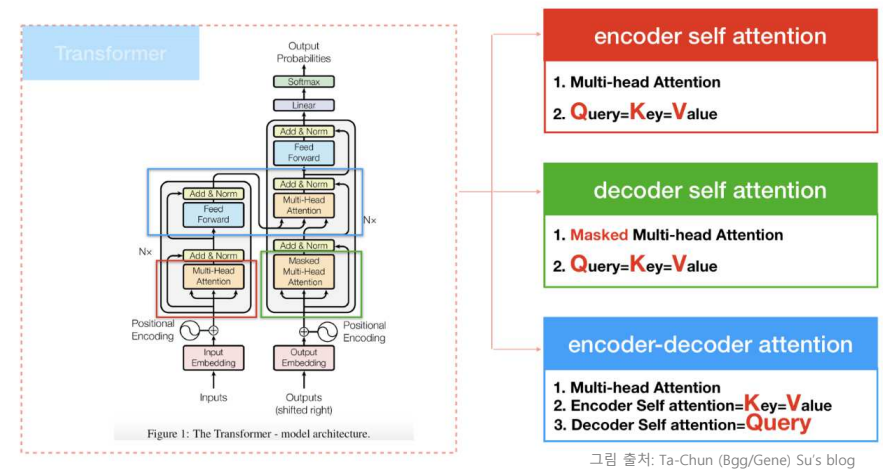


Edited by Harksoo Kim

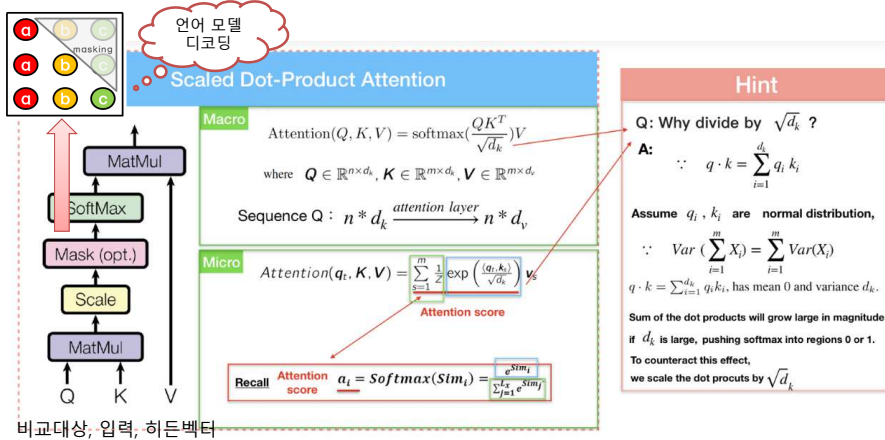
Transformer



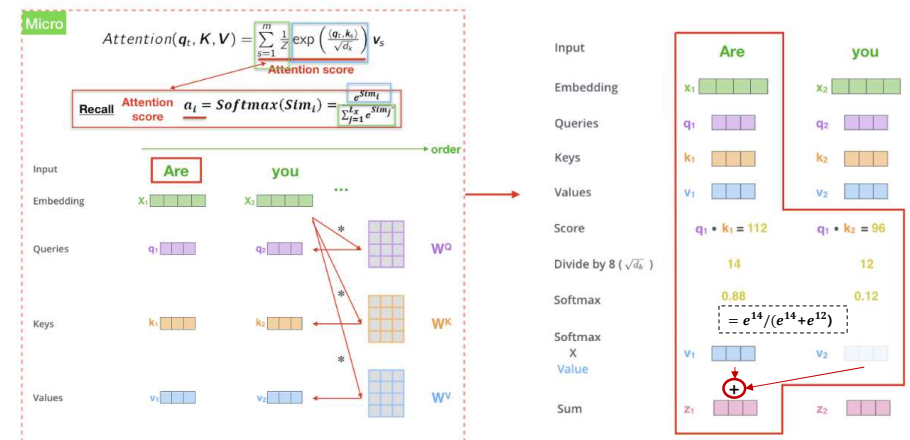
Attentions in Transformer



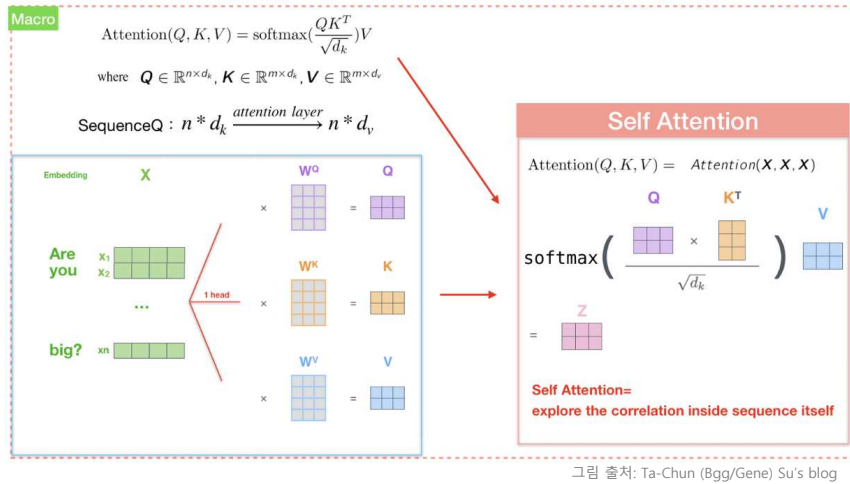
Scaled Dot-Product Attention



Calculation of Attentions



Calculation of Attentions



확인 문제

- 다음과 같이 단어 임베딩이 주어졌을 때, self-attention score를 계산시오. (소수점 이하 두 자리에서 반올림)

- are: [1,1], you: [2,1]

- root(2)=1.4로 계산

$$\text{Attention}(q_i, K, V) = \sum_{s=1}^m \frac{1}{Z} \exp\left(\frac{(q_i, k_s)}{\sqrt{d_k}}\right) v_s$$

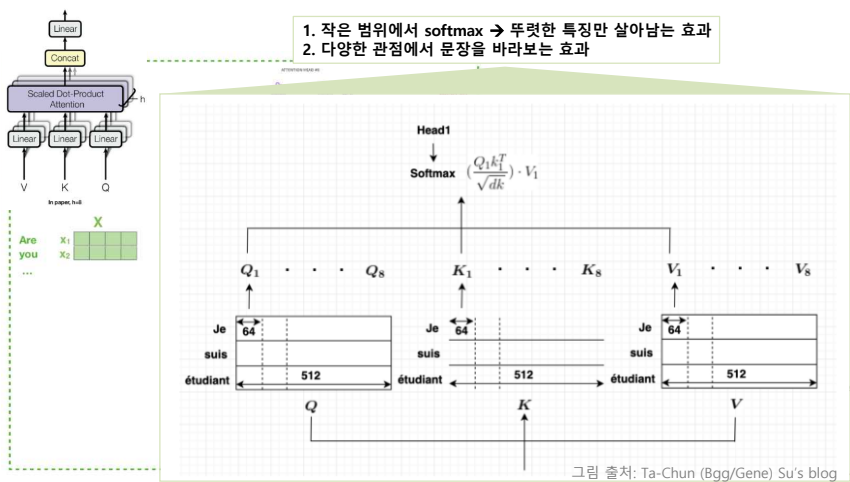
Scaled dot-product

	are	you
are		
you		

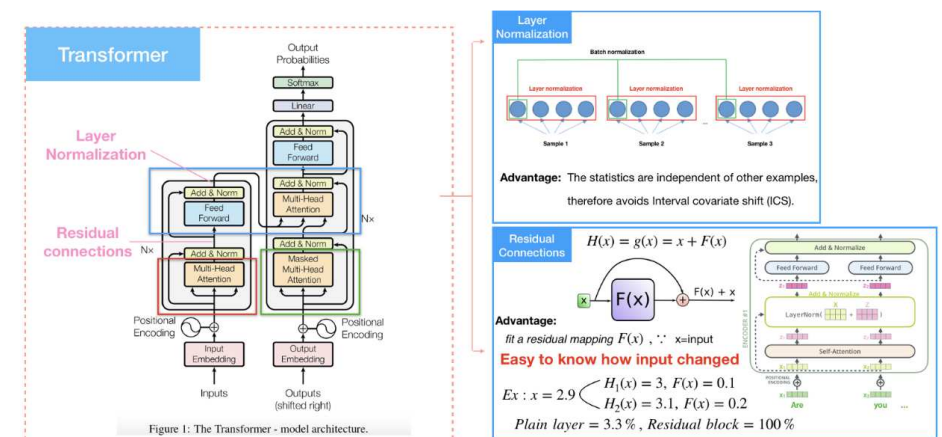
Self-attention score

	are	you
are		
you		

Multi-Head Attention

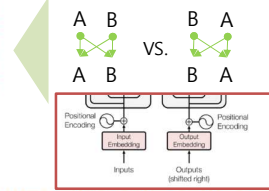


Layer Norm. & Residual Conn.



Position Encoding

Problem The multi-head attention network **cannot** naturally **make use of** the position of the words in the input sequence. The output of the multi-head attention network would be **the same** for the same sentences in different order.



Sol Positional Encoding

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

if shape of enc_dec = [T, d_model] pos: position of the word
→ then pos ∈ [0, T], i ∈ [0, d_model/2] t: Element i in d_model

Advantage: PE[pos+k] can be represented as a linear function of PE[pos], so the relative position between different embeddings can be easily inferred

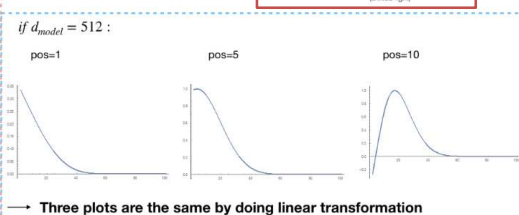
$$\sin(\alpha + \beta) = \sin\alpha\cos\beta + \cos\alpha\sin\beta$$

$$\cos(\alpha + \beta) = \cos\alpha\cos\beta - \sin\alpha\sin\beta$$

∴ Trigonometric Periodicity

∴ Learned the relation between relative and absolute position

Then Word embedding + positional encoding
(sum, concat...)



Three plots are the same by doing linear transformation

그림 출처: Ta-Chun (Bgg/Gene) Su's blog

Linear & Softmax

Linear + Softmax

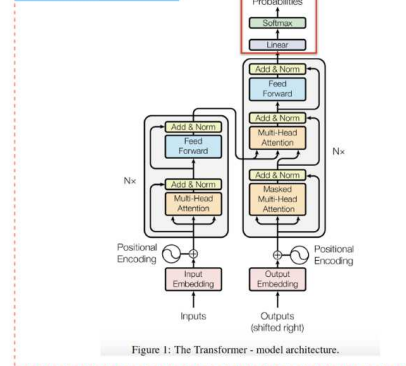


Figure 1: The Transformer - model architecture.

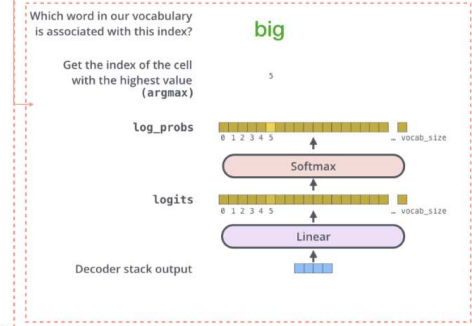


그림 출처: Ta-Chun (Bgg/Gene) Su's blog

실습

실습 코드 다운로드:
<https://github.com/KUNLP/Lecture>

• Transformer를 이용하여 1문 1답 챗봇 프로그램을 작성 하시오.

- 입력 데이터셋
 - 1문 1답 텍스트: 질문 wtt 응답 (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 있어? 이제 <SP> 수업 <SP> 끝났어! </SP>
난 <SP> 뭐해 대학생이야? </SP>
난 <SP> 지금 <SP> 과자 <SP> 먹어 <SP> 대학생 </SP>
중 <SP> 대학생 전공이 <SP> 뭐야? </SP>
전공이 <SP> 뭐야? 난 <SP> 집인 가보네 <SP> 부럽다 <SP> 나
난 <SP> 국어 국문학과 오 <SP> 좋은 <SP> 학문 <SP> 배우네 </SP>
집에 <SP> 가고 싶다 영 <SP> 무슨 <SP> 일 <SP> 하길래? </SP>
퇴근이 <SP> 언제야 평범한 <SP> 직장인이지 <SP> 뭐 </SP>
안녕하세요 안녕하세요! </SP>
밥 <SP> 먹었 나요! * * * 먹었습니다 </SP>
저는 <SP> 국밥 ㅋㅋㅋㅋㅋ 가성비가 <SP> 좋은 <SP> 국밥을

vocab.txt - Windows 메모장
파일(F) 편집(E) 서식(O) 보기(V) 도움말(H)
각
각
간
간
간
간
```

실습

from torch import torch 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 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
        self.mask = torch.zeros((self.max_length, self.max_length))
        self.mask = self.mask.triu_(1)

        # 전체 단어 분포로 변환하기 위한 linear
        self.projection_layer = nn.Linear(self.embedding_size, self.vocab_size)
```

실습

모델 설계: forward

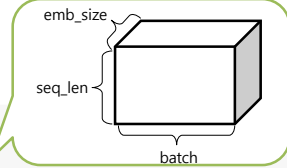
```
def forward(self, enc_inputs, dec_inputs):
    # enc_inputs: [batch, seq_len], dec_inputs: [batch, seq_len]
    # enc_input_features: [batch, seq_len, emb_size] -> [seq_len, batch, emb_size]
    enc_input_features = self.embeddings(enc_inputs).transpose(0, 1)

    # 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
```

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
```

전처리

```
# 문자 입력값을 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
```

실습

학습

```
def train(config):
    # Transformer Seq2Seq 모델 객체 생성
    model = TransformerChat(config).cuda()

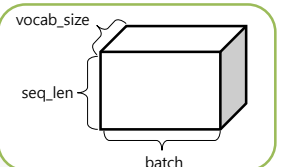
    # 데이터 읽기
    enc_inputs, dec_inputs, dec_outputs, word2idx, idx2word = load_dataset(config)

    for epoch in range(config['epoch'] + 1):
        for (step, batch) in enumerate(train_data_loader):
            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 <S>로 변환. "오늘 약속" -> "오늘 <S> 약속"
    input_sequence = " ".join([e if e != " " else "<S>" 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 확률을 갖는 인덱스 읽어오기
        dec_output_idx = torch.argmax(torch.tensor(dec_outputs))

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

        # </S> 실패할 경우 시, Decoding 종료
        if dec_output_idx == torch.tensor([config['end_token']]):
            break

        # 생성 토큰 추가
        response += idx2word[dec_output_idx]

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

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 exit를 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 [00:09<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 : 2.164565
Current Step : 800 / 1 Current Loss : 2.179307
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 : 2.046557
```

```
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!
바다 어디로 가?
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
```



Edited by Harksoo Kim

질의응답

Q&A

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



Edited by Harksoo Kim