

```
1 !pip install konlpy
```

```
Requirement already satisfied: konlpy in /usr/local/lib/python3.7/dist-packages (0.5.2)
Requirement already satisfied: colorama in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: beautifulsoup4==4.6.0 in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: lxml>=4.1.0 in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: numpy>=1.6 in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: tweepy>=3.7.0 in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: JPype1>=0.7.0 in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: requests-oauthlib>=0.7.0 in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: requests[socks]>=2.11.1 in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: six>=1.10.0 in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: typing-extensions; python_version < "3.8" in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: oauthlib>=3.0.0 in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: urllib3!=1.25.0,!1.25.1,<1.26,>=1.21.1 in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7/dist-packages (from konlpy)
Requirement already satisfied: PySocks!=1.5.7,>=1.5.6; extra == "socks" in /usr/local/lib/python3.7/dist-packages (from konlpy)
```

```
1 import numpy as np
2 import pandas as pd
3 import os
4 import re
5 import json
6 import enum
7 import matplotlib.pyplot as plt
8 import tensorflow as tf
9
10 from tqdm import tqdm
11 from konlpy.tag import Okt # 한글 형태소 활용
12 from konlpy.tag import Twitter
13 from sklearn.model_selection import train_test_split
14 from tensorflow.keras.callbacks import EarlyStopping, ModelCheckpoint
```

```
1 def load_data(path):
2     # 판다스를 통해서 데이터를 불러온다.
3     data_df = pd.read_csv(path, header=0)
4     # 질문과 답변 열을 가져와 question과 answer에 넣는다.
5     question, answer = list(data_df['Q']), list(data_df['A'])
6
7     return question, answer
```

```
1 def data_tokenizer(data):
2     # 토큰나이징 해서 담을 배열 생성
3     words = []
4     for sentence in data:
5         # FILTERS = "([~.,!?W\"'";:~])"
6         # 위 필터와 같은 값들을 정규화 표현식을
7         # 통해서 모두 "" 으로 변환 해주는 부분이다.
8         sentence = re.sub(CHANGE_FILTER, "", sentence)
9         for word in sentence.split():
10             words.append(word)
```

```

10         words.append(word)
11     # 토큰나이징과 정규표현식을 통해 만들어진
12     # 값들을 넘겨 준다.
13     return [word for word in words if word]

1 # 한글 텍스트를 토큰나이징 하기 위해 형태로소 분리하는 함수.
2 # Konlpy에서 제공하는 Okt를 사용해 형태소 기준으로 텍스트 데이터를 토큰나이징 한다.
3 def prepro_like_morphlized(data):
4     morph_analyzer = Okt()
5     result_data = list()
6     for seq in tqdm(data):
7         try :
8             morphlized_seq = " ".join(morph_analyzer.morphs(seq.replace(' ', '')))
9             result_data.append(morphlized_seq)
10        except Exception as exception :
11            print("오류 발생!!")
12    return result_data

1 def load_vocabulary(path, vocab_path, tokenize_as_morph=False):
2     # 사전을 담을 배열 준비한다.
3     vocabulary_list = []
4     # 사전을 구성한 후 파일로 저장 진행한다.
5     # 그 파일의 존재 여부를 확인한다.
6     if not os.path.exists(vocab_path):
7         # 이미 생성된 사전 파일이 존재하지 않으므로
8         # 데이터를 가지고 만들어야 한다.
9         # 그래서 데이터가 존재 하면 사전을 만들기 위해서
10        # 데이터 파일의 존재 여부를 확인한다.
11        if (os.path.exists(path)):
12            # 데이터가 존재하니 판다스를 통해서
13            # 데이터를 불러오자
14            data_df = pd.read_csv(path, encoding='utf-8')
15            # 판다스의 데이터 프레임을 통해서
16            # 질문과 답에 대한 열을 가져 온다.
17            question, answer = list(data_df['Q']), list(data_df['A'])
18            if tokenize_as_morph: # 형태로소에 따른 토큰나이저 처리
19                question = prepro_like_morphlized(question)
20                answer = prepro_like_morphlized(answer)
21            data = []
22            # 질문과 답변을 extend을
23            # 통해서 구조가 없는 배열로 만든다.
24            data.extend(question)
25            data.extend(answer)
26            # 토큰나이저 처리 하는 부분이다.
27            words = data_tokenizer(data)
28            # 공통적인 단어에 대해서는 모두
29            # 필요 없으므로 한개로 만들어 주기 위해서
30            # set해주고 이것들을 리스트로 만들어 준다.
31            words = list(set(words))
32            # 데이터 없는 내용중에 MARKER를 사전에
33            # 추가 하기 위해서 아래와 같이 처리 한다.
34            # 아래는 MARKER 값이며 리스트의 첫번째 부터
35            # 순서대로 넣기 위해서 인덱스 0에 추가한다.
36            # PAD = "<PADDING>"

```

```

37         # STD = "<START>"
38         # END = "<END>"
39         # UNK = "<UNKWON>"
40         words[:0] = MARKER
41         # 사전을 리스트로 만들었으니 이 내용을
42         # 사전 파일을 만들어 넣는다.
43         with open(vocab_path, 'w', encoding='utf-8') as vocabulary_file:
44             for word in words:
45                 vocabulary_file.write(word + '\n')
46
47         # 사전 파일이 존재하면 여기에서
48         # 그 파일을 불러서 배열에 넣어 준다.
49         with open(vocab_path, 'r', encoding='utf-8') as vocabulary_file:
50             for line in vocabulary_file:
51                 vocabulary_list.append(line.strip())
52
53         # 배열에 내용을 키와 값이 있는
54         # 딕셔너리 구조로 만든다.
55         char2idx, idx2char = make_vocabulary(vocabulary_list)
56         # 두가지 형태의 키와 값이 있는 형태를 리턴한다.
57         # (예) 단어: 인덱스 , 인덱스: 단어)
58         return char2idx, idx2char, len(char2idx)

```

```

1 def make_vocabulary(vocabulary_list):
2     # 리스트를 키가 단어이고 값이 인덱스인
3     # 딕셔너리를 만든다.
4     char2idx = {char: idx for idx, char in enumerate(vocabulary_list)}
5     # 리스트를 키가 인덱스이고 값이 단어인
6     # 딕셔너리를 만든다.
7     idx2char = {idx: char for idx, char in enumerate(vocabulary_list)}
8     # 두개의 딕셔너리를 넘겨 준다.
9     return char2idx, idx2char

```

```

1 def enc_processing(value, dictionary, tokenize_as_morph=False):
2     # 인덱스 값들을 가지고 있는
3     # 배열이다.(누적된다.)
4     sequences_input_index = []
5     # 하나의 인코딩 되는 문장의
6     # 길이를 가지고 있다.(누적된다.)
7     sequences_length = []
8     # 형태소 토크나이징 사용 유무
9     if tokenize_as_morph:
10         value = prepro_like_morphlized(value)
11
12     # 한줄씩 불러온다.
13     for sequence in value:
14         # FILTERS = "([~.,!?W'':;)(])"
15         # 정규화를 사용하여 필터에 들어 있는
16         # 값들을 "" 으로 치환 한다.
17         sequence = re.sub(CHANGE_FILTER, "", sequence)
18         # 하나의 문장을 인코딩 할때
19         # 가지고 있기 위한 배열이다.
20         sequence_index = []
21         # 문장을 스페이스 단위로

```

```

1 def dec_output_processing(value, dictionary, tokenize_as_morph=False):
2     # 인덱스 값들을 가지고 있는
3     # 배열이다.(누적된다)
4     sequences_output_index = []
5     # 하나의 디코딩 입력 되는 문장의
6     # 길이를 가지고 있다.(누적된다)
7     sequences_length = []
8     # 형태소 토크나이징 사용 유무
9     if tokenize_as_morph:
10         value = prepro_like_morphlized(value)
11     # 한줄씩 불러온다.
12     for sequence in value:
13         # FILTERS = "([~.,!?W'':;)(])"
14         # 정규화를 사용하여 필터에 들어 있는
15         # 값들을 "" 으로 치환 한다.
16         sequence = re.sub(CHANGE_FILTER, "", sequence)
17         # 하나의 문장을 디코딩 할때 가지고
18         # 있기 위한 배열이다.
19         sequence_index = []
20         # 디코딩 입력의 처음에는 START가 와야 하므로
21         # 그 값을 넣어 주고 시작한다.
22         # 문장에서 스페이스 단위별로 단어를 가져와서 딕셔너리의
23         # 값인 인덱스를 넣어 준다.
24         sequence_index = [dictionary[STD]] + [dictionary[word] if word in dictionary else dictio
25         # 문장 제한 길이보다 길어질 경우 뒤에 토큰을 자르고 있다.
26         if len(sequence_index) > MAX_SEQUENCE:
27             sequence_index = sequence_index[0:MAX_SEQUENCE]

```

```

27         sequence_index = sequence_index[:MAX_SEQUENCE]
28     # 하나의 문장에 길이를 넣어주고 있다.
29     sequences_length.append(len(sequence_index))
30     # max_sequence_length보다 문장 길이가
31     # 작다면 빈 부분에 PAD(0)를 넣어준다.
32     sequence_index += (MAX_SEQUENCE - len(sequence_index)) * [dictionary[PAD]]
33     # 인덱스화 되어 있는 값을
34     # sequences_output_index 넣어 준다.
35     sequences_output_index.append(sequence_index)
36 # 인덱스화된 일반 배열을 넘파이 배열로 변경한다.
37 # 이유는 텐서플로우 dataset에 넣어 주기 위한
38 # 사전 작업이다.
39 # 넘파이 배열에 인덱스화된 배열과 그 길이를 넘겨준다.
40 return np.asarray(sequences_output_index), sequences_length

1 def dec_target_processing(value, dictionary, tokenize_as_morph=False):
2     # 인덱스 값들을 가지고 있는
3     # 배열이다.(누적된다)
4     sequences_target_index = []
5     # 형태소 토크나이징 사용 유무
6     if tokenize_as_morph:
7         value = prepro_like_morphlized(value)
8     # 한줄씩 불러온다.
9     for sequence in value:
10         # FILTERS = "([~.,!?W\"'':;])"
11         # 정규화를 사용하여 필터에 들어 있는
12         # 값들을 "" 으로 치환 한다.
13         sequence = re.sub(CHANGE_FILTER, "", sequence)
14         # 문장에서 스페이스 단위별로 단어를 가져와서
15         # 딕셔너리의 값인 인덱스를 넣어 준다.
16         # 디코딩 출력의 마지막에 END를 넣어 준다.
17         sequence_index = [dictionary[word] if word in dictionary else dictionary[UNK] for word in sequence]
18         # 문장 제한 길이보다 길어질 경우 뒤에 토큰을 자르고 있다.
19         # 그리고 END 토큰을 넣어 준다
20         if len(sequence_index) >= MAX_SEQUENCE:
21             sequence_index = sequence_index[:MAX_SEQUENCE - 1] + [dictionary[END]]
22         else:
23             sequence_index += [dictionary[END]]
24         # max_sequence_length보다 문장 길이가
25         # 작다면 빈 부분에 PAD(0)를 넣어준다.
26         sequence_index += (MAX_SEQUENCE - len(sequence_index)) * [dictionary[PAD]]
27         # 인덱스화 되어 있는 값을
28         # sequences_target_index에 넣어 준다.
29         sequences_target_index.append(sequence_index)
30 # 인덱스화된 일반 배열을 넘파이 배열로 변경한다.
31 # 이유는 텐서플로우 dataset에 넣어 주기 위한 사전 작업이다.
32 # 넘파이 배열에 인덱스화된 배열과 그 길이를 넘겨준다.
33 return np.asarray(sequences_target_index)

1 # 시각화 함수
2 def plot_graphs(history, string):
3     plt.plot(history.history[string])
4     plt.plot(history.history['val_'+string], '')
5     plt.xlabel("Epochs")

```

```

6     plt.ylabel(string)
7     plt.legend([string, 'val_'+string])
8     plt.show()

1 def create_padding_mask(seq):
2     seq = tf.cast(tf.math.equal(seq, 0), tf.float32)
3
4     # add extra dimensions to add the padding
5     # to the attention logits.
6     return seq[:, tf.newaxis, tf.newaxis, :] # (batch_size, 1, 1, seq_len)

1 def create_look_ahead_mask(size):
2     mask = 1 - tf.linalg.band_part(tf.ones((size, size)), -1, 0)
3     return mask # (seq_len, seq_len)

1 def create_masks(inp, tar):
2     # Encoder padding mask
3     enc_padding_mask = create_padding_mask(inp)
4
5     # Used in the 2nd attention block in the decoder.
6     # This padding mask is used to mask the encoder outputs.
7     dec_padding_mask = create_padding_mask(inp)
8
9     # Used in the 1st attention block in the decoder.
10    # It is used to pad and mask future tokens in the input received by
11    # the decoder.
12    look_ahead_mask = create_look_ahead_mask(tf.shape(tar)[1])
13    dec_target_padding_mask = create_padding_mask(tar)
14    combined_mask = tf.maximum(dec_target_padding_mask, look_ahead_mask)
15
16    return enc_padding_mask, combined_mask, dec_padding_mask

1 # pos/100002i/dim 값을 만드는 함수
2 def get_angles(pos, i, d_model):
3     angle_rates = 1 / np.power(10000, (2 * i//2) / np.float32(d_model))
4     return pos * angle_rates

1 def positional_encoding(position, d_model):
2     angle_rads = get_angles(np.arange(position)[:, np.newaxis],
3                             np.arange(d_model)[np.newaxis, :],
4                             d_model)
5
6     # apply sin to even indices in the array; 2i
7     angle_rads[:, 0::2] = np.sin(angle_rads[:, 0::2])
8
9     # apply cos to odd indices in the array; 2i+1
10    angle_rads[:, 1::2] = np.cos(angle_rads[:, 1::2])
11
12    pos_encoding = angle_rads[np.newaxis, ...]
13
14    return tf.cast(pos_encoding, dtype=tf.float32)

```

```

1 def scaled_dot_product_attention(q, k, v, mask):
2     """Calculate the attention weights.
3     q, k, v must have matching leading dimensions.
4     k, v must have matching penultimate dimension, i.e.: seq_len_k = seq_len_v.
5     The mask has different shapes depending on its type(padding or look ahead)
6     but it must be broadcastable for addition.
7
8     Args:
9     q: query shape == (... , seq_len_q, depth)
10    k: key shape == (... , seq_len_k, depth)
11    v: value shape == (... , seq_len_v, depth_v)
12    mask: Float tensor with shape broadcastable
13          to (... , seq_len_q, seq_len_k). Defaults to None.
14
15    Returns:
16    output, attention_weights
17    """
18
19    matmul_qk = tf.matmul(q, k, transpose_b=True) # (... , seq_len_q, seq_len_k)
20
21    # scale matmul_qk
22    dk = tf.cast(tf.shape(k)[-1], tf.float32)
23    scaled_attention_logits = matmul_qk / tf.math.sqrt(dk)
24
25    # add the mask to the scaled tensor.
26    # 매우 작은 값을 넣어줌으로써 다음 레이어인 softmax에서 무시할 수 있도록 함 (softmax: 매우작
27    if mask is not None:
28        scaled_attention_logits += (mask * -1e9)
29
30    # softmax is normalized on the last axis (seq_len_k) so that the scores
31    # add up to 1.
32    attention_weights = tf.nn.softmax(scaled_attention_logits, axis=-1) # (... , seq_len_q, seq_
33
34    output = tf.matmul(attention_weights, v) # (... , seq_len_q, depth_v)
35
36    return output, attention_weights

```

```

1 class MultiHeadAttention(tf.keras.layers.Layer):
2     def __init__(self, **kwargs):
3         super(MultiHeadAttention, self).__init__()
4         self.num_heads = kwargs['num_heads']
5         self.d_model = kwargs['d_model']
6
7         # 나머지가 발생하면 에러가 나게끔 함
8         assert self.d_model % self.num_heads == 0
9
10        self.depth = self.d_model // self.num_heads
11
12        self.wq = tf.keras.layers.Dense(kwargs['d_model'])
13        self.wk = tf.keras.layers.Dense(kwargs['d_model'])
14        self.wv = tf.keras.layers.Dense(kwargs['d_model'])
15
16        self.dense = tf.keras.layers.Dense(kwargs['d_model'])

```

```

17
18 # key, query, value에 대한 벡터를 헤드 수만큼 분리할 수 있게 하는 함수 [batch, sequence, fea
19 def split_heads(self, x, batch_size):
20     """Split the last dimension into (num_heads, depth).
21     Transpose the result such that the shape is (batch_size, num_heads, seq_len, depth)
22     """
23     x = tf.reshape(x, (batch_size, -1, self.num_heads, self.depth)) # 피쳐 차원을 헤드 수 만
24     return tf.transpose(x, perm=[0, 2, 1, 3]) # 시퀀스, 헤드 자리를 바꿈: [batch, head, seq
25
26 def call(self, v, k, q, mask):
27     batch_size = tf.shape(q)[0]
28
29     q = self.wq(q) # (batch_size, seq_len, d_model)
30     k = self.wk(k) # (batch_size, seq_len, d_model)
31     v = self.wv(v) # (batch_size, seq_len, d_model)
32
33     q = self.split_heads(q, batch_size) # (batch_size, num_heads, seq_len_q, depth)
34     k = self.split_heads(k, batch_size) # (batch_size, num_heads, seq_len_k, depth)
35     v = self.split_heads(v, batch_size) # (batch_size, num_heads, seq_len_v, depth)
36
37     # scaled_attention.shape == (batch_size, num_heads, seq_len_q, depth)
38     # attention_weights.shape == (batch_size, num_heads, seq_len_q, seq_len_k)
39     scaled_attention, attention_weights = scaled_dot_product_attention(
40         q, k, v, mask)
41
42     scaled_attention = tf.transpose(scaled_attention, perm=[0, 2, 1, 3]) # (batch_size, seq
43
44     concat_attention = tf.reshape(scaled_attention,
45                                   (batch_size, -1, self.d_model)) # (batch_size, seq_len_q,
46
47     output = self.dense(concat_attention) # (batch_size, seq_len_q, d_model)
48
49     return output, attention_weights

```



```

1 def point_wise_feed_forward_network(**kwargs):
2     return tf.keras.Sequential([
3         tf.keras.layers.Dense(kwargs['dff'], activation='relu'), # (batch_size, seq_len, dff)
4         tf.keras.layers.Dense(kwargs['d_model']) # (batch_size, seq_len, d_model)
5     ])

```



```

1 class EncoderLayer(tf.keras.layers.Layer):
2     def __init__(self, **kwargs):
3         super(EncoderLayer, self).__init__()
4
5         self.mha = MultiHeadAttention(**kwargs)
6         self.ffn = point_wise_feed_forward_network(**kwargs)
7
8         self.layernorm1 = tf.keras.layers.LayerNormalization(epsilon=1e-6)
9         self.layernorm2 = tf.keras.layers.LayerNormalization(epsilon=1e-6)
10
11         self.dropout1 = tf.keras.layers.Dropout(kwargs['rate'])
12         self.dropout2 = tf.keras.layers.Dropout(kwargs['rate'])
13
14     def call(self, x, mask):

```



```

15     attn_output, _ = self.mha(x, x, x, mask) # self attention (batch_size, input_seq_len, d
16     attn_output = self.dropout1(attn_output)
17     out1 = self.layernorm1(x + attn_output) # (batch_size, input_seq_len, d_model)
18
19     ffn_output = self.ffn(out1) # (batch_size, input_seq_len, d_model)
20     ffn_output = self.dropout2(ffn_output)
21     out2 = self.layernorm2(out1 + ffn_output) # (batch_size, input_seq_len, d_model)
22
23     return out2

```

```

1 class DecoderLayer(tf.keras.layers.Layer):
2     def __init__(self, **kwargs):
3         super(DecoderLayer, self).__init__()
4
5         self.mha1 = MultiHeadAttention(**kwargs)
6         self.mha2 = MultiHeadAttention(**kwargs)
7
8         self.ffn = point_wise_feed_forward_network(**kwargs)
9
10        self.layernorm1 = tf.keras.layers.LayerNormalization(epsilon=1e-6)
11        self.layernorm2 = tf.keras.layers.LayerNormalization(epsilon=1e-6)
12        self.layernorm3 = tf.keras.layers.LayerNormalization(epsilon=1e-6)
13
14        self.dropout1 = tf.keras.layers.Dropout(kwargs['rate'])
15        self.dropout2 = tf.keras.layers.Dropout(kwargs['rate'])
16        self.dropout3 = tf.keras.layers.Dropout(kwargs['rate'])
17
18
19    def call(self, x, enc_output, look_ahead_mask, padding_mask):
20        # enc_output.shape == (batch_size, input_seq_len, d_model)
21        attn1, attn_weights_block1 = self.mha1(x, x, x, look_ahead_mask) # (batch_size, target_
22        attn1 = self.dropout1(attn1)
23        out1 = self.layernorm1(attn1 + x)
24
25        attn2, attn_weights_block2 = self.mha2(
26            enc_output, enc_output, out1, padding_mask) # (batch_size, target_seq_len, d_model)
27        attn2 = self.dropout2(attn2)
28        out2 = self.layernorm2(attn2 + out1) # (batch_size, target_seq_len, d_model)
29
30        ffn_output = self.ffn(out2) # (batch_size, target_seq_len, d_model)
31        ffn_output = self.dropout3(ffn_output)
32        out3 = self.layernorm3(ffn_output + out2) # (batch_size, target_seq_len, d_model)
33
34        return out3, attn_weights_block1, attn_weights_block2

```

```

1 class Encoder(tf.keras.layers.Layer):
2     def __init__(self, **kwargs):
3         super(Encoder, self).__init__()
4
5         self.d_model = kwargs['d_model']
6         self.num_layers = kwargs['num_layers']
7
8         self.embedding = tf.keras.layers.Embedding(kwargs['input_vocab_size'], self.d_model)
9         self.pos_encoding = positional_encoding(kwargs['maximum_position_encoding'])

```

```

9         self.pos_encoding = positional_encoding(kargs['maximum_position_encoding'],
10                                                  self.d_model)
11
12
13         self.enc_layers = [EncoderLayer(**kargs)
14                             for _ in range(self.num_layers)]
15
16         self.dropout = tf.keras.layers.Dropout(kargs['rate'])
17
18     def call(self, x, mask):
19
20         seq_len = tf.shape(x)[1]
21
22         # adding embedding and position encoding.
23         x = self.embedding(x) # (batch_size, input_seq_len, d_model)
24         x *= tf.math.sqrt(tf.cast(self.d_model, tf.float32))
25         x += self.pos_encoding[:, :seq_len, :]
26
27         x = self.dropout(x)
28
29         for i in range(self.num_layers):
30             x = self.enc_layers[i](x, mask)
31
32         return x # (batch_size, input_seq_len, d_model)

```



```

1 class Decoder(tf.keras.layers.Layer):
2     def __init__(self, **kargs):
3         super(Decoder, self).__init__()
4
5         self.d_model = kargs['d_model']
6         self.num_layers = kargs['num_layers']
7
8         self.embedding = tf.keras.layers.Embedding(kargs['target_vocab_size'], self.d_model)
9         self.pos_encoding = positional_encoding(kargs['maximum_position_encoding'], self.d_model)
10
11         self.dec_layers = [DecoderLayer(**kargs)
12                             for _ in range(self.num_layers)]
13         self.dropout = tf.keras.layers.Dropout(kargs['rate'])
14
15     def call(self, x, enc_output, look_ahead_mask, padding_mask):
16         seq_len = tf.shape(x)[1]
17         attention_weights = {}
18
19         x = self.embedding(x) # (batch_size, target_seq_len, d_model)
20         x *= tf.math.sqrt(tf.cast(self.d_model, tf.float32))
21         x += self.pos_encoding[:, :seq_len, :]
22
23         x = self.dropout(x)
24
25         for i in range(self.num_layers):
26             x, block1, block2 = self.dec_layers[i](x, enc_output, look_ahead_mask, padding_mask)
27
28             attention_weights['decoder_layer{}_block1'.format(i+1)] = block1
29             attention_weights['decoder_layer{}_block2'.format(i+1)] = block2
30

```

```

31     # x.shape == (batch_size, target_seq_len, d_model)
32     return x, attention_weights

1 class Transformer(tf.keras.Model):
2     def __init__(self, **kargs):
3         super(Transformer, self).__init__(name=kargs['model_name'])
4         self.end_token_idx = kargs['end_token_idx']
5
6         self.encoder = Encoder(**kargs)
7         self.decoder = Decoder(**kargs)
8
9         self.final_layer = tf.keras.layers.Dense(kargs['target_vocab_size'])
10
11     def call(self, x):
12         inp, tar = x
13
14         enc_padding_mask, look_ahead_mask, dec_padding_mask = create_masks(inp, tar)
15         enc_output = self.encoder(inp, enc_padding_mask) # (batch_size, inp_seq_len, d_model)
16
17         # dec_output.shape == (batch_size, tar_seq_len, d_model)
18         dec_output, _ = self.decoder(
19             tar, enc_output, look_ahead_mask, dec_padding_mask)
20
21         final_output = self.final_layer(dec_output) # (batch_size, tar_seq_len, target_vocab_si
22
23         return final_output
24
25     def inference(self, x):
26         inp = x
27         tar = tf.expand_dims([STD_INDEX], 0)
28
29         enc_padding_mask, look_ahead_mask, dec_padding_mask = create_masks(inp, tar)
30         enc_output = self.encoder(inp, enc_padding_mask)
31
32         predict_tokens = list()
33         for t in range(0, MAX_SEQUENCE):
34             dec_output, _ = self.decoder(tar, enc_output, look_ahead_mask, dec_padding_mask)
35             final_output = self.final_layer(dec_output)
36             outputs = tf.argmax(final_output, -1).numpy()
37             pred_token = outputs[0][-1]
38             if pred_token == self.end_token_idx:
39                 break
40             predict_tokens.append(pred_token)
41             tar = tf.expand_dims([STD_INDEX] + predict_tokens, 0)
42             _, look_ahead_mask, dec_padding_mask = create_masks(inp, tar)
43
44         return predict_tokens

1 def loss(real, pred):
2     mask = tf.math.logical_not(tf.math.equal(real, 0))
3     loss_ = loss_object(real, pred)
4
5     mask = tf.cast(mask, dtype=loss_.dtype)
6     loss *= mask

```

```

6     return tf.reduce_mean(loss_)
7
8     return tf.reduce_mean(loss_)
9
10 def accuracy(real, pred):
11     mask = tf.math.logical_not(tf.math.equal(real, 0))
12     mask = tf.expand_dims(tf.cast(mask, dtype=pred.dtype), axis=-1)
13     pred *= mask
14     acc = train_accuracy(real, pred)
15
16     return tf.reduce_mean(acc)

```

```

1 from google.colab import drive
2 drive.mount('/content/drive')

```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/c



```

1 os.chdir('/content/drive/MyDrive/시립대_딥러닝/과제')

```

```

1 FILTERS = "([~.,!?W\"'':;)])"
2 PAD = "<PAD>"
3 STD = "<SOS>"
4 END = "<END>"
5 UNK = "<UNK>"
6
7 PAD_INDEX = 0
8 STD_INDEX = 1
9 END_INDEX = 2
10 UNK_INDEX = 3
11
12 MARKER = [PAD, STD, END, UNK]
13 CHANGE_FILTER = re.compile(FILTERS) # ~.,!?W\"'':;)( 들과 매치되는지 확인하기 위함
14
15 MAX_SEQUENCE = 50

```

```

1 # PATH = 'data_in/ChatBotData_csv2.csv'
2 PATH = 'data_in/BlueHouse_okt_100_03.csv'
3 VOCAB_PATH = 'data_in/vocabulary.txt'

```

```

1 inputs, outputs = load_data(PATH)
2 char2idx, idx2char, vocab_size = load_vocabulary(PATH, VOCAB_PATH, tokenize_as_morph=False)

```

```

1 char2idx

```

```

1 index_inputs, input_seq_len = enc_processing(inputs, char2idx, tokenize_as_morph=True)
2 index_outputs, output_seq_len = dec_output_processing(outputs, char2idx, tokenize_as_morph=True)
3 index_targets = dec_target_processing(outputs, char2idx, tokenize_as_morph=True)

```

73%|██████████ | 12947/17773 [00:06<00:02, 1891.98it/s]

75%|██████████ | 13321/17773 [00:06<00:02, 2221.10it/s]

77%|██████████ | 13636/17773 [00:06<00:01, 2385.80it/s]

79%|██████████ | 13997/17773 [00:06<00:01, 2655.59it/s]

81%|██████████ | 14320/17773 [00:06<00:01, 2769.34it/s]

```
1 data_configs = {}  
2 data_configs['char2idx'] = char2idx  
3 data_configs['idx2char'] = idx2char
```

```

4 data_configs['vocab_size'] = vocab_size
5 data_configs['pad_symbol'] = PAD
6 data_configs['std_symbol'] = STD
7 data_configs['end_symbol'] = END
8 data_configs['unk_symbol'] = UNK

```

```

1 # 경로설정
2 DATA_IN_PATH = './data_in/'
3 DATA_OUT_PATH = './data_out/'
4 TRAIN_INPUTS = 'train_inputs.npy'
5 TRAIN_OUTPUTS = 'train_outputs.npy'
6 TRAIN_TARGETS = 'train_targets.npy'
7 DATA_CONFIGS = 'data_configs.json'

```

```

1 np.save(open(DATA_IN_PATH + TRAIN_INPUTS, 'wb'), index_inputs)
2 np.save(open(DATA_IN_PATH + TRAIN_OUTPUTS, 'wb'), index_outputs)
3 np.save(open(DATA_IN_PATH + TRAIN_TARGETS, 'wb'), index_targets)
4 json.dump(data_configs, open(DATA_IN_PATH + DATA_CONFIGS, 'w'))

```

```

1 char2idx
2 idx2char

```

```

934: '해가는',
935: '관계',
936: '없어지면',
937: '변한',
938: '씨야겠다',
939: '남동생',
940: '누',
941: '데려가주세요',
942: '꼭시키고 싶어',
943: '싸우셨나',
944: '외쳐',
945: '울적한',
946: '상책',
947: '일어',
948: '싸야',
949: '싸웠다',
950: '종합',
951: '동상',
952: '어마',
953: '마세요',
954: '태닝',
955: '나가겠어',
956: '달라졌어',
957: '새집',
958: '똑똑할까',
959: '사줘',
960: '사귀여 자친구가 환승이 별했어',
961: '미치겠네',
962: '깊어져',
963: '사용',
964: '여사',
965: '어플',
966: '변해',
967: '줄겠다',

```

```

968: '클',
969: '무관심한것',
970: '돌아오는길',
971: '있었습니다',
972: '어리석게',
973: '꽃',
974: '있게다가',
975: '만족할줄을',
976: '드는것도',
977: '익숙해지는',
978: '가오니',
979: '평상시',
980: '거부',
981: '이었던만큼',
982: '정글',
983: '소중히',
984: '해야하고',
985: '당연하다고',
986: '독립',
987: '있었겠죠',
988: '최고',
989: '씻기도',
990: '변한거',
991: '만나고싶으면',
992: '해봐도',

```

```

1 # 랜덤시드 고정
2 SEED_NUM = 1234
3 tf.random.set_seed(SEED_NUM)

```

```

1 # 파일로드
2 index_inputs = np.load(open(DATA_IN_PATH + TRAIN_INPUTS, 'rb'))
3 index_outputs = np.load(open(DATA_IN_PATH + TRAIN_OUTPUTS, 'rb'))
4 index_targets = np.load(open(DATA_IN_PATH + TRAIN_TARGETS, 'rb'))
5 prepro_configs = json.load(open(DATA_IN_PATH + DATA_CONFIGS, 'r'))

```

```

1 char2idx = prepro_configs['char2idx']
2 end_index = prepro_configs['end_symbol']
3 model_name = 'transformer'
4 vocab_size = prepro_configs['vocab_size']
5 BATCH_SIZE = 2
6 MAX_SEQUENCE = 50
7 EPOCHS = 2
8 VALID_SPLIT = 0.1

```

```

1 kargs = {'model_name': model_name,
2         'num_layers': 2,
3         'd_model': 512, # 512 / 8 = 64
4         'num_heads': 8,
5         'dff': 2048, # 512 * 4
6         'input_vocab_size': vocab_size,
7         'target_vocab_size': vocab_size,
8         'maximum_position_encoding': MAX_SEQUENCE,
9         'end_token_idx': char2idx[end_index],
10        'rate': 0.1
11        }

```

```

1 enc_padding_mask, look_ahead_mask, dec_padding_mask = create_masks(index_inputs, index_outputs)
2 pos_encoding = positional_encoding(50, 512)
3 print (pos_encoding.shape)

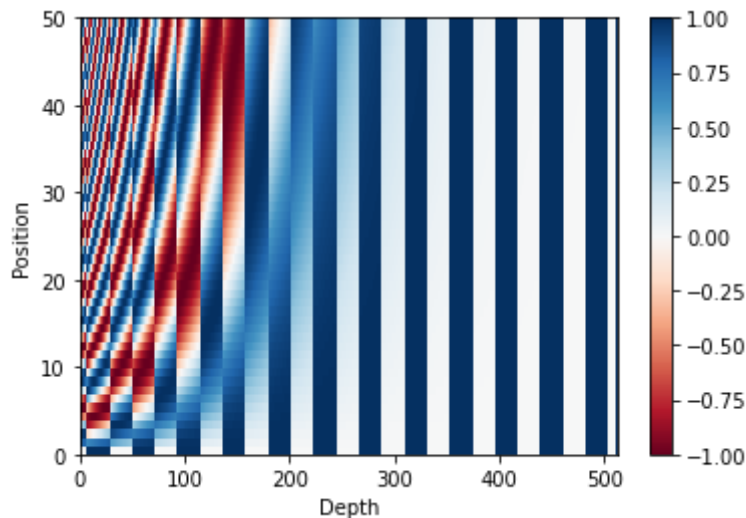
```

```
(1, 50, 512)
```

```

1 plt.pcolormesh(pos_encoding[0], cmap='RdBu')
2 plt.xlabel('Depth')
3 plt.xlim((0, 512))
4 plt.ylabel('Position')
5 plt.colorbar()
6 plt.show()

```



```

1 loss_object = tf.keras.losses.SparseCategoricalCrossentropy(
2     from_logits=True, reduction='none')
3
4 train_accuracy = tf.keras.metrics.SparseCategoricalAccuracy(name='accuracy')

```

```

1 model = Transformer(**kwargs)
2 model.compile(optimizer=tf.keras.optimizers.Adam(1e-4),
3               loss=loss,
4               metrics=[accuracy])

```

```

1 # overfitting을 막기 위한 earlystop 추가
2 earlystop_callback = EarlyStopping(monitor='val_accuracy', min_delta=0.0001, patience=10)
3 # min_delta: the threshold that triggers the termination (acc should at least improve 0.0001)
4 # patience: no improvment epochs (patience = 1, 1번 이상 상승이 없으면 종료)
5
6 checkpoint_path = DATA_OUT_PATH + model_name + '/weights.h5'
7 checkpoint_dir = os.path.dirname(checkpoint_path)
8
9 # Create path if exists
10 if os.path.exists(checkpoint_dir):
11     print("{} -- Folder already exists\n".format(checkpoint_dir))
12 else:
13     os.makedirs(checkpoint_dir, exist_ok=True)

```



```

13 os.makedirs(checkpoint_dir, exist_ok=True)
14 print("{} -- Folder create complete Wn".format(checkpoint_dir))
15
16 cp_callback = ModelCheckpoint(
17     checkpoint_path, monitor='val_accuracy', verbose=1, save_best_only=True, save_weights_only=True)

./data_out/transformer -- Folder already exists

```

```

1 history = model.fit([index_inputs, index_outputs], index_targets,
2                     batch_size=BATCH_SIZE, epochs=EPOCHS,
3                     validation_split=VALID_SPLIT, callbacks=[earlystop_callback, cp_callback])

```

Epoch 1/2

7969/7969 [=====] - 4824s 604ms/step - loss: 0.0493 - accuracy: 0.97

Epoch 00001: val_accuracy improved from -inf to 0.98173, saving model to ./data_out/transformer

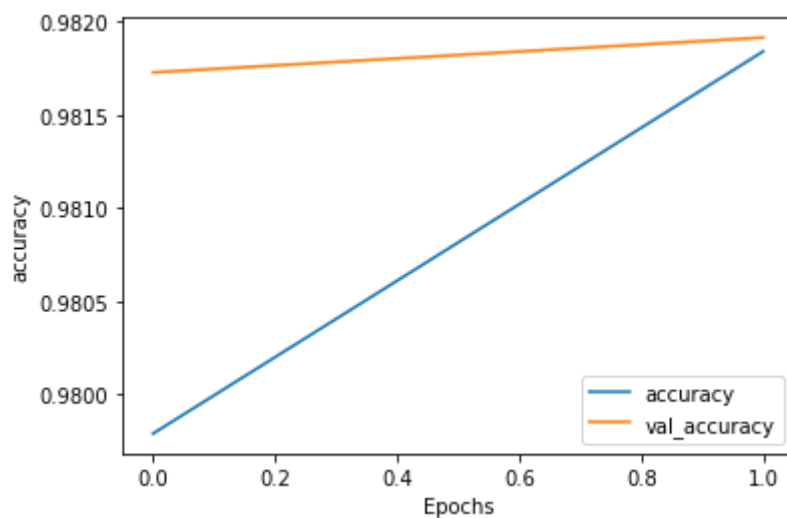
Epoch 2/2

7969/7969 [=====] - 4756s 597ms/step - loss: 0.0438 - accuracy: 0.98

Epoch 00002: val_accuracy improved from 0.98173 to 0.98192, saving model to ./data_out/transformer



```
1 plot_graphs(history, 'accuracy')
```



```
1 plot_graphs(history, 'loss')
```



0043-1

```
1 text = "양성평등"
2 test_index_inputs, _ = enc_processing([text], char2idx)
3 print(test_index_inputs)
4 print(model.inference(test_index_inputs))
5 outputs = model.inference(test_index_inputs)
6
7 print(' '.join([idx2char[str(o)] for o in outputs]))
```

```
[ [3 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
    0 0 0 0 0 0 0 0 0 0 0 0 0 0]]  
[3, 13561, 3, 13561, 3]  
<UNK> - <UNK> - <UNK>
```