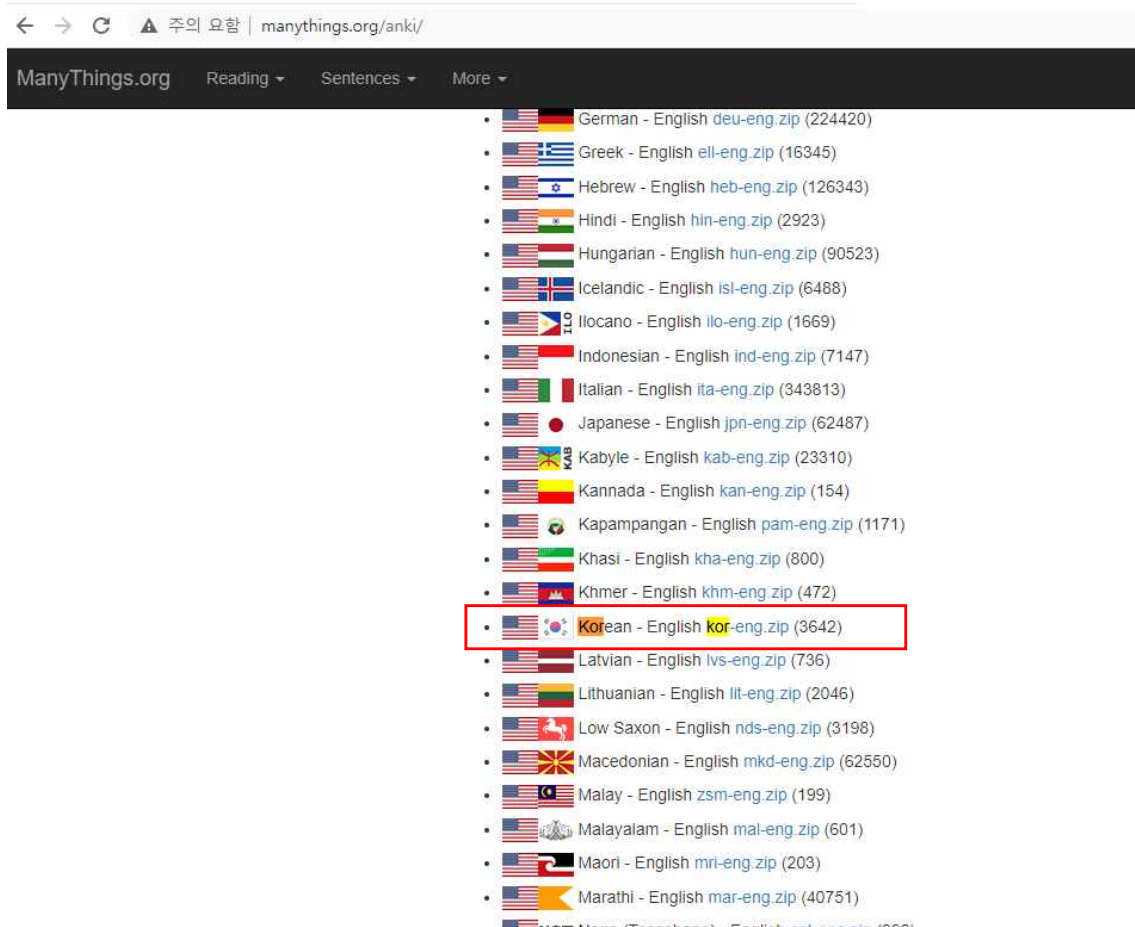


딥러닝 기말과제

20152658 빅데이터경영통계전공 원동욱

1. 데이터 수집

교수님이 주신 데이터와 'manythings' 사이트의 데이터를 활용



2. 전처리

2-1. 외부 내부 공통 데이터 전처리

```
1 def unicode_to_ascii(s):
2     return ''.join(c for c in unicodedata.normalize('NFD', s)
3                     if unicodedata.category(c) != 'Mn')

1 def preprocess_sentence(w):
2     w = unicode_to_ascii(w.lower().strip())

3
4     w = re.sub(r"([?!,])", r" \1 ", w)
5     w = re.sub(r'" " ', " ", w)
6
7     w = re.sub(r"^[^a-zA-Z가-힣]?!\+", " ", w)
8     w = re.sub(r'[-+=, #/₩: ^$@+₩ ※~&% · 』 ₩₩ ₩(₩)₩[₩] ₩<₩> `₩' … 》 ]', ' ', w)
9
10    w = w.strip()
11
12    w = '<s> ' + w + ' </s>'
13    return w
```

<단어와 구두점 사이에 공백을 만들고 특수문자는 정규표현식으로 제거>

<시작을 의미하는 토큰 <s>와 종료를 의미하는 토큰 </s> 추가>

2-2. 외부 데이터 전처리

```
1 encoder_input, decoder_input = [], []
2
3 with open('kor.txt', 'r', encoding = 'utf-8') as lines:
4     for i, line in enumerate(lines):
5
6         # source 데이터와 target 데이터 분리
7         src_line, tar_line, _ = line.strip().split('₩t')
8
9         # source 데이터 전처리
10        src_line_input = [w for w in preprocess_sentence(src_line).split()]
11
12        # target 데이터 전처리
13        tar_line_input = [w for w in preprocess_sentence(tar_line).split()]
14
15        encoder_input.append(src_line_input)
16        decoder_input.append(tar_line_input)
```

<영어 데이터와 한글 데이터를 분리하여 따로 전처리 수행>

```
1 d = [' '.join(decoder_input[i]) for i in range(len(decoder_input))]
2 e = [' '.join(encoder_input[i]) for i in range(len(encoder_input))]
```

```
1 d[:5]
```

```
['<s> 가 . </s>',
 '<s> 안녕 . </s>',
 '<s> 뛰어 ! </s>',
 '<s> 뛰어 . </s>',
 '<s> 누구 ? </s>']
```

```
1 e[:5]
```

```
['<s> go . </s>',
 '<s> hi . </s>',
 '<s> run ! </s>',
 '<s> run . </s>',
 '<s> who ? </s>']
```

<나뉘어진 단어 join을 수행하여 문장으로 구성>

```
1 new
```

	eng	kor
0	<s> go . </s>	<s> 가 . </s>
1	<s> hi . </s>	<s> 안녕 . </s>
2	<s> run ! </s>	<s> 뛰어 ! </s>
3	<s> run . </s>	<s> 뛰어 . </s>
4	<s> who ? </s>	<s> 누구 ? </s>
...
3637	<s> science fiction has undoubtedly been the i...	<s> 공상 과학 소설은 의심의 여지 없이 ㄴ...
3638	<s> i started a new blog . ill do my best not ...	<s> 난 블로그를 시작했어 . 블로그를 초ㄴ...
3639	<s> i think its a shame that some foreign lang...	<s> 몇몇 외국어 선생님이 한 번도 원어...
3640	<s> if someone who doesnt know your background...	<s> 만일 네 사정도 잘 모르는 사람이 원ㄴ...
3641	<s> doubtless there exists in this world preci...	<s> 의심의 여지 없이 세상에는 어떤 남자ㄴ...

<하나의 DataFrame 형성>

2-3. 내부 데이터 전처리

sentences = [
("He came to London by way of Siberia."	"그 사람은 시베리아를 거쳐서 런던에 왔어.")	NaN
("Is there bus service to the airport?"	"공항까지 가는 버스가 있나요?"	NaN
("Japan and South Korea are neighbors."	"일본과 남한은 이웃 국가다.")	NaN
("My favorite singer is Kylie Minogue."	"내가 좋아하는 가수는 카일리 미노그야.")	NaN
	"내가 좋아하는 가수는 카일리 미노그입니다.")	NaN
...		
("I can not answer your question."	"나는 네 질문에 답할 수 없어.")	NaN
("I really hate myself right now."	"바로 지금 난 내 자신이 싫어.")	NaN
("Paris is the capital of France."	"파리는 프랑스의 수도다.")	NaN
("Please take me to the hospital."	"병원으로 데려다 주세요.")	NaN
1	NaN	NaN

<DataFrame 정리>

```

1 def data_modify(data) :
2     df = data.reset_index()
3
4     del df['sentences = []]
5
6     df = df.iloc[:len(df)-1]
7     df.rename(columns = {'level_0' : 'eng', 'level_1' : 'kor'}, inplace = True)
8     return df

```



```

1 d = d[:1500 - len(df)]
2 e = e[:1500 - len(df)]

1 new = pd.DataFrame(d,e).reset_index().rename(columns = {'index' : 'eng', 0 : 'kor'})

1 df = pd.concat([df,new]).reset_index()

1 del df['index']

1 df.head()

```

	eng	kor
0	<s> he came to london by way of siberia . </s>	<s> 그 사람은 시베리아를 거쳐서 런던에 와...
1	<s> is there bus service to the airport ? </s>	<s> 공항까지 가는 버스가 있나요 ? </s>
2	<s> japan and south korea are neighbors . </s>	<s> 일본과 남한은 이웃 국가다 . </s>
3	<s> my favorite singer is kylie minogue . </s>	<s> 내가 좋아하는 가수는 카일리 미노그야
4	<s> my favorite singer is kylie minogue . </s>	<s> 내가 좋아하는 가수는 카일리 미노그입느...

<학습 속도로 인하여 외부 데이터 포함 1500개만 선정>

3. 토큰화

SentencePiece란?

실무에서 바로 적용 가능하고, BPE를 포함한 기타 서브워드 토크나이징 알고리즘들을 내장 또한, 사전 토큰화 작업없이 단어 분리 토큰화를 수행 → 언어에 종속되지 않는다.

3-1. 영어 SentencePiece

```
English

1 with open('eng_translation.txt', 'w', encoding = 'utf8') as f:
2     f.write("\n".join(df['eng']))

1 spm.SentencePieceTrainer.Train('--input=eng_translation.txt --model_prefix=english --vocab_size=1600 --user_defined_symbols=<s>,</s>')
True

1 sp_eng = spm.SentencePieceProcessor()
2 vocab_file = "english.model"
3 sp_eng.load(vocab_file)
True
```

<vocab_size는 1600으로 지정>

<시작·출력 토큰인 <s>,</s>는 나누어지지 않게 정의>

input에 학습시킬 파일, model_prefix는 모델명, vocab_size는 단어 집합의 크기
user_defined_symbols는 사용자 정의 토큰이다.

```
1 print(df['eng'][0])
2 print(sp_eng.encode_as_pieces(df['eng'][0]))
3 print(sp_eng.encode_as_ids(df['eng'][0]))

<s> he came to london by way of siberia . </s>
['_', '<s>', '_he', '_came', '_to', '_lon', 'don', '_', 'by', '_', 'way', '_of', '_si', 'ber', 'i', 'a',
'_', ' ', ' ', ' ', ' ', '</s>']
[3, 1, 47, 219, 21, 812, 701, 3, 205, 3, 105, 77, 316, 189, 51, 44, 3, 4, 3, 2]
```

<예시>

위의 user_defined_symbols에 시작 출력 토큰을 지정해주었기에 나누어지지 않았다.

3-2. 한글 SentencePiece

Kor

```
1 with open('kor_translation.txt','w',encoding = 'utf8') as f:
2     f.write('\n'.join(df['kor']))

1 spm.SentencePieceTrainer.Train('--input=kor_translation.txt --model_prefix=kor --user_defined_symbols=<s>,</s> --vocab_size=2800')

True

1 sp_kor = spm.SentencePieceProcessor()
2 vocab_file = "kor.model"
3 sp_kor.load(vocab_file)
```

<vocab_size는 2800으로 지정>

<시작·출력 토큰인 <s>,</s>는 나누어지지 않게 정의>

```
1 print(df['kor'][900])
2 print(sp_kor.encode_as_pieces(df['kor'][900]))
3 print(sp_kor.encode_as_ids(df['kor'][900]))
```

<s> 우리가 사과할게 . </s>

['_ ', '<s>', '_우리가', '_사과', '할게', '._ ', '_ ', '</s>']
[2147, 1, 257, 290, 993, 3, 2147, 2]

<예시>

한글 또한 마찬가지

3-3. Padding(패딩)

padding

```
] : 1 eng_tokenizer = pad_sequences(eng_tokenizer,padding = 'post')
    2 kor_tokenizer = pad_sequences(kor_tokenizer,padding = 'post')
```

```
] : 1 print('eng_tokenizer의 크기 : {}'.format(eng_tokenizer.shape))
    2 print('kor_tokenizer의 크기 : {}'.format(kor_tokenizer.shape))
```

eng_tokenizer의 크기 : (1500, 28)

kor_tokenizer의 크기 : (1500, 21)

길이가 가장 긴 문서의 길이를 기준으로 패딩

<Padding 후 Shape>

3-4. 데이터 분리

```
1 indices = np.arange(eng_tokenizer.shape[0])
2 np.random.shuffle(indices)

1 encoder_input = eng_tokenizer[indices]
2 decoder_input = kor_tokenizer[indices]

1 encoder_train, encoder_test, decoder_train, decoder_test = train_test_split(encoder_input, decoder_input, test

1 print(decoder_train.shape)
2 print(encoder_train.shape)
3 print(decoder_test.shape)
4 print(encoder_test.shape)

(1200, 21)
(1200, 28)
(300, 21)
(300, 28)
```

<적절한 분포를 갖도록 shuffle 진행 후 Train · Test 8:2 비율>

4. 인코더 및 디코더 모델 작성 및 학습

4-1. Dataset 및 모델 구축

```
1 BUFFER_SIZE = len(encoder_train)
2 BATCH_SIZE = 64
3 steps_per_epoch = len(encoder_train)//BATCH_SIZE
4 embedding_dim = 256
5 units = 1024
6 vocab_inp_size = max([max(l) for l in encoder_input])+1
7 vocab_tar_size = max([max(l) for l in decoder_input])+1
8
9 dataset = tf.data.Dataset.from_tensor_slices((encoder_train, decoder_train)).shuffle(BUFFER_SIZE)
10 dataset = dataset.batch(BATCH_SIZE, drop_remainder=True)
```

```
1 example_input_batch, example_target_batch = next(iter(dataset))
2 example_input_batch.shape, example_target_batch.shape
```

(TensorShape([64, 28]), TensorShape([64, 21]))

4-2 인코더 및 디코더

```
1 class Encoder(tf.keras.Model):
2     def __init__(self, vocab_size, embedding_dim, enc_units, batch_sz):
3         super(Encoder, self).__init__()
4         self.batch_sz = batch_sz
5         self.enc_units = enc_units
6         self.embedding = tf.keras.layers.Embedding(vocab_size, embedding_dim)
7         self.gru = tf.keras.layers.GRU(self.enc_units,
8                                         return_sequences=True,
9                                         return_state=True,
10                                        recurrent_initializer='glorot_uniform')
11
12     def call(self, x, hidden):
13         x = self.embedding(x)
14         output, state = self.gru(x, initial_state = hidden)
15         return output, state
16
17     def initialize_hidden_state(self):
18         return tf.zeros((self.batch_sz, self.enc_units))
```

```
1 encoder = Encoder(vocab_inp_size, embedding_dim, units, BATCH_SIZE)
2
3 # sample input
4 sample_hidden = encoder.initialize_hidden_state()
5 sample_output, sample_hidden = encoder(example_input_batch, sample_hidden)
6 print ('Encoder output shape: (batch size, sequence length, units) {}'.format(sample_output.shape))
7 print ('Encoder Hidden state shape: (batch size, units) {}'.format(sample_hidden.shape))
```

Encoder output shape: (batch size, sequence length, units) (64, 28, 1024)

Encoder Hidden state shape: (batch size, units) (64, 1024)

<인코더 설정>

```

1 class BahdanauAttention(tf.keras.layers.Layer):
2     def __init__(self, units):
3         super(BahdanauAttention, self).__init__()
4         self.W1 = tf.keras.layers.Dense(units)
5         self.W2 = tf.keras.layers.Dense(units)
6         self.V = tf.keras.layers.Dense(1)
7
8     def call(self, query, values):
9         query_with_time_axis = tf.expand_dims(query, 1)
10
11         score = self.V(tf.nn.tanh(self.W1(query_with_time_axis) + self.W2(values)))
12
13         attention_weights = tf.nn.softmax(score, axis=1)
14
15         context_vector = attention_weights * values
16         context_vector = tf.reduce_sum(context_vector, axis=1)
17
18         return context_vector, attention_weights

```

```

1 attention_layer = BahdanauAttention(10)
2 attention_result, attention_weights = attention_layer(sample_hidden, sample_output)
3
4 print("Attention result shape: (batch size, units) {}".format(attention_result.shape))
5 print("Attention weights shape: (batch_size, sequence_length, 1) {}".format(attention_weights.shape))

```

Attention result shape: (batch size, units) (64, 1024)

Attention weights shape: (batch_size, sequence_length, 1) (64, 28, 1)

<바다나우 어텐션 정의>

어텐션 함수의 Query가 디코더 셀의 t 시점의 은닉 상태가 아니라 t-1시점의 은닉 상태이다.

```

1 class Decoder(tf.keras.Model):
2     def __init__(self, vocab_size, embedding_dim, dec_units, batch_sz):
3         super(Decoder, self).__init__()
4         self.batch_sz = batch_sz
5         self.dec_units = dec_units
6         self.embedding = tf.keras.layers.Embedding(vocab_size, embedding_dim)
7         self.gru = tf.keras.layers.GRU(self.dec_units,
8                                         return_sequences=True,
9                                         return_state=True,
10                                         recurrent_initializer='glorot_uniform')
11
12         self.fc = tf.keras.layers.Dense(vocab_size)
13
14         self.attention = BahdanauAttention(self.dec_units)
15
16     def call(self, x, hidden, enc_output):
17
18         context_vector, attention_weights = self.attention(hidden, enc_output)
19
20         x = self.embedding(x)
21
22         x = tf.concat([tf.expand_dims(context_vector, 1), x], axis=-1)
23
24         output, state = self.gru(x)
25
26         output = tf.reshape(output, (-1, output.shape[2]))
27
28         x = self.fc(output)
29
30         return x, state, attention_weights

```

```

1 decoder = Decoder(vocab_tar_size, embedding_dim, units, BATCH_SIZE)
2
3 sample_decoder_output, _, _ = decoder(tf.random.uniform((BATCH_SIZE, 1)),
4                                       sample_hidden, sample_output)
5
6 print('Decoder output shape: (batch_size, vocab size) {}'.format(sample_decoder_output.shape))

```

Decoder output shape: (batch_size, vocab size) (64, 2799)

<디코더 설정>

4-3 모델

```
1 optimizer = tf.keras.optimizers.Adam()
2 loss_object = tf.keras.losses.SparseCategoricalCrossentropy(
3     from_logits=True, reduction='none')
4
5 def loss_function(real, pred):
6     mask = tf.math.logical_not(tf.math.equal(real, 0))
7     loss_ = loss_object(real, pred)
8
9     mask = tf.cast(mask, dtype=loss_.dtype)
10    loss_ *= mask
11
12    return tf.reduce_mean(loss_)
```

```
1 checkpoint_dir = './training_checkpoints'
2 checkpoint_prefix = os.path.join(checkpoint_dir, "ckpt")
3 checkpoint = tf.train.Checkpoint(optimizer=optimizer,
4                                   encoder=encoder,
5                                   decoder=decoder)
```

<옵티마이저 및 손실 함수 정의>

```
1 @tf.function
2 def train_step(inp, targ, enc_hidden):
3     loss = 0
4     with tf.GradientTape() as tape:
5         enc_output, enc_hidden = encoder(inp, enc_hidden)
6
7         dec_hidden = enc_hidden
8
9         dec_input = tf.expand_dims([sp_kor.piece_to_id('</s>')] * BATCH_SIZE, 1)
10
11         for t in range(1, targ.shape[1]):
12             predictions, dec_hidden, _ = decoder(dec_input, dec_hidden, enc_output)
13
14             loss += loss_function(targ[:, t], predictions)
15
16         # using teacher forcing
17         dec_input = tf.expand_dims(targ[:, t], 1)
18
19     batch_loss = (loss / int(targ.shape[1]))
20
21     variables = encoder.trainable_variables + decoder.trainable_variables
22
23     gradients = tape.gradient(loss, variables)
24
25     optimizer.apply_gradients(zip(gradients, variables))
26
27     return batch_loss
```

<체크 포인트>

4-4 학습

```
1 EPOCHS = 30
2
3 for epoch in range(EPOCHS):
4     start = time.time()
5
6     enc_hidden = encoder.initialize_hidden_state()
7     total_loss = 0
8
9     for (batch, (inp, targ)) in enumerate(dataset.take(steps_per_epoch)):
10        batch_loss = train_step(inp, targ, enc_hidden)
11        total_loss += batch_loss
12
13        if batch % 100 == 0:
14            print('Epoch {} Batch {} Loss {:.4f}'.format(epoch + 1,
15                                                         batch,
16                                                         batch_loss.numpy()))
17
18
19    print('Epoch {} Loss {:.4f}'.format(epoch + 1, total_loss / steps_per_epoch))
20
21    print('Time taken for 1 epoch {} sec\n'.format(time.time() - start))
```

```
Epoch 27 Batch 0 Loss 0.2151
Epoch 27 Loss 0.2261
Time taken for 1 epoch 32.2049036026001 sec
```

```
Epoch 28 Batch 0 Loss 0.1804
Epoch 28 Loss 0.2018
Time taken for 1 epoch 32.409356117248535 sec
```

```
Epoch 29 Batch 0 Loss 0.1686
Epoch 29 Loss 0.1767
Time taken for 1 epoch 32.050315380096436 sec
```

```
Epoch 30 Batch 0 Loss 0.1299
Epoch 30 Loss 0.1598
Time taken for 1 epoch 32.92041778564453 sec
```

<Loss 값 : 0.1299>

5. 평가

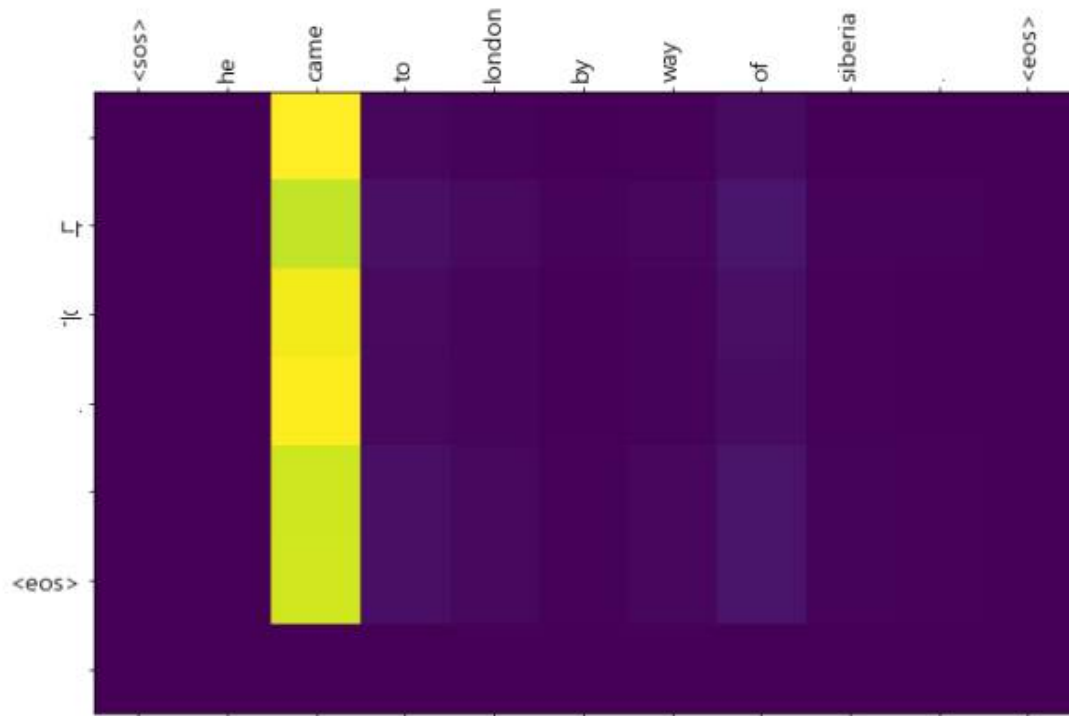
```
1 def evaluate(sentence):
2
3     max_length_targ = decoder_input.shape[1]
4     max_length_inp = encoder_input.shape[1]
5
6     attention_plot = np.zeros((max_length_targ, max_length_inp))
7
8     sentence = preprocess_sentence(sentence)
9
10    inputs = sp_eng.encode_as_ids(sentence)
11    inputs = tf.keras.preprocessing.sequence.pad_sequences([inputs],
12                                                         maxlen=max_length_inp,
13                                                         padding='post')
14    inputs = tf.convert_to_tensor(inputs)
15
16    result = ''
17
18    hidden = [tf.zeros((1, units))]
19    enc_out, enc_hidden = encoder(inputs, hidden)
20
21    dec_hidden = enc_hidden
22    dec_input = tf.expand_dims([sp_kor.piece_to_id('<sos>')], 0)
23
24    for t in range(max_length_targ):
25        predictions, dec_hidden, attention_weights = decoder(dec_input,
26                                                         dec_hidden,
27                                                         enc_out)
28
29        # storing the attention weights to plot later on
30        attention_weights = tf.reshape(attention_weights, (-1, ))
31        attention_plot[t] = attention_weights.numpy()
32
33        predicted_id = tf.argmax(predictions[0]).numpy()
34
35        result += sp_kor.decode_ids([int(predicted_id)]) + ' '
36
37        if sp_kor.decode_ids([int(predicted_id)]) == '<eos>':
38            return result, sentence, attention_plot
39
40    # the predicted ID is fed back into the model
41    dec_input = tf.expand_dims([predicted_id], 0)
42
43    return result, sentence, attention_plot
```

```
1 def plot_attention(attention, sentence, predicted_sentence):
2     fig = plt.figure(figsize=(10,10))
3     ax = fig.add_subplot(1, 1, 1)
4     ax.matshow(attention, cmap='viridis')
5
6     fontdict = {'fontsize': 14}
7
8     ax.set_xticklabels([''] + sentence, fontdict=fontdict, rotation=90)
9     ax.set_yticklabels([''] + predicted_sentence, fontdict=fontdict)
10
11    ax.xaxis.set_major_locator(ticker.MultipleLocator(1))
12    ax.yaxis.set_major_locator(ticker.MultipleLocator(1))
13
14    plt.show()
```

```
1 def translate(sentence):
2     result, sentence, attention_plot = evaluate(sentence)
3
4     print('Input: %s' % (sentence))
5     print('Predicted translation: {}'.format(result))
6
7     attention_plot = attention_plot[:len(result.split(' ')), :len(sentence.split(' '))]
8     plot_attention(attention_plot, sentence.split(' '), result.split(' '))
```

5-1. 모든 데이터를 가지고 한 모델링 예제

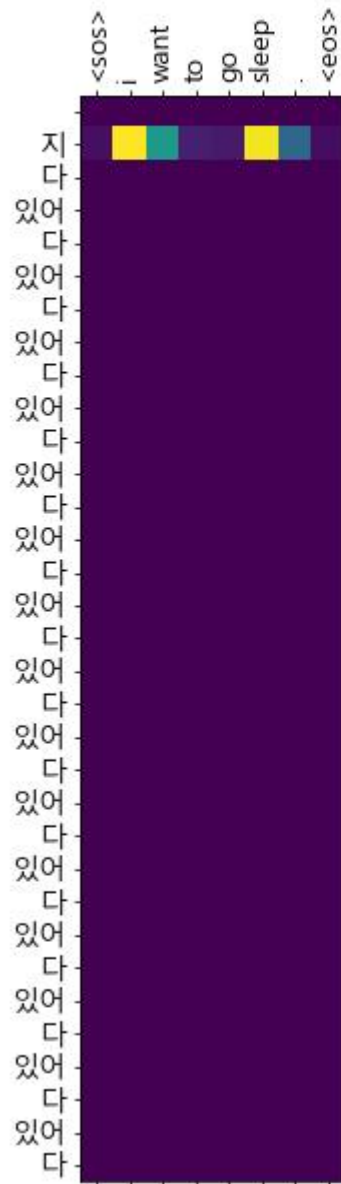
Input: <sos> he came to london by way of siberia . <eos>
 Predicted translation: 다 온 . 는 <eos>



```
1 translate('i want to go sleep.')
```

Input: <eos> i want to go sleep . <eos>

Predicted translation: 지 다 있어 다 있어 다
다 있어 다 있어 다 있어 다 있어 다 있어 다 오

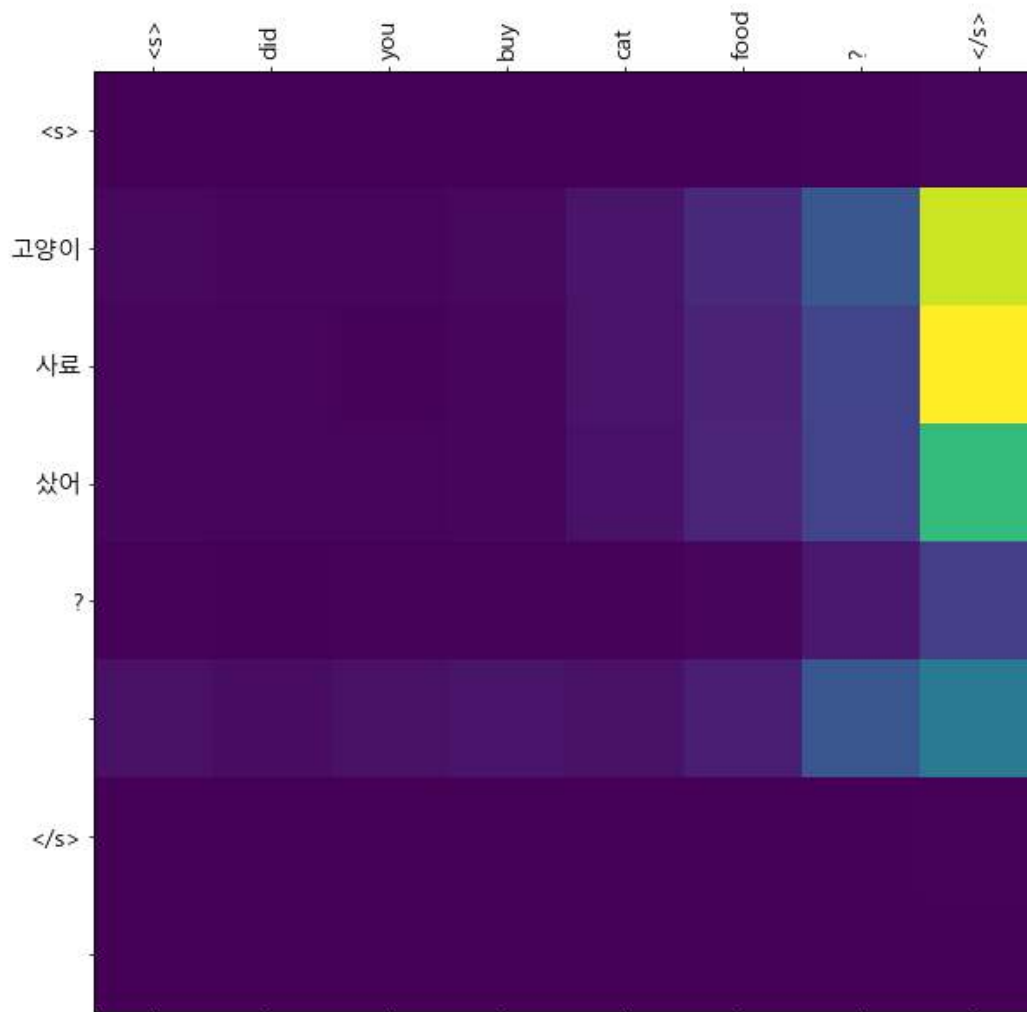


5-2. 잘 된 예제

```
1 translate(u'did you buy cat food ?')
```

Input: <s> did you buy cat food ? </s>

Predicted translation: <s> 고양이 사료 샀어 ? </s>

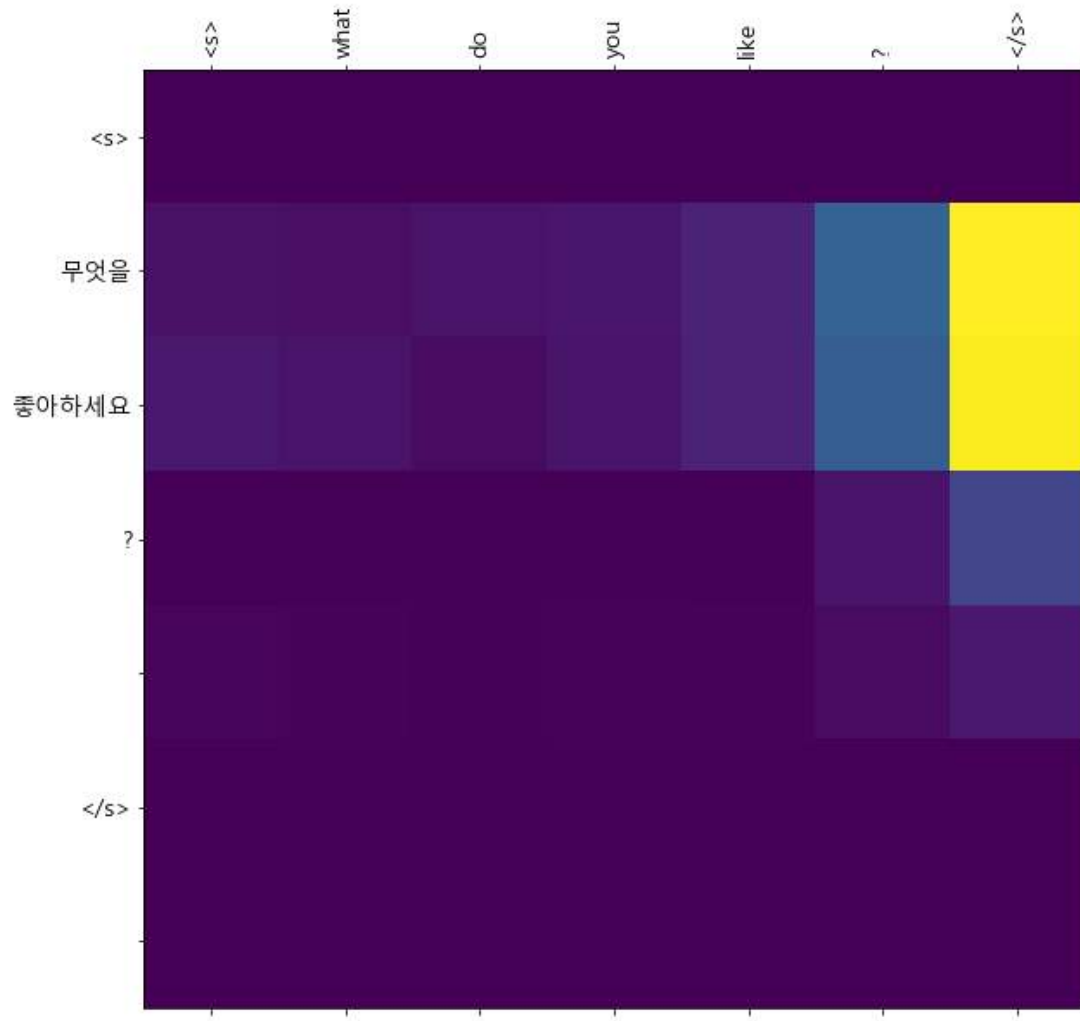


<복원 테스트>

```
1 translate(u'what do you like ?')
```

Input: <s> what do you like ? </s>

Predicted translation: <s> 무엇을 좋아하세요 ? </s>



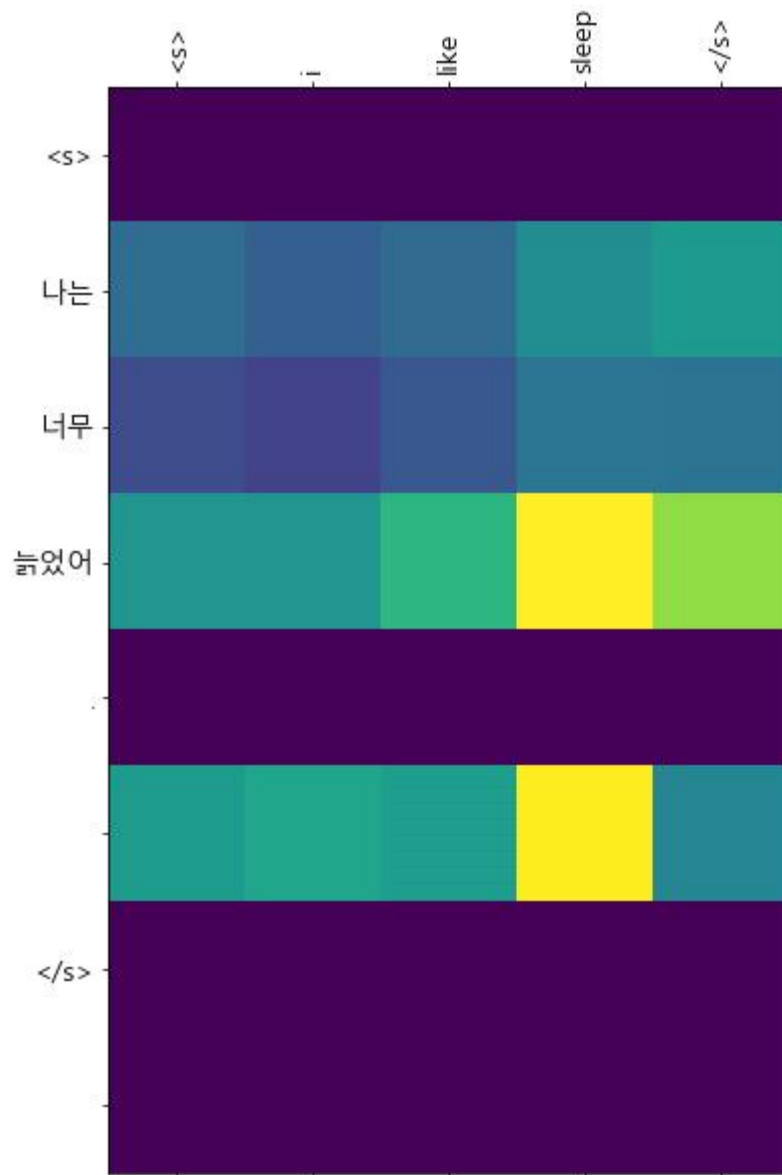
<복원 테스트>

5-3. 번역 오류

```
1 translate(u' i like sleep |')
```

Input: <s> i like sleep </s>

Predicted translation: <s> 나는 너무 늙었어 . </s>



```
1 translate(u' do spiders scare you ? ')
```

Input: <s> do spiders scare you ? </s>

Predicted translation: <s> 노래하는 거 좋아해 ? </s>

