

자연어처리_과제#6

1. 강의자료 12장에서 다룬 프로그램은 인터넷 자료(wikidocs.net/book/2155, 딥 러닝을 이용한 자연어 처리 입문)의 14-02절에서 찾을 수 있다. 인터넷 자료에 나온 순서대로 프로그램을 구현하여 12장 강의 자료에서 제시한 것과 유사하게 10개의 샘플을 수행하여 결과를 제시하라. 10개의 샘플은 자신의 학번 마지막 2 자리를 이용하여 생성한 숫자를 이용하여 선택한다. 예를들면 자신의 학번이 201803371 이면 [171, 271, ..., 1071] 등 10개의 숫자를 제시한다.

프로그램 구현

```
import os
import re
import shutil
import zipfile

import numpy as np
import pandas as pd
import tensorflow as tf
import unicodedata
import urllib3
from tensorflow.keras.layers import Embedding, GRU, Dense
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.preprocessing.text import Tokenizer

import requests

headers = {
    'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64)
AppleWebKit/537.36 (KHTML, like Gecko) Chrome/91.0.4472.124 Safari/537.36'
}

# ZIP 파일 다운로드
def download_zip(url, output_path):
    response = requests.get(url, headers=headers, stream=True)
    if response.status_code == 200:
        with open(output_path, 'wb') as f:
            for chunk in response.iter_content(chunk_size=8192):
                f.write(chunk)
        print(f"ZIP file downloaded to {output_path}")
    else:
        print(f"Failed to download. HTTP Response Code:
{response.status_code}")

url = "http://www.manythings.org/anki/fra-eng.zip"
output_path = "fra-eng.zip"
download_zip(url, output_path)

path = os.getcwd()
zipfilename = os.path.join(path, output_path)
```

```

with zipfile.ZipFile(zipfilename, 'r') as zip_ref:
    zip_ref.extractall(path)

num_samples = 33000

def to_ascii(s):
    # 프랑스어 악센트(accent) 삭제
    # 예시 : 'déjà diné' -> déjà dine
    return ''.join(c for c in unicodedata.normalize('NFD', s)
                    if unicodedata.category(c) != 'Mn')

def preprocess_sentence(sent):
    # 악센트 제거 함수 호출
    sent = to_ascii(sent.lower())

    # 단어와 구두점 사이에 공백 추가.
    # ex) "I am a student." => "I am a student ."
    sent = re.sub(r"([?.!,&])", r" \1", sent)

    # (a-z, A-Z, ".", "?", "!", ",",) 이들을 제외하고는 전부 공백으로 변환.
    sent = re.sub(r"[^a-zA-Z!.,?]+", r" ", sent)

    # 다수 개의 공백을 하나의 공백으로 치환
    sent = re.sub(r"\s+", " ", sent)
    return sent

```

```

# 전처리 테스트
en_sent = u"Have you had dinner?"
fr_sent = u"Avez-vous déjà diné?"

print('전처리 전 영어 문장 :', en_sent)
print('전처리 후 영어 문장 :', preprocess_sentence(en_sent))
print('전처리 전 프랑스어 문장 :', fr_sent)
print('전처리 후 프랑스어 문장 :', preprocess_sentence(fr_sent))

```

```

전처리 전 영어 문장 : Have you had dinner?
전처리 후 영어 문장 : have you had dinner ?
전처리 전 프랑스어 문장 : Avez-vous déjà diné?
전처리 후 프랑스어 문장 : avez vous deja dine ?

```

```

def load_preprocessed_data():
    encoder_input, decoder_input, decoder_target = [], [], []

    with open("fra.txt", "r") as lines:
        for i, line in enumerate(lines):
            # source 데이터와 target 데이터 분리
            src_line, tar_line, _ = line.strip().split('\t')

```

```

# source 데이터 전처리
src_line = [w for w in preprocess_sentence(src_line).split()]

# target 데이터 전처리
tar_line = preprocess_sentence(tar_line)
tar_line_in = [w for w in ("<sos> " + tar_line).split()]
tar_line_out = [w for w in (tar_line + " <eos>").split()]

encoder_input.append(src_line)
decoder_input.append(tar_line_in)
decoder_target.append(tar_line_out)

if i == num_samples - 1:
    break

return encoder_input, decoder_input, decoder_target

sents_en_in, sents_fra_in, sents_fra_out = load_preprocessed_data()
print('인코더의 입력 :', sents_en_in[:5])
print('디코더의 입력 :', sents_fra_in[:5])
print('디코더의 레이블 :', sents_fra_out[:5])

```

```

인코더의 입력 : [['go', '.'], ['go', '.'], ['go', '.'], ['go', '.'], ['hi',
'.']]
디코더의 입력 : [['<sos>', 'va', '!'], ['<sos>', 'marche', '.'], ['<sos>',
'en', 'route', '!'], ['<sos>', 'bouge', '!'], ['<sos>', 'salut', '!']]
디코더의 레이블 : [['va', '!', '<eos>'], ['marche', '.', '<eos>'], ['en',
'route', '!', '<eos>'], ['bouge', '!', '<eos>'], ['salut', '!', '<eos>']]

```

```

tokenizer_en = Tokenizer(filters="", lower=False)
tokenizer_en.fit_on_texts(sents_en_in)
encoder_input = tokenizer_en.texts_to_sequences(sents_en_in)
encoder_input = pad_sequences(encoder_input, padding="post")

tokenizer_fra = Tokenizer(filters="", lower=False)
tokenizer_fra.fit_on_texts(sents_fra_in)
tokenizer_fra.fit_on_texts(sents_fra_out)

decoder_input = tokenizer_fra.texts_to_sequences(sents_fra_in)
decoder_input = pad_sequences(decoder_input, padding="post")

decoder_target = tokenizer_fra.texts_to_sequences(sents_fra_out)
decoder_target = pad_sequences(decoder_target, padding="post")

print('인코더의 입력의 크기(shape) :', encoder_input.shape)
print('디코더의 입력의 크기(shape) :', decoder_input.shape)
print('디코더의 레이블의 크기(shape) :', decoder_target.shape)

```

```
인코더의 입력의 크기(shape) : (33000, 7)
디코더의 입력의 크기(shape) : (33000, 16)
디코더의 레이블의 크기(shape) : (33000, 16)
```

```
src_vocab_size = len(tokenizer_en.word_index) + 1
tar_vocab_size = len(tokenizer_fra.word_index) + 1
print("영어 단어 집합의 크기 : {:d}, 프랑스어 단어 집합의 크기 :
{:d}".format(src_vocab_size, tar_vocab_size))
```

```
영어 단어 집합의 크기 : 4485, 프랑스어 단어 집합의 크기 : 7878
```

```
src_to_index = tokenizer_en.word_index
index_to_src = tokenizer_en.index_word
tar_to_index = tokenizer_fra.word_index
index_to_tar = tokenizer_fra.index_word

indices = np.arange(encoder_input.shape[0])
np.random.shuffle(indices)
print('랜덤 시퀀스 :',indices)
```

```
랜덤 시퀀스 : [ 8402 27938 30789 ... 3752 449 16217]
```

```
encoder_input = encoder_input[indices]
decoder_input = decoder_input[indices]
decoder_target = decoder_target[indices]

encoder_input[30997]
# >>> array([ 6, 84, 13, 106, 1, 0, 0], dtype=int32)

decoder_input[30997]
# >>> array([ 2, 26, 27, 5, 11, 908, 1, 0, 0, 0, 0, 0,
0, 0, 0, 0], dtype=int32)

decoder_target[30997]
# >>> array([ 26, 27, 5, 11, 908, 1, 3, 0, 0, 0, 0, 0,
0, 0, 0, 0], dtype=int32)

n_of_val = int(33000*0.1)
print('검증 데이터의 개수 :',n_of_val)
```

검증 데이터의 개수 : 3300

```
encoder_input_train = encoder_input[:-n_of_val]
decoder_input_train = decoder_input[:-n_of_val]
decoder_target_train = decoder_target[:-n_of_val]

encoder_input_test = encoder_input[-n_of_val:]
decoder_input_test = decoder_input[-n_of_val:]
decoder_target_test = decoder_target[-n_of_val:]

print('훈련 source 데이터의 크기 :', encoder_input_train.shape)
print('훈련 target 데이터의 크기 :', decoder_input_train.shape)
print('훈련 target 레이블의 크기 :', decoder_target_train.shape)
print('테스트 source 데이터의 크기 :', encoder_input_test.shape)
print('테스트 target 데이터의 크기 :', decoder_input_test.shape)
print('테스트 target 레이블의 크기 :', decoder_target_test.shape)
```

```
훈련 source 데이터의 크기 : (29700, 7)
훈련 target 데이터의 크기 : (29700, 16)
훈련 target 레이블의 크기 : (29700, 16)
테스트 source 데이터의 크기 : (3300, 7)
테스트 target 데이터의 크기 : (3300, 16)
테스트 target 레이블의 크기 : (3300, 16)
```

```
from tensorflow.keras.layers import Input, LSTM, Embedding, Dense, Masking
from tensorflow.keras.models import Model

embedding_dim = 64
hidden_units = 64

# 인코더
encoder_inputs = Input(shape=(None,))
enc_emb = Embedding(src_vocab_size, embedding_dim)(encoder_inputs) # 임베딩 층
enc_masking = Masking(mask_value=0.0)(enc_emb) # 패딩 0은 연산에서 제외
encoder_lstm = LSTM(hidden_units, return_state=True) # 상태값 리턴을 위해 return_state는 True
encoder_outputs, state_h, state_c = encoder_lstm(enc_masking) # 은닉 상태와 셀 상태를 리턴
encoder_states = [state_h, state_c] # 인코더의 은닉 상태와 셀 상태를 저장

# 디코더
decoder_inputs = Input(shape=(None,))
dec_emb_layer = Embedding(tar_vocab_size, hidden_units) # 임베딩 층
dec_emb = dec_emb_layer(decoder_inputs) # 패딩 0은 연산에서 제외
dec_masking = Masking(mask_value=0.0)(dec_emb)
```

```

# 상태값 리턴을 위해 return_state는 True, 모든 시점에 대해서 단어를 예측하기 위해
return_sequences는 True
decoder_lstm = LSTM(hidden_units, return_sequences=True,
return_state=True)

# 인코더의 은닉 상태를 초기 은닉 상태(initial_state)로 사용
decoder_outputs, _, _ = decoder_lstm(dec_masking,
initial_state=encoder_states)

# 모든 시점의 결과에 대해서 소프트맥스 함수를 사용한 출력층을 통해 단어 예측
decoder_dense = Dense(tar_vocab_size, activation='softmax')
decoder_outputs = decoder_dense(decoder_outputs)

# 모델의 입력과 출력을 정의.
model = Model([encoder_inputs, decoder_inputs], decoder_outputs)

model.compile(optimizer='adam', loss='sparse_categorical_crossentropy',
metrics=['acc'])

model.fit(x=[encoder_input_train, decoder_input_train],
y=decoder_target_train, \
validation_data=([encoder_input_test, decoder_input_test],
decoder_target_test),
batch_size=128, epochs=50)

```

```

Epoch 1/50
233/233 _____ 13s 28ms/step - acc: 0.6049 - loss: 5.0482 -
val_acc: 0.6199 - val_loss: 1.9874
Epoch 2/50
233/233 _____ 6s 27ms/step - acc: 0.6657 - loss: 1.8791 -
val_acc: 0.7404 - val_loss: 1.6994
Epoch 3/50
233/233 _____ 6s 26ms/step - acc: 0.7431 - loss: 1.6459 -
val_acc: 0.7523 - val_loss: 1.5686
Epoch 4/50
233/233 _____ 6s 27ms/step - acc: 0.9190 - loss: 0.3350 -
val_acc: 0.8713 - val_loss: 0.7310
...
... (생략)
...
Epoch 48/50
233/233 _____ 10s 28ms/step - acc: 0.9234 - loss: 0.3134 -
val_acc: 0.8719 - val_loss: 0.7308
Epoch 49/50
233/233 _____ 10s 29ms/step - acc: 0.9248 - loss: 0.3089 -
val_acc: 0.8714 - val_loss: 0.7299
Epoch 50/50
233/233 _____ 10s 29ms/step - acc: 0.9261 - loss: 0.3028 -
val_acc: 0.8721 - val_loss: 0.7309
<keras.src.callbacks.history.History at 0x7899bb065150>

```

```

# 인코더
encoder_model = Model(encoder_inputs, encoder_states)

# 디코더 설계 시작
# 이전 시점의 상태를 보관할 텐서
decoder_state_input_h = Input(shape=(hidden_units,))
decoder_state_input_c = Input(shape=(hidden_units,))
decoder_states_inputs = [decoder_state_input_h, decoder_state_input_c]

# 훈련 때 사용했던 임베딩 층을 재사용
dec_emb2 = dec_emb_layer(decoder_inputs)

# 다음 단어 예측을 위해 이전 시점의 상태를 현 시점의 초기 상태로 사용
decoder_outputs2, state_h2, state_c2 = decoder_lstm(dec_emb2,
initial_state=decoder_states_inputs)
decoder_states2 = [state_h2, state_c2]

# 모든 시점에 대해서 단어 예측
decoder_outputs2 = decoder_dense(decoder_outputs2)

# 수정된 디코더
decoder_model = Model(
    [decoder_inputs] + decoder_states_inputs,
    [decoder_outputs2] + decoder_states2)

def decode_sequence(input_seq):
    # 입력으로부터 인코더의 마지막 시점의 상태(은닉 상태, 셀 상태)를 얻음
    states_value = encoder_model.predict(input_seq)

    # <SOS>에 해당하는 정수 생성
    target_seq = np.zeros((1,1))
    target_seq[0, 0] = tar_to_index['<sos>']

    stop_condition = False
    decoded_sentence = ''

    # stop_condition이 True가 될 때까지 루프 반복
    # 구현의 간소화를 위해서 이 함수는 배치 크기를 1로 가정합니다.
    while not stop_condition:
        # 이점 시점의 상태 states_value를 현 시점의 초기 상태로 사용
        output_tokens, h, c = decoder_model.predict([target_seq] +
states_value)

        # 예측 결과를 단어로 변환
        sampled_token_index = np.argmax(output_tokens[0, -1, :])
        sampled_char = index_to_tar[sampled_token_index]

        # 현재 시점의 예측 단어를 예측 문장에 추가
        decoded_sentence += ' ' + sampled_char

        # <eos>에 도달하거나 정해진 길이를 넘으면 중단.

```

```

    if (sampled_char == '<eos>' or
        len(decoded_sentence) > 50):
        stop_condition = True

    # 현재 시점의 예측 결과를 다음 시점의 입력으로 사용하기 위해 저장
    target_seq = np.zeros((1,1))
    target_seq[0, 0] = sampled_token_index

    # 현재 시점의 상태를 다음 시점의 상태로 사용하기 위해 저장
    states_value = [h, c]

    return decoded_sentence

# 원문의 정수 시퀀스를 텍스트 시퀀스로 변환
def seq_to_src(input_seq):
    sentence = ''
    for encoded_word in input_seq:
        if(encoded_word != 0):
            sentence = sentence + index_to_src[encoded_word] + ' '
    return sentence

# 번역문의 정수 시퀀스를 텍스트 시퀀스로 변환
def seq_to_tar(input_seq):
    sentence = ''
    for encoded_word in input_seq:
        if(encoded_word != 0 and encoded_word != tar_to_index['<sos>'] and
           encoded_word != tar_to_index['<eos>']):
            sentence = sentence + index_to_tar[encoded_word] + ' '
    return sentence

```

10개의 샘플을 수행

```

for seq_index in [195, 295, 395, 495, 595, 695, 795, 895, 995, 1095]:
    input_seq = encoder_input_train[seq_index: seq_index + 1]
    decoded_sentence = decode_sequence(input_seq)

    print("입력문장 :", seq_to_src(encoder_input_train[seq_index]))
    print("정답문장 :", seq_to_tar(decoder_input_train[seq_index]))
    print("번역문장 :", decoded_sentence[1:-5])
    print("-"*50)

```

```

1/1 _____ 0s 23ms/step
1/1 _____ 0s 27ms/step
1/1 _____ 0s 19ms/step
1/1 _____ 0s 19ms/step
1/1 _____ 0s 17ms/step
1/1 _____ 0s 17ms/step

```


1/1 _____ 0s 20ms/step

1/1 _____ 0s 17ms/step

입력문장 : what a fool i was !

정답문장 : quel idiot j ai ete !

번역문장 : quel idiot j ai ete !

1/1 _____ 0s 17ms/step

1/1 _____ 0s 17ms/step

1/1 _____ 0s 21ms/step

1/1 _____ 0s 20ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 20ms/step

입력문장 : did you hear that ?

정답문장 : avez vous entendu cela ?

번역문장 : l as tu entendu ?

1/1 _____ 0s 23ms/step

1/1 _____ 0s 18ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 20ms/step

1/1 _____ 0s 18ms/step

1/1 _____ 0s 18ms/step

입력문장 : i m shaken .

정답문장 : je suis remue .

번역문장 : je suis sidere .

1/1 _____ 0s 17ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 20ms/step

1/1 _____ 0s 22ms/step

1/1 _____ 0s 22ms/step

1/1 _____ 0s 20ms/step

1/1 _____ 0s 19ms/step

입력문장 : i was invited .

정답문장 : j ai ete invitee .

번역문장 : j ai ete invitee .

1/1 _____ 0s 18ms/step

1/1 _____ 0s 21ms/step

1/1 _____ 0s 25ms/step

1/1 _____ 0s 18ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 19ms/step

입력문장 : get started .

정답문장 : commence .

번역문장 : continuez de nouveau .

1/1 _____ 0s 16ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 18ms/step

1/1 _____ 0s 18ms/step

1/1 _____ 0s 18ms/step

입력문장 : i ate some cheese .

정답문장 : j ai mange du fromage .

번역문장 : j ai mange un fleurs .

1/1 _____ 0s 16ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 28ms/step

1/1 _____ 0s 22ms/step

1/1 _____ 0s 22ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 19ms/step

입력문장 : when can we eat ?

정답문장 : quand pouvons nous manger ?

번역문장 : quand pouvons nous manger ?

1/1 _____ 0s 19ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 25ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 18ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 21ms/step

입력문장 : he tricked me .

정답문장 : il m a arnaquee .

번역문장 : il m a eue .

1/1 _____ 0s 20ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 23ms/step

1/1 _____ 0s 26ms/step

1/1 _____ 0s 20ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 19ms/step

입력문장 : take care .

정답문장 : prends soin de toi .

번역문장 : prends soin de toi .

1/1 _____ 0s 22ms/step

1/1 _____ 0s 20ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 18ms/step

1/1 _____ 0s 19ms/step

1/1 _____ 0s 18ms/step

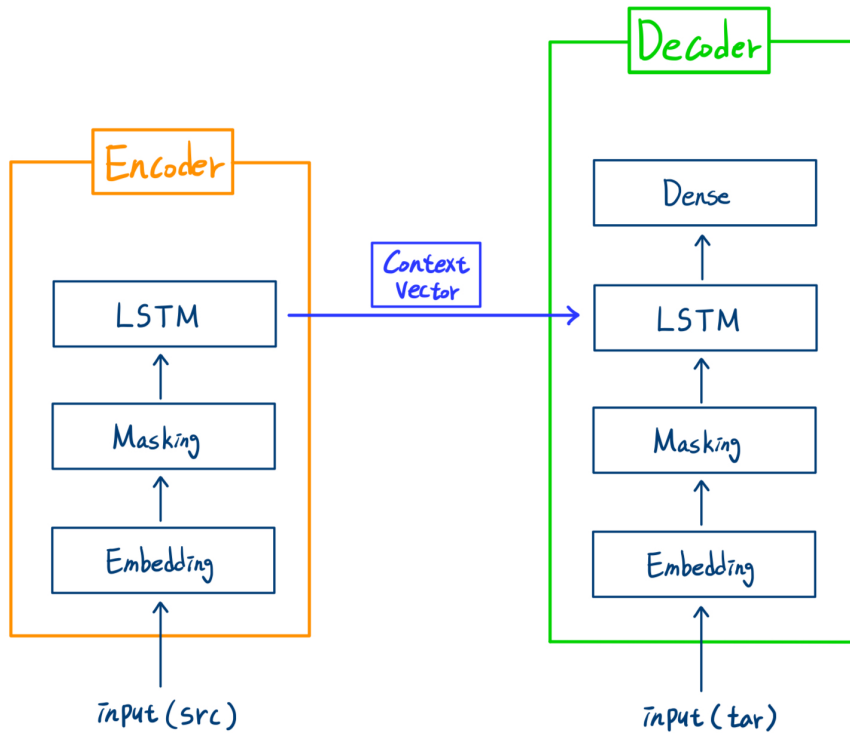
1/1 _____ 0s 20ms/step

입력문장 : he doesn t want it .

정답문장 : il n en veut pas .

번역문장 : il n est pas la .

2. 이 프로그램에서 구현한 **encoder**와 **decoder**가 어떤 구조를 가지는지 그림으로 나타내라. 각 부분에서 필요한 파라미터 개수가 어떻게 되는지 계산하라.

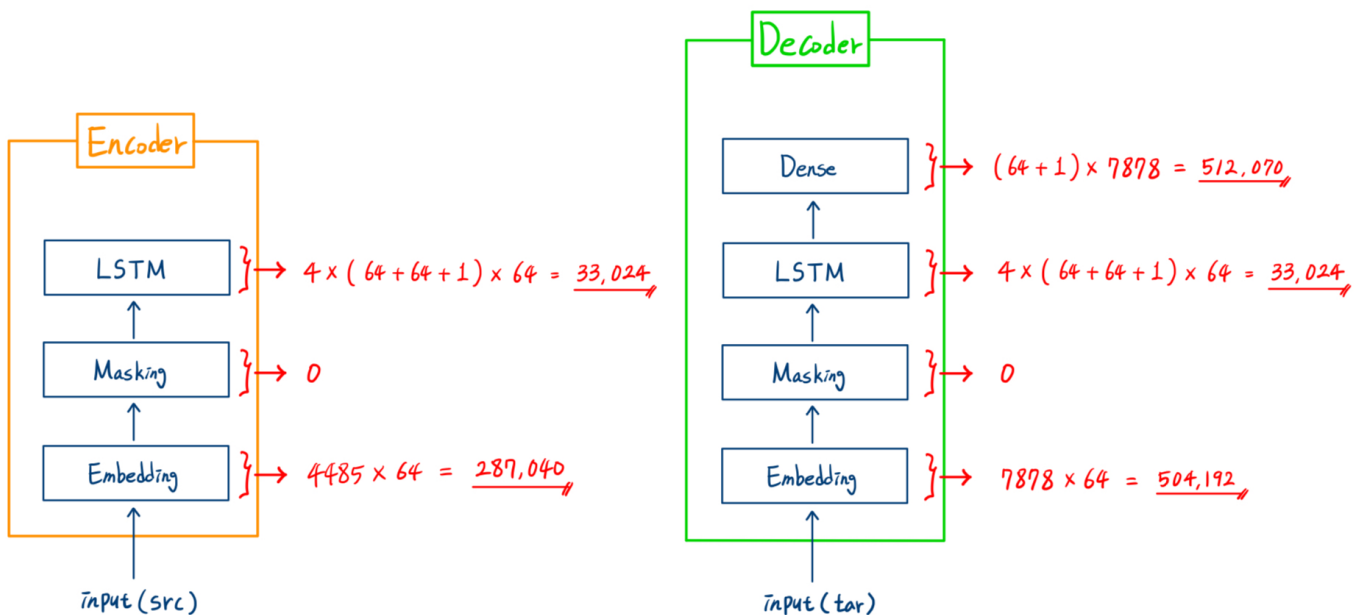


(src) 영어 단어 집합의 크기 : 4485

(tar) 프랑스어 단어 집합의 크기 : 7878

embedding_dim = 64

hidden_units = 64



3. 이 프로그램에서 **context vector**가 어떤 형태로 표시되는지 설명하라. 프로그램에서 **context vector**에 해당하는 변수가 어떤 것이고 어떻게 전달되는지 설명하라.

context vector란 인코더-디코더 구조에서 인코더의 최종 출력인 **은닉상태**와 **셀상태**를 말하는데, 이는 인코더가 입력 시퀀스를 요약한 정보로, 디코더의 초기 상태로 전달되어 디코더가 번역 문장을 생성하는 데 사용된다.

이 프로그램에서는 인코더 내 LSTM 레이어의 출력값 중 **은닉상태**를 **state_h**에 저장하고, **셀상태**를 **state_c**에 저장했다.

이 프로그램에서 context vector의 은닉상태와 셀상태는 둘 다 크기가 hidden_units와 같은 1차원 텐서인데, 프로그램에서 hidden_units = 64 라고 지정해주었고, 훈련할때 batch_size = 128 라고 지정했으므로, 은닉상태 **state_h**의 shape는 **(128, 64)**의 형태를 갖고, 셀상태 **state_c**의 shape 또한 **(128, 64)**의 형태를 갖는다.

이 프로그램에선 state_h와 state_c를 한 리스트에 묶어 encoder_states에 저장했으므로, **context vector**에 해당하는 변수는 **encoder_states**이고, **[(128, 64), (128, 64)]**의 형태를 갖는다.

이 **encode_states**가 디코더의 LSTM의 **initial_state**로 전달됨으로써 context vector가 디코더에 전달된다.