# Model Description

Advanced AI（CSE750200)

Term project

2019310290 Sangman Jung

# LSTMmodel.py

```python
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, LSTM, InputLayer, Bidirectional, TimeDistributed, Embedding, Dropout
from tensorflow.keras.optimizers import Adam

model = Sequential()
model.add(Embedding(vocab_size, 128, input_length = max_len, mask_zero = True))
model.add(Bidirectional(LSTM(256, return_sequences = True)))
model.add(Dropout(0.5))
model.add(TimeDistributed(Dense(tag_size, activation = ('relu'))))
model.add(TimeDistributed(Dense(tag_size, activation = ('softmax'))))
model.compile(loss = 'categorical_crossentropy',optimizer = Adam(0.001),metrics = ['accuracy'])
```

Library : Keras

Hypothesis model : Bidirectional LSTM

Activation : dence = Relu , output = Softmax

Loss function : Categorical cross entropy

Optimizer : Adam, lr = 0.001

# Model accuracy

```
C:\Users\92nor>python train.py --train_file train.txt
 system memory.
2674570/2674570 [==============================] - 2573s 962us/sample - loss: 0.2090 - accuracy: 0.7784
Epoch 2/2
2674570/2674570 [==============================] - 2645s 989us/sample - loss: 0.0189 - accuracy: 0.9694
```

**My laptop (only CPU) / training accuracy : 0.9694, batch size = 10000, epochs = 2**

```
1 model.fit(x_train, y_train, batch_size = 10000, epochs = 2,validation_data=(x_test,y_test))

Train on 2139656 samples, validate on 534914 samples
Epoch 1/2
2139656/2139656 [==============================] - 57s 27us/sample - loss: 0.2573 - acc: 0.7373 - val_loss: 0.0310 - val_acc: 0.9573
Epoch 2/2
2139656/2139656 [==============================] - 55s 26us/sample - loss: 0.0251 - acc: 0.9627 - val_loss: 0.0199 - val_acc: 0.9680
<tensorflow.python.keras.callbacks.History at 0x7fb8a1c56128>
```

**Google CoLab (GPU) / accuracy of (train, valid) : (0.9627,0.9680), used train.txt only (train_test_split in sklearn), batch size and epochs is the same as above.**

# train.py (1)

```
1  import pandas as pd
2  import numpy as np
3  from tensorflow.keras.preprocessing.text import Tokenizer
4  from tensorflow.keras.preprocessing.sequence import pad_sequences
5  from tensorflow.keras.utils import to_categorical
6  import argparse
7
8  # argument
9  parser = argparse.ArgumentParser()
10 parser.add_argument('--train_file',type=str)
11 parser.add_argument('--input_file',type=str)
12 parser.add_argument('--output_file',type=str)
13 args = parser.parse_args()
14
15 # tokenizer
16 def tokenize(samples):
17     tokenizer = Tokenizer()
18     tokenizer.fit_on_texts(samples)
19     return tokenizer
20
```

**Library import**

**Argument parser in terminal**

**Define the tokenizer**

# train.py (2)

```python
21  # data load
22  data = pd.read_csv(args.train_file,names=['Text','Label'],sep='\s+',quoting=3)
23  new_data = pd.concat([data.Label.str.split('+',expand = True)], axis = 1)
24
25  # change to list
26  a = pd.Series.tolist(new_data[0:-1])
27  tmp = []
28
29  # data split
30  for i in range(len(a)):
31      d = []
32      for j in range(len(a[0])):
33          if a[i][j] != None:
34              b = a[i][j].split('/')
35              b = tuple(b)
36              d.append(b)
37      tmp.append(d)
38
39  # consider '/'
40  for i in range(len(tmp)):
41      for j in range(len(tmp[i])):
42          if tmp[i][j] == ('', '', 'SP'):
43              tmp[i][j] = ('/','SP')
```

**Data load and change to list type using Pandas**

**Split the loaded data**

**This code is the exception if '/' exists in the word**

# train.py (3)

```python
45    # obtain the text and pos tags (separation)
46    texts, pos_tags = [],[]
47    for t in tmp:
48        text, tag = zip(*t)
49        texts.append(list(text))
50        pos_tags.append(list(tag))
51
52    # tokenizing
53    tx_token = tokenize(texts)
54    pos_token = tokenize(pos_tags)
55
56    with open('train_tag.txt', 'w') as file:
57        for key, item in zip(pos_token.index_word.keys(),pos_token.index_word.values()):
58            print(key,' ',item, file = file)
59
60    # size of vocab & tag set
61    vocab_size = len(tx_token.word_index) + 1
62    tag_size = len(pos_token.word_index) + 1
63
64    # change to the sequence
65    x_train = tx_token.texts_to_sequences(texts)
66    y_train = pos_token.texts_to_sequences(pos_tags)
67
68    # zero padding
69    max_len = 10 # according to the histogram
70    x_train = pad_sequences(x_train, padding = 'post', maxlen = max_len)
71    y_train = pad_sequences(y_train, padding = 'post', maxlen = max_len)
```

**Separate words and tags**

**Save the tokenized sets in order to make the index of words and tags**

**It would be used in test.py as  tag index**

**Convert to sequence data and add zero padding**

# train.py (4)

```python
73   # one-hot encoding
74   y_train = to_categorical(y_train, num_classes = tag_size)
75
76   # LSTMmodel
77   from tensorflow.keras.models import Sequential
78   from tensorflow.keras.layers import Dense, LSTM, InputLayer, Bidirectional, TimeDistributed, Embedding, Dropout
79   from tensorflow.keras.optimizers import Adam
80   from tensorflow.keras.models import load_model
81
82   model = Sequential()
83   model.add(Embedding(vocab_size, 128, input_length = max_len, mask_zero = True))
84   model.add(Bidirectional(LSTM(256, return_sequences = True)))
85   model.add(Dropout(0.5))
86   model.add(TimeDistributed(Dense(tag_size, activation = ('relu'))))
87   model.add(TimeDistributed(Dense(tag_size, activation = ('softmax'))))
88   model.compile(loss = 'categorical_crossentropy',optimizer = Adam(0.001),metrics = ['accuracy'])
89
90   # model training
91   model.fit(x_train, y_train, batch_size = 10, epochs = 2)
92
93   model.save('model.h5')
```

**One hot encoding & model import part**

**Save the train model**

# test.py (1)

```
21    # Data load
22    data = pd.read_csv(args.input_file,names=['Text','Label'], sep = '\s+',quoting = 3)
23    data.head()
24
25    # toList
26    data_value = data.Text.values
27    data_to_list = list(data_value)
28
29    # Preprocessing
30    komoran = Komoran()
31    texts_test = []
32    for i in range(len(data_to_list)):
33        postagging = komoran.morphs(data_to_list[i])
34        texts_test.append(postagging)
35
36    tx_token = tokenize(texts_test)
37
38    x_test = tx_token.texts_to_sequences(texts_test)
39
40    max_len = 10 # according to the histogram
41    x_test = pad_sequences(x_test, padding = 'post', maxlen = max_len)
```

**Package importing and data loading is similar to train.py, but the different thing is this used konlpy for tokenizing**

# test.py (2)

```python
43      # load the model
44      from tensorflow.keras.models import load_model
45      model = load_model('model.h5')
46
47      # index sets
48      index_to_word = tx_token.index_word
49      index_to_tag = {}
50      with open('train_tag.txt') as file:
51          for line in file:
52              (key, val) = line.split()
53              index_to_tag[int(key)] = val
54
55      # prediction
56      y_hat = model.predict(x_test)
57      y_hat = np.argmax(y_hat,-1)
58
59      # save the result as the text file : result.txt
60      out = open(args.output_file,'w')
61      outputlist = []
62      for i in range(len(x_test)):
63          outputlist.clear()
64          for w, p in zip(x_test[i],y_hat[i]):
65              if w != 0: # except 'pad'
66                  outputlist.append(index_to_word[w])
67                  outputlist.append('/' + index_to_tag[p].upper()+'+')
68          tmpstring = ''.join(outputlist)
69          print(tmpstring[:-1],file = out)
70      out.close()
```

**Load the model after excute train.py**

**Load POS tag index set in train.py**

**Model prediction part**

**Save as .txt file named by result.txt**