

ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH

TRƯỜNG ĐẠI HỌC BÁCH KHOA

KHOA KHOA HỌC - KỸ THUẬT MÁY TÍNH



KỸ THUẬT DỮ LIỆU

Bài tập lớn

Sử dụng CNN, LSTM và kết hợp để thực hiện sentiment analysis

Sinh viên thực hiện: Phạm Trần Đức Trọng - 2270519

Ho Chi Minh, 4-2023



Mục lục

1	Chuẩn bị dữ liệu	2
2	Tiền xử lý dữ liệu	2
3	LSTM	3
4	CNN	4
5	Kết hợp CNN và LSTM	5

1 Chuẩn bị dữ liệu

Trong bài tập lớn này tôi Sử dụng tập dữ liệu VLSP để thực hiện việc train và test cho các mô hình CNN, LSTM và kết hợp của 2 mô hình

2 Tiền xử lý dữ liệu

Thực hiện clean data đối với các loại data có những ký tự đặc biệt, emoji, xóa các từ bị duplicate, khoảng trắng...

```
1 def text_cleaner(text):
2     # lower case text
3     text = text.lower()
4     # remove emojis
5     text = re.sub(emoji_pattern, " ", text)
6     # reduce repeated characters
7     text = re.sub(r'([a-z]+?)\1+', r'\1', text)
8     # remove multiple spaces
9     text = re.sub(r'\s+', ' ', text)
10    # make sure punctuation is in the right letter
11    text = text_normalize(text)
12    # remove punctuations
13    text = text.translate(str.maketrans('', '', string.punctuation))
14    return text
```

Listing 1: Text cleaner

Lấy ngẫu nhiên 5 câu trong tập dữ liệu để xem kết quả:

1. mua phê vải hộp viên viên pin đồng hồ.
2. gồm thẻ xứng tầm Surface Pro Chất Full HD Thiết kế màn hình viền to thời Pin đối thiết dạng Tablet Pin

3. chỉ thua máy Ultrabook Dell XPS Vaio T gọi Tablet lại Tablet h pin tối
 4. Note hổng nè chán Note
 5. Ram màn hình độ phân giải sản phẩm Apple thiếp Ram chạy ok nâng sợ tốn tiền thiết lăm
- Cá nhân S ram G Iphone lăm Ram Iphone nâng chạy ISO Beta trơn tru Apple ko chạy
đua cứng tiếng nâng Ram màn hình dung

3 LSTM

```
1 def lstm_model():
2     x = embedding_layer(input_layer)
3     x = Bidirectional(LSTM(128, return_sequences=True, dropout=0.2,
4         recurrent_dropout=0.2))(x)
5     x = Dropout(0.2)(x)
6     x = Bidirectional(LSTM(128, return_sequences=True, dropout=0.2,
7         recurrent_dropout=0.2))(x)
8     x = Dropout(0.2)(x)
9     x = GlobalMaxPool1D()(x)
10    x = Dense(128, activation="relu")(x)
11    x = Dropout(0.2)(x)
12    x = Dense(3, activation="softmax")(x)
13    model = Model(inputs=input_layer, outputs=x)
14    model.compile(loss='binary_crossentropy', optimizer=adam, metrics=['accuracy',
15        ])
16    model.summary()
17    return model
```

Listing 2: lstm model

4 CNN

```
1 def cnn_model():
2     from keras.models import Model
3     from keras.layers import *
4     from keras.callbacks import EarlyStopping
5     from tensorflow.keras.optimizers import Adam
6     from keras.models import Model
7     from keras import regularizers
8     sequence_length = data.shape[1]
9     filter_sizes = [3,4,5]
10    num_filters = 100
11    drop = 0.5
12
13    inputs = Input(shape=(sequence_length,))
14    embedding = embedding_layer(inputs)
15
16    ##### LSTM ONLY #####
17    reshape = Reshape((sequence_length, EMBEDDING_DIM))(embedding)
18
19    ##### SINGLE LSTM #####
20    lstm_0 = LSTM(512)(reshape)
21
22    # YOU WANNA ADD MORE LSTM LAYERS? UNCOMMENT THIS #
23    lstm_2 = LSTM(1024, return_sequences=True)(reshape)
24    lstm_1 = LSTM(512, return_sequences=True)(lstm_2)
25    lstm_0 = LSTM(256)(lstm_1)
26
27    #####
28
29
```

```
30 ##### CRNN #####
31 reshape = Reshape((sequence_length, EMBEDDING_DIM))(embedding)
32
33 conv_0 = Conv1D(num_filters, (filter_sizes[0], ), padding="same", activation='
relu', kernel_regularizer=regularizers.l2(0.01))(reshape)
34 conv_1 = Conv1D(num_filters, (filter_sizes[1], ), padding="same", activation='
relu', kernel_regularizer=regularizers.l2(0.01))(reshape)
35 conv_2 = Conv1D(num_filters, (filter_sizes[2], ), padding="same", activation='
relu', kernel_regularizer=regularizers.l2(0.01))(reshape)
36
37 conv_0 = MaxPool1D(300)(conv_0)
38 conv_1 = MaxPool1D(300)(conv_1)
39 conv_2 = MaxPool1D(300)(conv_2)
40 # Reshape output to match RNN dimension
41 conv_0 = Reshape((-1, num_filters))(conv_0)
42 conv_1 = Reshape((-1, num_filters))(conv_1)
43 conv_2 = Reshape((-1, num_filters))(conv_2)
44
45 concat = concatenate([conv_0, conv_1, conv_2])
46 concat = Flatten()(concat)
47
48 dropout = Dropout(drop)(concat)
49 output = Dense(units=3, activation='softmax', kernel_regularizer=regularizers.
l2(0.01))(dropout)
```

Listing 3: cnn model

5 Kết hợp CNN và LSTM

```
1 def lstm_cnn_model():
2     from keras.models import Model
```

```
3  from keras.layers import *
4  from keras.callbacks import EarlyStopping
5  from tensorflow.keras.optimizers import Adam
6  from keras.models import Model
7  from keras import regularizers
8  sequence_length = data.shape[1]
9  filter_sizes = [3,4,5]
10 num_filters = 100
11 drop = 0.5
12
13 inputs = Input(shape=(sequence_length,))
14 embedding = embedding_layer(inputs)
15
16 ##### LSTM ONLY #####
17 reshape = Reshape((sequence_length, EMBEDDING_DIM))(embedding)
18
19 ##### SINGLE LSTM #####
20 lstm_0 = LSTM(512)(reshape)
21
22 # YOU WANNA ADD MORE LSTM LAYERS? UNCOMMENT THIS #
23 lstm_2 = LSTM(1024, return_sequences=True)(reshape)
24 lstm_1 = LSTM(512, return_sequences=True)(lstm_2)
25 lstm_0 = LSTM(256)(lstm_1)
26
27 #####
28
29
30 ##### CRNN #####
31 reshape = Reshape((sequence_length, EMBEDDING_DIM))(embedding)
32
33 conv_0 = Conv1D(num_filters, (filter_sizes[0], ),padding="same",activation=''
```

```
relu',kernel_regularizer=regularizers.l2(0.01))(reshape)
34 conv_1 = Conv1D(num_filters, (filter_sizes[1], ),padding="same",activation='
relu',kernel_regularizer=regularizers.l2(0.01))(reshape)
35 conv_2 = Conv1D(num_filters, (filter_sizes[2], ),padding="same",activation='
relu',kernel_regularizer=regularizers.l2(0.01))(reshape)
36
37 conv_0 = MaxPool1D(300)(conv_0)
38 conv_1 = MaxPool1D(300)(conv_1)
39 conv_2 = MaxPool1D(300)(conv_2)
40 # Reshape output to match RNN dimension
41 conv_0 = Reshape((-1, num_filters))(conv_0)
42 conv_1 = Reshape((-1, num_filters))(conv_1)
43 conv_2 = Reshape((-1, num_filters))(conv_2)
44
45 concat = concatenate([conv_0, conv_1, conv_2])
46 concat = Flatten()(concat)
47
48 dropout = Dropout(drop)(concat)
49 output = Dense(units=3, activation='softmax',kernel_regularizer=regularizers.
l2(0.01))(dropout)
```

Listing 4: lstm cnn model