

HAM / SPAM CLASSIFICATION USING Bidirectional LSTM

- Junheui Lee (2021020357)

< requirements for project >

For analysis) tensorflow, keras, sklearn

For preprocessing) pandas, numpy, random, nltk

For visualization) plotly, matplotlib, seaborn, wordcloud

1. Introduction

특정한 메시지가 메일함에 도착했을 때, 메시지가 스팸으로 분류될 경우 스팸 메일함으로 보내지게 됩니다. 저는 어떠한 메시지가 스팸인지 아닌지 한 학기 동안 배운 딥러닝으로 잘 구분할 수 있을지에 대한 궁금증이 생겨, UCI Machine Learning Repository에 있는 SMS Spam Collection Dataset을 download 받아 분석을 진행했습니다.

2. EDA, Visualization, Preprocessing

```
1 spam[spam.duplicated()] ## 중복 데이터
```

	label	message
102	ham	As per your request 'Melle Melle (Oru Minnamin...
153	ham	As per your request 'Melle Melle (Oru Minnamin...
206	ham	As I entered my cabin my PA said," Happy B'd...
222	ham	Sorry, I'll call later
325	ham	No calls..messages..missed calls
...
5524	spam	You are awarded a SiPix Digital Camera! call 0...
5535	ham	I know you are thinkin malaria. But relax, chi...
5539	ham	Just sleeping, and surfing
5553	ham	Hahaha..use your brain dear
5558	ham	Sorry, I'll call later

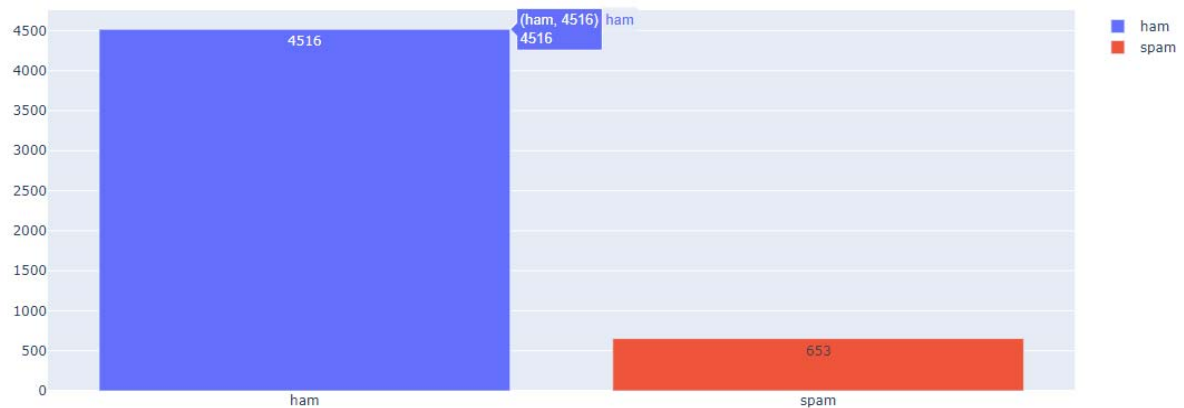
403 rows x 2 columns

spam data는 5572개의 메시지가 있었고, 각각의 메시지는 spam인지 아닌지에 대한 label이 달려 있습니다. 그 중 403개의 메시지가 동일한 내용의 메시지임을 확인해, 중복되는 메시지를 제거했습니다.

```
1 spam.groupby('label').describe().T
```

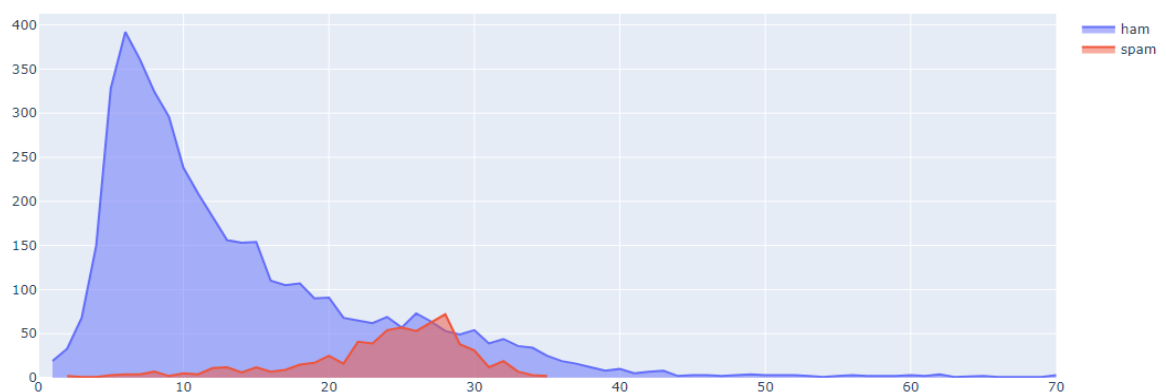
	label	ham	spam
message	count	4516	653
	unique	4516	653
	top	I'm leaving my house now...	Ur ringtone service has changed! 25 Free credi...
	freq	1	1

Count of spam, ham labels



중복 메시지를 제거한 결과 전체 데이터는 5572개에서 5169개가 되었고, HAM인 데이터는 4516개, SPAM인 데이터는 653개임을 확인했습니다.

Distribution of length of ham, spam email



SPAM인 메시지는 HAM인 메시지에 비해 비교적 긴 길이를 가지는 것을 알 수 있었습니다. 이를 통해 SPAM인 메시지와 아닌 메시지는 명확한 차이를 가지고 있는 것을 알 수 있었습니다.

```
1 spam['message']
```

```
0      Go until jurong point, crazy.. Available only ...
1              Ok lar... Joking wif u oni...
2      Free entry in 2 a wkly comp to win FA Cup fina...
3      U dun say so early hor... U c already then say...
4      Nah I don't think he goes to usf, he lives aro...
      ...
5567     This is the 2nd time we have tried 2 contact u...
5568             Will i_ b going to esplanade fr home?
5569     Pity, * was in mood for that. So...any other s...
5570     The guy did some bitching but I acted like i'd...
5571             Rofl. Its true to its name
Name: message, Length: 5169, dtype: object
```

SPAM data의 메시지에는 분석에 불필요한 요소들이 많이 포함되어 있는 것을 확인할 수 있었습니다. 그래서 핵심적인 단어들만 추출하여 메시지 분석이 용이하도록 전처리를 진행했습니다.

```
def clean_corpus(corpus):
    # make text to lower text
    corpus = str(corpus).lower()
    # remove text in square brackets
    corpus = re.sub('\[.*?\]', '', corpus)
    # remove homepage type
    corpus = re.sub('https?://\S+|www\.\S+', '', corpus)
    # remove continus '.'
    corpus = re.sub('<.*?>+', '', corpus)
    # remove special characters
    corpus = re.sub('%s' % re.escape(string.punctuation), '', corpus)
    # remove enter
    corpus = re.sub('\n', '', corpus)
    # remove words containing digits
    corpus = re.sub('\w*\d\w*', '', corpus)

    return corpus
```

```
stop_words = stopwords.words('english')
my_stop_words = ['u', 'im']
stop_words = stop_words + my_stop_words

def remove_stopwords(corpus):
    corpus = ' '.join(word for word in corpus.split(' ') if word not in stop_words)
    return corpus
```

```
snow_stemmer = nltk.SnowballStemmer("english")

def stemming_corpus(corpus):
    corpus = ' '.join(snow_stemmer.stem(word) for word in corpus.split(' '))
    return corpus
```

```

1 def preprocessing_corpus(corpus):
2     # Clean punctuation, urls, and so on
3     corpus = clean_corpus(corpus)
4     # Remove stopwords
5     corpus = remove_stopwords(corpus)
6     # Stem all the words in the sentence
7     corpus = stemming_corpus(corpus)
8
9     return corpus

1 spam['message_clean'] = spam['message'].apply(preprocessing_corpus)

1 spam.head()

```

	label	message	message_len	message_clean
0	ham	Go until jurong point, crazy.. Available only ...	20	go jurong point crazi avail bugi n great world...
1	ham	Ok lar...Joking wif u oni...	6	ok lar joke wif oni
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	28	free entri wkly comp win fa cup final tkts m...
3	ham	U dun say so early hor... U c already then say...	11	dun say earli hor c already say
4	ham	Nah I don't think he goes to usf, he lives aro...	13	nah dont think goe usf live around though

전처리의 과정은 다음과 같습니다.

우선, 다음과 같은 과정을 통해 분류에 유의하지 않은 문자들을 제거했습니다.

- I) 대문자와 소문자가 섞여 있어 소문자로 통일해 주었습니다.
- II) 대괄호 안에 있는 문자들을 제거해 주었습니다.
- III) https:와 같은 홈페이지를 제거해 주었습니다.
- IV) 연속된 '.' 및 특수문자, enter를 제거해 주었습니다.
- V) 숫자를 포함한 문자를 제거했습니다.

다음에, 문장에 많이 사용되지만 의미가 없는 stopwords를 제거했습니다. stopwords를 제거함으로써, 메시지의 핵심 문장들을 모델이 더 잘 파악할 수 있도록 했습니다.

마지막으로, stemming을 진행했습니다. 사람들이 사용하는 언어에서 단어는 문법적 요소나 맥락에 따라 다양한 형태로 변화하는 것을 알 수 있습니다. 미래 / 현재 / 과거 시제가 존재하는 것을 예시로 들 수 있습니다. stemming을 진행함으로써 변화된 단어의 원형을 찾아줄 수 있고, 단어의 원형만 사용해서 전처리 전에 비해 메시지의 복잡함을 줄여주는 메시지 정규화의 효과를 얻을 수 있었습니다.

전처리를 진행한 결과, 전처리를 마친 메시지가 기존의 메시지에 비해 분석에 용이한 형태가 된 것을 message_clean에서 확인할 수 있었습니다.

```
1 from sklearn.preprocessing import LabelEncoder
2
3 LE = LabelEncoder()
4 LE.fit(spam['label'])
5
6 spam['label_encoded'] = LE.transform(spam['label'])
```

5

3. Modeling

GLOVE Word Embedding

```
1 texts = spam['message_clean']
2 target = spam['label_encoded']
```

```
1 word_tokenizer = Tokenizer() # keras tokenizer
2 word_tokenizer.fit_on_texts(texts)
3
4 vocab_length = len(word_tokenizer.word_index) + 1
5 vocab_length
```

```
def embed_corpus(corpus):
    return word_tokenizer.texts_to_sequences(corpus)

longest_train = max(texts, key=lambda sentence: len(word_tokenize(sentence)))
maxlen = len(word_tokenize(longest_train))

train_padded_sentences = pad_sequences(
    embed_corpus(texts),
    maxlen,
    padding='post'
)
```

```
1 embeddings_dictionary = dict()
2 embedding_dim = 100
3
4 with open('glove.6B.100d.txt', encoding='UTF8') as f:
5     for line in f.readlines():
6         records = line.split()
7         word = records[0]
8         vector_dimensions = np.asarray(records[1:], dtype='float32')
9         embeddings_dictionary[word] = vector_dimensions
10
11 f.close()
12
13 print(len(embeddings_dictionary))
```

```
1 embedding_matrix = np.zeros((vocab_length, embedding_dim))
2
3 for word, index in word_tokenizer.word_index.items():
4     embedding_vector = embeddings_dictionary.get(word)
5     if embedding_vector is not None:
6         embedding_matrix[index] = embedding_vector
```

전체 corpus에서 단어들의 동시 발생 빈도를 기반으로 단어의 의미를 살려 vector화 시켜주는 GLOVE embedding을 사용하기 위해 메시지를 tokenizing 시킨 후, 400000 개의 단어로 학습된 GLOVE pre-trained weights를 가져와 적용했습니다.

```

1 model = Sequential()
2
3 model.add(Embedding(
4     input_dim=embedding_matrix.shape[0],
5     output_dim=embedding_matrix.shape[1],
6     weights = [embedding_matrix],
7     input_length=maxlen
8 ))
9
10 model.add(Bidirectional(LSTM(maxlen, return_sequences = True, recurrent_dropout=0.2)))
11
12 model.add(GlobalMaxPool1D())
13 model.add(BatchNormalization())
14 model.add(Dropout(0.5))
15 model.add(Dense(maxlen, activation = "relu"))
16 model.add(Dropout(0.5))
17 model.add(Dense(maxlen, activation = "relu"))
18 model.add(Dropout(0.5))
19 model.add(Dense(1, activation = 'sigmoid')) # binary classification

```

```

# for freezing glove weights
model.layers[0].set_weights([embedding_matrix])
model.layers[0].trainable=False

```

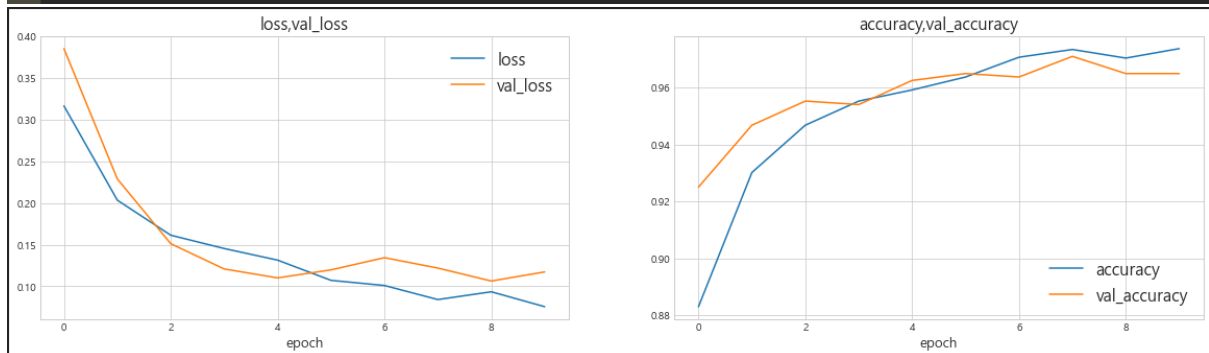
SPAM / HAM 분류를 위해 Bidirectional LSTM model을 구축했습니다.

과적합을 막기 위해 GlobalMaxPooling, BatchNormalization, dropout, recurrent_dropout을 적용했으며, GLOVE weights의 활용을 위해, 해당 weights 부분을 non-trainable하게 설정했습니다. 마지막으로 SPAM / HAM 분류를 위해 sigmoid 활성화함수를 사용했습니다.

```

1 model.compile(optimizer='rmsprop', loss='binary_crossentropy', metrics=['accuracy'])
2
3 checkpoint = ModelCheckpoint('model_bidirectional_lstm.h5', monitor = 'val_loss', verbose = 1, save_best_only = True)
4 reduce_lr = ReduceLRonPlateau( monitor = 'val_loss', factor = 0.2, verbose = 1, patience = 2, min_lr = 0.001)
5
6 history = model.fit(
7     partial_x_train, partial_y_train,
8     epochs = 10, batch_size = 32,
9     validation_data = (x_val, y_val),
10    verbose = 1,
11    callbacks = [reduce_lr, checkpoint]
12 )

```



안정적인 결과 복원을 위해 seed를 42로 고정하고 모델링을 진행했습니다. ReducedLROnPlateau scheduler를 통해 validation loss의 개선 정도에 따라 learning rate를 조정하여 최적의 모델 학습이 진행될 수 있도록 했습니다. 10 epoch, 32 batch size를 사용해 모델을 적합했고, learning curve를 통해 epoch 경과에 따라 안정적인 metric의 향상이 있음을 확인할 수 있었습니다.

```
1 model.compile(optimizer='rmsprop', loss='binary_crossentropy', metrics=['accuracy'])
2
3 checkpoint = ModelCheckpoint('model_bidirectional_lstm.h5', monitor = 'val_loss', verbose = 1, save_best_only = True)
4 reduce_lr = ReduceLROnPlateau( monitor = 'val_loss', factor = 0.2, verbose = 1, patience = 2, min_lr = 0.001)
5
6 history = model.fit(
7     x_train, y_train,
8     epochs = 10, batch_size = 32,
9     validation_data = (x_test, y_test),
10    verbose = 1,
11    callbacks = [reduce_lr, checkpoint]
12 )
```

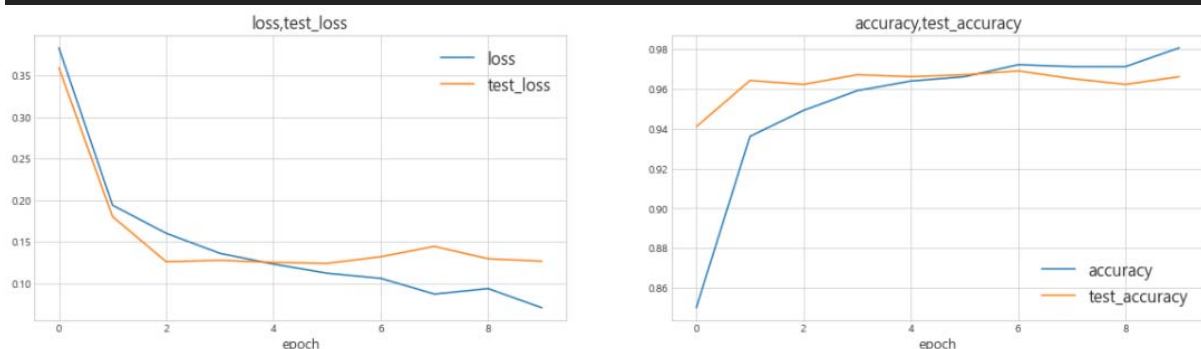
Epoch 9/10
4135/4135 [=====] - 18s 4ms/step - loss: 0.0935 - accuracy: 0.9712 - val_loss: 0.1294 - val_accuracy: 0.9623

Epoch 00009: val_loss did not improve from 0.12379

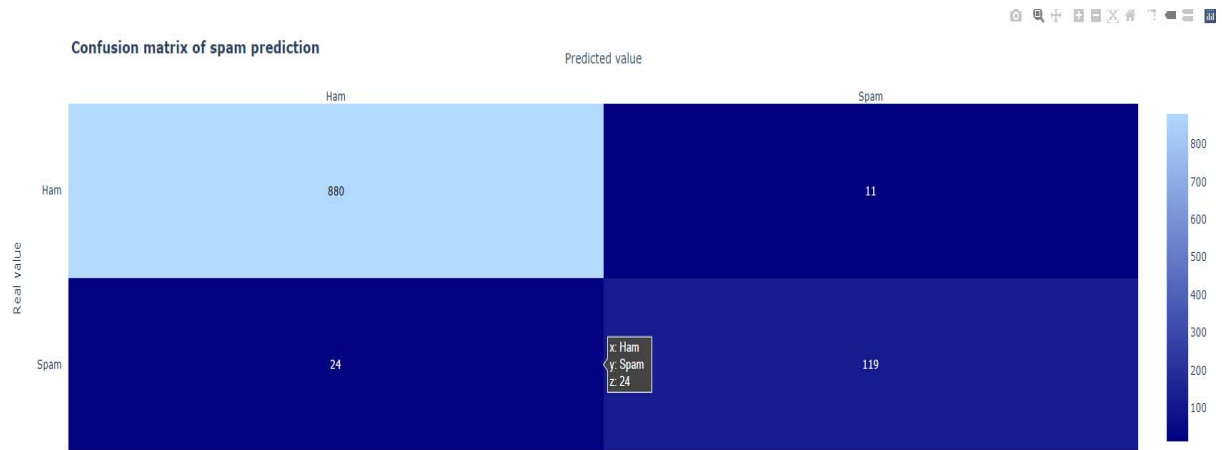
Epoch 10/10
4135/4135 [=====] - 18s 4ms/step - loss: 0.0704 - accuracy: 0.9807 - val_loss: 0.1264 - val_accuracy: 0.9662

Epoch 00010: ReduceLROnPlateau reducing learning rate to 0.001.

Epoch 00010: val_loss did not improve from 0.12379



Validation set을 통해 hyperparameter tuning이 완료된 모델을 앞과 동일한 10 epoch, 32 batch size를 사용하여 전체 train set에 적합하였습니다. 그 결과 training accuracy는 0.9807, test accuracy는 0.9662로 미세한 차이를 보였고, learning curve를 통해 최종 모델에서 overfitting이 완화된 것을 확인할 수 있었습니다.



```
1 from sklearn.metrics import zero_one_loss
2 misclassification_rate = zero_one_loss(y_test, y_pred)
3 misclassification_rate
```

0.03384912959381048

test set의 design matrix로 SPAM/HAM classification prediction을 한 결과, 모델이 대부분의 classification을 올바르게 한 것을 confusion matrix를 통해 확인할 수 있었습니다. 모델의 misclassification rate는 3.3%로 준수한 성능을 보였습니다.