

## HW5: Fake News Detection 2

1. Data Preprocessing
2. Model Training
3. Results

全部展開

[回到顶部](#)

移至底部

## HW5: Fake News Detection 2

309554001 劉雨恩

## 1. Data Preprocessing

1. 將 `train.csv`、`test.csv` 和 `sample_submission.csv` 的資料 (`id`, `text`, `label`) 提取出來，並消除 HTML tag 及 stop words (利用 `spacy` 提供的停頓詞列表)。

```

1 import pandas as pd
2 import spacy
3
4 spacy.load('en_core_web_sm')
5 spacy_stopwords = spacy.lang.en.stop_words.STOP_WORDS
6
7 def data_preprocess(data_path, mode, stop_words):
8
9     def rm_tags_stops(text, stop_words):
10         # remove tags
11         re_tag = re.compile(r'<[^>]+>')
12         out_text = re_tag.sub(' ', text.lower())
13
14         # remove stop words
15         out_text = (" ").join([word for word in out_text.split(' ') if word not in stop_words])
16
17         return out_text
18
19     # read data
20     fp = open(f'{data_path}{mode}.csv')
21     data_lines = fp.readlines()
22     fp.close()
23
24     # data preprocess
25     data_df = pd.DataFrame([], columns=['id', 'text', 'label'])
26     avg_text_len = 0
27     if mode == 'train':
28         for i, l_i in enumerate(data_lines):
29             if i == 0:
30                 continue
31             try:
32                 dict_i = {}
33                 text_i, label = l_i.split('\t')
34                 text_i = rm_tags_stops(text_i.strip(), stop_words)
35                 avg_text_len += len(text_i.split(' '))
36
37                 dict_i['id'] = [str(i)]
38                 dict_i['text'] = [text_i]
39                 dict_i['label'] = [int(label.strip())]
40                 data_df = pd.concat([data_df, pd.DataFrame.from_dict(dict_i, orient='columns')])
41             except:
42                 pass
43     else:
44         label_data = pd.read_csv(f'{data_path}sample_submission.csv')
45
46         for i, (test_li, label_i) in enumerate(zip(data_lines, label_data['label'])):
47             if i == 0:
48                 continue
49             dict_i = {}
50             id, text_i = test_li.split('\t')
51             text_i = rm_tags_stops(text_i.strip(), stop_words)
52             avg_text_len += len(text_i.split(' '))
53
54             dict_i['id'] = [id.strip()]
55             dict_i['text'] = [text_i]
56             dict_i['label'] = [int(label_i)]
57             data_df = pd.concat([data_df, pd.DataFrame.from_dict(dict_i, orient='columns')])
58     print(f'{mode} avg text len: {avg_text_len / len(data_df)}')
59
60     return data_df['text'], np.array(data_df['label'].values)

```

2. 使用 `Tokenizer` 模組建立 `token`，建立一個 `3800` (`hyper_params['dict_len']`) 字的字典：讀取所有訓練文檔資料之後，會依照每個英文字在資料出現的次數進行排序，並將前 `3800` 名的英文單字加進字典中。
3. 透過 `texts_to_sequences` 可以將訓練和測試集資料中的文檔轉換為數字列表。
4. 每一篇影評文字字數不固定，但後續進行深度學習模型訓練時長度必須固定，因此需要截長補短 `sequences.pad_sequences`：長度小於 `380` (`hyper_params['content_len']`) 的，前面的數字補 `0`；長度大於 `380` 的，截去前面的數字。

```
1 from keras.preprocessing import sequence
2 from keras.preprocessing.text import Tokenizer
3
4 # create token
5 token = Tokenizer(num_words=hyper_params['dict_len'])
6 token.fit_on_texts(text_train.values)
7 x_train = token.texts_to_sequences(text_train.values)
8 x_train = sequence.pad_sequences(x_train, maxlen=hyper_params['content_len'])
9 x_test = token.texts_to_sequences(text_test.values)
10 x_test = sequence.pad_sequences(x_test, maxlen=hyper_params['content_len'])
```

## 2. Model Training

- 訓練兩種模型 SimpleRNN 和 LSTM，並將每個 Epoch 的 training loss、training accuracy 和 testing accuracy 記錄下來。
- Hyperparameters
  - epoch: 100
  - batch size: 4987
  - drop out: 0.7
  - optimizer: Adam (learning rate: 1e-3)

```

1 import keras
2 from keras.models import Sequential
3 from keras.layers.core import Dense, Dropout
4 from keras.layers.embeddings import Embedding
5 from keras.layers.recurrent import SimpleRNN, LSTM
6
7 hyper_params = {'epoch': 100, 'batch_size': 4987, 'drop_out': 0.7,
8                 'dict_len': 3800, 'content_len': 380}
9
10 def train(target_path, model_type, params, x_train, y_train, x_test, y_test):
11     # set up model
12     model = Sequential()
13     model.add(Embedding(output_dim=128,
14                         input_dim=params['dict_len'],
15                         input_length=params['content_len']))
16     model.add(Dropout(params['drop_out']))
17     if model_type == 'RNN':
18         model.add(SimpleRNN(units=128))
19     else:
20         model.add(LSTM(units=128))
21     model.add(Dense(units=128, activation='relu'))
22     model.add(Dropout(params['drop_out']))
23     model.add(Dense(units=1, activation='sigmoid'))
24     model.summary()
25
26     # define model
27     model.compile(loss='binary_crossentropy',
28                 optimizer=keras.optimizers.Adam(learning_rate=1e-3),
29                 metrics=['accuracy'])
30
31     # training
32     fp_train_loss = open(f'{target_path}{model_type}_train_loss.txt', 'w')
33     fp_train_acc = open(f'{target_path}{model_type}_train_acc.txt', 'w')
34     fp_test_acc = open(f'{target_path}{model_type}_test_acc.txt', 'w')
35
36     for epoch_i in range(1, params['epoch'] + 1):
37         train_history = model.fit(x_train, y_train,
38                                 epochs=1, batch_size=params['batch_size'],
39                                 shuffle=True, workers=12,
40                                 verbose=0)
41         train_loss = train_history.history['loss'][0]
42         train_acc = train_history.history['accuracy'][0]
43         print(f'\nEpoch {epoch_i}/{params["epoch"]}')
44         print(f'Loss: {train_loss}, Accuracy: {train_acc}')
45         fp_train_loss.write(f'{train_loss}\n')
46         fp_train_acc.write(f'{train_acc}\n')
47
48     score = model.evaluate(x_test, y_test, verbose=0, workers=12)
49     print(f'Testing accuracy: {score[1]}')
50     fp_test_acc.write(f'{score[1]}\n')
51
52     fp_train_loss.close()
53     fp_train_acc.close()
54     fp_test_acc.close()

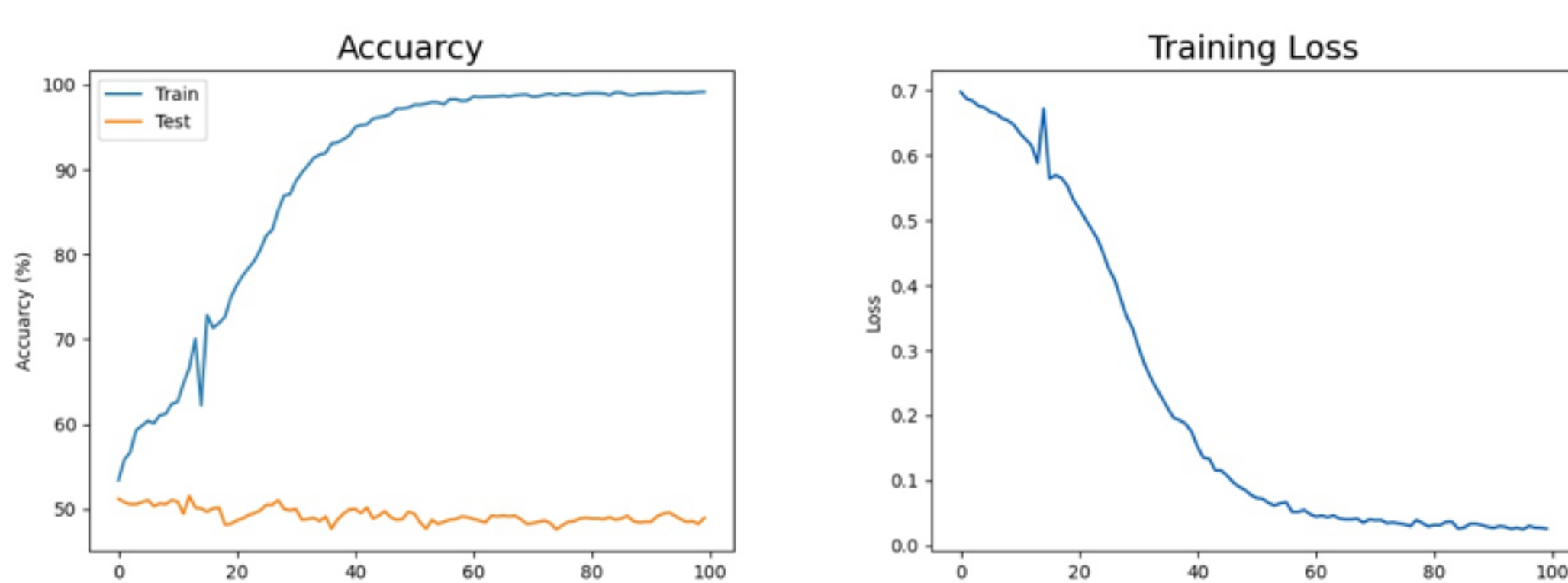
```

### 3. Results

最後雖然兩種模型 training accuracy 都有提升至將近百分百、loss 也都趨近 0，但 testing accuracy 都維持在 50% 左右，有 overfitting 的現象。

- RNN

	Train	Test
Accuracy	99.16%	51.52%



- LSTM

	Train	Test
Accuracy	97.03%	50.80%

