

# HW4: Fake News Detection

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。

```

1 import pandas as pd
2
3 # read data
4 fp = open(f'{data_path}{mode}.csv')
5 data_lines = fp.readlines()
6 fp.close()
7
8 # data preprocess
9 data_df = pd.DataFrame([], columns=['id', 'text', 'label'])
10 if mode == 'train':
11     for i, l_i in enumerate(data_lines):
12         if i == 0:
13             continue
14         try:
15             dict_i = {}
16             text_i, label = l_i.split('\t')
17
18             dict_i['id'] = [str(i)]
19             dict_i['text'] = [text_i.strip()]
20             dict_i['label'] = [int(label.strip())]
21             data_df = pd.concat([data_df, \
22                                 pd.DataFrame.from_dict(dict_i, orient='columns')])
23         except:
24             pass
25 else:
26     label_data = pd.read_csv(f'{data_path}sample_submission.csv')
27
28     for i, (test_li, label_i) in \
29         enumerate(zip(data_lines, label_data['label'])):
30         if i == 0:
31             continue
32         dict_i = {}
33         id, text_i = test_li.split('\t')
34
35         dict_i['id'] = [id.strip()]
36         dict_i['text'] = [text_i.strip()]
37         dict_i['label'] = [int(label_i)]
38         data_df = pd.concat([data_df, \
39                             pd.DataFrame.from_dict(dict_i, orient='columns')])

```

2. 利用 `spacy` 提供的停頓詞列表來去除停頓詞，並使用 `TfidfVectorizer` 將文字資料型態轉換成向量。

```

1 import spacy
2 import pandas as pd
3 from sklearn.feature_extraction.text import TfidfVectorizer
4
5 spacy.load('en_core_web_sm')
6 spacy_stopwords = spacy.lang.en.stop_words.STOP_WORDS
7
8 def data_preprocess(data_path, mode, stop_words):
9     ...
10     tv = TfidfVectorizer(stop_words=stop_words, max_features=10000)
11     X_tf = tv.fit_transform(data_df['text']).toarray()

```

## 2. Model Training

訓練三種模型 XGBClassifier、GradientBoostingClassifier 和 LGBMClassifier，並計算 Accuracy、Precision、Recall 和 F-measure。

[illegible]

### 3. Results

	<b>XGBoost</b>	<b>GBDT</b>	<b>LightGBM</b>
Accuracy	0.497	0.502	0.497
Precision	0.4	0.424	0.4
Recall	0.036	0.023	0.036
F-measure	0.066	0.043	0.066