

災難推文辨識

小組成員和負責工作：

109590041 范遠皓 撰寫程式和測試

109590043 柯瑞霖 撰寫程式和報告

環境

使用的語言：Python

所需套件：

requirements.txt 檔案內

[pytorch 安裝](#)

安裝辦法：

Zip 檔內有 requirements.txt

終端機輸入以下指令安裝

```
pip3 install -r requirements.txt
```

資料讀取：

```
13 data = pd.read_csv('nlp-getting-started/train.csv')
14 test_data = pd.read_csv('nlp-getting-started/test.csv')
```

資料預處理：

將 keyword 標籤內有包含 '%20' 的替換掉
移除網址、@someone

```
16 def keyword_preprocess(text):
17     """Clean keyword by removing '%20'"""
18     if pd.notnull(text):
19         text = text.replace("%20", " ")
20     else:
21         text = ''
22     return text
23
24 def remove_url(text):
25     url_pattern = re.compile(r'https?://t\.\co/[\^s]*)')
26     new_text = url_pattern.sub('', text)
27     return new_text
28
29 def remove_at(text):
30     at_pattern = re.compile(r'@[^\s]*')
31     new_text = at_pattern.sub('', text)
32     return new_text
33
34 def text_preprocess(text):
35     """Clean text by removing url and @someone"""
36     text = remove_url(text)
37     text = remove_at(text)
38     return text
39
40 # remove url and @ from text
41 data['text'] = data['text'].apply(text_preprocess)
42 test_data['text'] = test_data['text'].apply(text_preprocess)
43
44 # remove %20 from keyword
45 data['keyword'] = data['keyword'].apply(keyword_preprocess)
46 test_data['keyword'] = test_data['keyword'].apply(keyword_preprocess)
47
```

將 keyword 加入 text(tweet)

```
48 # combine keyword and text
49 data['keyword_text'] = data.apply(lambda row: row['keyword'] + ' ' + row['text'], axis=1)
50 test_data['keyword_text'] = test_data.apply(lambda row: row['keyword'] + ' ' + row['text'], axis=1)
```

將資料轉成 dataset

```
52 train_data_dict = {
53     "text": data["keyword_text"].tolist(),
54     "label": data["target"].tolist()
55 }
56
57 test_data_dict = {
58     "text": test_data["keyword_text"].tolist()
59 }
60
61 train_dataset = Dataset.from_dict(train_data_dict)
62 test_dataset = Dataset.from_dict(test_data_dict)
```

載入 BERT 模型：

初始化一個 BERT tokenizer，並將文本轉換成 BERT 模型所需的輸入格式，載入預訓練模型 BERT，用於訓練一個基於 BERT 模型的序列分類器

```
64 checkpoint = "bert-base-uncased"
65 tokenizer = AutoTokenizer.from_pretrained(checkpoint)
66
67 def tokenize_function(example):
68     return tokenizer(example["text"], truncation=True)
69
70 tokenized_train_dataset = train_dataset.map(tokenize_function, batched=True)
71 tokenized_test_dataset = test_dataset.map(tokenize_function, batched=True)
72
73 # use dynamic padding
74 data_collator = DataCollatorWithPadding(tokenizer=tokenizer)
75
76 model = AutoModelForSequenceClassification.from_pretrained(checkpoint, num_labels=2)
77
```

調整訓練參數並輸出結果：


```
78 training_args = TrainingArguments(  
79     "test-trainer",  
80     report_to='none',  
81     num_train_epochs=2,  
82     save_strategy = "epoch"  
83 )  
84  
85 trainer = Trainer(  
86     model,  
87     training_args,  
88     train_dataset=tokenized_train_dataset,  
89     data_collator=data_collator,  
90     tokenizer=tokenizer,  
91 )  
92  
93 trainer.train()  
94  
95 predictions = trainer.predict(tokenized_test_dataset)  
96 preds = np.argmax(predictions.predictions, axis=-1)  
97  
98 submission = pd.DataFrame({'id':test_data['id'],'target':preds})  
99 submission.to_csv('nlp-getting-started/mySubmission.csv', index=False)
```

執行結果：

有預處理

	submission1.csv Complete · 19h ago	0.83634
---	--	----------------

沒有預處理

	submission_test.csv Complete · now	0.84063
---	--	----------------

找出受災地區：

使用了 spaCy 套件的 en_core_web_sm 模型，檢測每條貼文中的地點

```
import en_core_web_sm
nlp = en_core_web_sm.load()
def location_detect(text):
    doc = nlp(text)
    data = [(X.text, X.label_) for X in doc.ents]
    for word, pos in data:
        if pos == 'GPE':
            return(word)
test_data_location['location'] = test_data_location['text'].apply(location_detect)
# print(test_data_location['location'])
result = pd.DataFrame({'text':test_data_location['text'],'target':preds,'location':test_data_location['location']})
result = result[result['target'] == 1]
result = result[result['location'].notnull()]
print(result)
```

	text	target	\
4	Typhoon Soudelor kills 28 in China and Taiwan	1	
15	Birmingham Wholesale Market is ablaze BBC News...	1	
34	Accident on A27 near Lewes is it Kingston Ro...	1	
36	For Legal and Medical Referral Service Call u...	1	
52	'We are still living in the aftershock of Hiro...	1	
...	
3238	Wreckage 'Conclusively Confirmed' as From MH37...	1	
3239	Wreckage 'Conclusively Confirmed' as From MH37...	1	
3254	Officials: Alabama home quarantined over possi...	1	
3257	The death toll in a #IS-suicide car bombing on...	1	
3260	Green Line derailment in Chicago	1	

	location
4	China
15	Birmingham
34	Lewes
36	Legal
52	Hiroshima
...	...
3238	Wreckage
3239	Wreckage
3254	Alabama
3257	the Village of Rajman
3260	Chicago

[441 rows x 3 columns]