自然語言處理 HW1

組員 :

110590450 歐佳昀(60%) => 作業架構、流程、分析、改善

110590452 莊于潔(40%) => 作業模型挑選、參數調整

| 環境：使用 Jupyter Notebook | 語言：python |
|---|---|

Classification results (Best *CRF)

- Accuracy = 0.936

- Precision = 0.927

- Recall = 0.936

- F-measure = 0.928

```
================================================

*CRF :
Accuracy: 0.936
Weighted Average Precision: 0.927
Weighted Average Recall: 0.936
Weighted Average F1 Score: 0.928


================================================
```

經使用不同模型後，分類表現 spaCy < NLTK < CRF

| spaCy train - test | spaCy train - (test+dev) |
|---|---|
| - Precision = 0.613 | - Precision = 0.641 |
| - Recall = 0.497 | - Recall = 0.587 |
| - F-measure = 0.540 | - F-measure = 0.611 |
| - Accuracy = 0.434 | - Accuracy = 0.521 |
| ```<br>================================<br><br>spaCy :<br>Accuracy:0.434<br>Weighted Average Precision: 0.613<br>Weighted Average Recall: 0.497<br>Weighted Average F1 Score: 0.540<br><br>================================<br>``` | ```<br>================================<br><br>spaCy :<br>Accuracy:0.521<br>Weighted Average Precision: 0.641<br>Weighted Average Recall: 0.587<br>Weighted Average F1 Score: 0.611<br><br>================================<br>``` |
| spaCy k -fold (5-fold) | spaCy (train for command - dev for command - test): |
| - Precision = 0.3357 | - Precision = 0.879 |
| - Recall = 0.3619 | - Recall = 0.825 |
| - F-measure = 0.3465 | - F-measure = 0.850 |
| - Accuracy = 0.2267 | - Accuracy = 0.781 |

```
==============================
spaCy :
Accuracy: 0.22675619834710745
Weighted Average Precision: 0.3357
Weighted Average Recall: 0.3619
Weighted Average F1 Score: 0.3465

==============================
```

```
==============================
spaCy :
Accuracy: 0.781
Weighted Average Precision: 0.879
Weighted Average Recall: 0.825
Weighted Average F1 Score: 0.850

==============================
```

NLTK (train - test):
- Precision = 0.882
- Recall = 0.891
- F-measure = 0.878
- Accuracy = 0.891

```
==============================
NLTK :
Accuracy:0.891
Weighted Average Precision: 0.882
Weighted Average Recall: 0.891
Weighted Average F1 Score: 0.878

==============================
```

*CRF :
- Precision = 0.936
- Recall = 0.927
- F-measure = 0.936
- Accuracy = 0.928

```
==============================
*CRF :
Accuracy: 0.936
Weighted Average Precision: 0.927
Weighted Average Recall: 0.936
Weighted Average F1 Score: 0.928

==============================
```

資料預處理

Test: Section F Development: second half of Section H Training: everything else

- 論文作者之 github:
  https://github.com/GateNLP/broad_twitter_corpus/blob/master/README.md

篩去不合 iob 格式之資料列
*CRF : 大小寫轉換