# Natural Language Processing and Text Mining:

# HW#2

小組成員和負責工作：

109590041 范遠皓 測試和報告

109590043 柯瑞霖 撰寫程式

# 環境

使用的語言：Python

所需套件：
       click==8.1.3
       colorama==0.4.6
       joblib==1.2.0
       nltk==3.8.1
       numpy==1.24.2
       pandas==1.5.3
       python-dateutil==2.8.2
       pytz==2023.2
       regex==2023.3.23
       scikit-learn==1.2.2
       scipy==1.10.1
       six==1.16.0
       threadpoolctl==3.1.0
       tqdm==4.65.0

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

資料讀取並進行前置處理

```
66    #  處理完的數據
67    train_sents_f = []
68    for file_name in train_file_names:
69        train_sents = all_preprocess(file_name)
70        train_sents_f.append(train_sents)
71
72    test_sents_f = []
73    for file_name in test_file_names:
74        test_sents = all_preprocess(file_name)
75        test_sents_f.append(test_sents)
76
```

NLTK 前置處理:

對資料進行詞性標記，然後提取其中的名詞短語 (NP)，並將其特徵值插入原始資料中。

```
29    # 預處理-詞性標記
30    def preprocess(sent):
31        sent = nltk.word_tokenize(sent)
32        sent = nltk.pos_tag(sent) #詞性標記
33        return sent
34
35    # NLTK  預處理
36    def nltk_preprocess(data):
37        data_without_second_value = [(t[0]) for t in data]
38        data_tostrig = ' '.join(data_without_second_value)
39        data_preprocess = preprocess(data_tostrig)
40        pattern = 'NP: {<DT>?<JJ>*<NN>}' #正規表達式提取NP短語
41        cp = nltk.RegexpParser(pattern)
42        data_cs = cp.parse(data_preprocess)
43        data_iob_tagged = tree2conlltags(data_cs)
44        return data_iob_tagged
45
46    # 將 NLTK 取出的 NP 短語特徵值插入原始數據
47    def nltk_insert_data(nltk_preprocess_data, data_o):
48        i = 0
49        result = []
50        for data in data_o:
51            temp_l = list(data)
52            temp_l.insert(1, nltk_preprocess_data[i][1])
53            temp_t = tuple(temp_l)
54            result.append(temp_t)
55            if (i < len(data_o)-1):
56                i+=1
57        return result
58
59    # 整個預處理過程
60    def all_preprocess(file_name):
61        data_o = read_conll_file(file_name)
62        nltk_preprocess_result = nltk_preprocess(data_o)
63        nltk_insert_data_result = nltk_insert_data(nltk_preprocess_result, data_o)
64        return nltk_insert_data_result
```

提取特徵後放入 Perceptron 模型進行分類訓練，最後輸出結果。

```
125    # 特徵和標籤提取
126    X_train = extract_features(train_sents_f)
127    y_train = extract_labels(train_sents_f)
128    X_test = extract_features(test_sents_f)
129    y_test = extract_labels(test_sents_f)
130
131    # 特徵向量化
132    vec = DictVectorizer()
133    X_train = vec.fit_transform(X_train)
134    X_test = vec.transform(X_test)
135
136    # 感知模型訓練和預測
137    clf = Perceptron(verbose=10, n_jobs=-1, max_iter=5)
138    clf.fit(X_train, y_train)
139    y_pred = clf.predict(X_test)
140
141    # 報告結果輸出
142    print(classification_report(y_test, y_pred))
```

特徵提取

```
77    # 特徵提取函數
78    def word2features(sent, i):
79        word = sent[i][0]
80        features = {
81            'bias': 1.0,
82            'word.lower()': word.lower(),
83            'word[-3:]': word[-3:],
84            'word[-2:]': word[-2:],
85            'word.isupper()': word.isupper(),
86            'word.istitle()': word.istitle(),
87            'word.isdigit()': word.isdigit(),
88        }
89        if i > 0:
90            word1 = sent[i-1][0]
91            features.update({
92                '-1:word.lower()': word1.lower(),
93                '-1:word.istitle()': word1.istitle(),
94                '-1:word.isupper()': word1.isupper(),
95            })
96        else:
97            features['BOS'] = True
98        if i < len(sent)-1:
99            word1 = sent[i+1][0]
100            features.update({
101                '+1:word.lower()': word1.lower(),
102                '+1:word.istitle()': word1.istitle(),
103                '+1:word.isupper()': word1.isupper(),
104            })
105        else:
106            features['EOS'] = True
107        return features
108
109    # 特徵提取函數-應用於訓練和測試數據
110    def extract_features(sentences):
111        X = []
112        for sent in sentences:
113            for i in range(len(sent)):
114                X.append(word2features(sent, i))
115        return X
116
```

執行結果：

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| B-LOC        | 0.78      | 0.22   | 0.34     | 309     |
| B-ORG        | 0.30      | 0.30   | 0.30     | 736     |
| B-PER        | 0.85      | 0.65   | 0.74     | 2323    |
| I-LOC        | 0.20      | 0.31   | 0.24     | 114     |
| I-ORG        | 0.16      | 0.13   | 0.14     | 127     |
| I-PER        | 0.64      | 0.34   | 0.44     | 201     |
| O            | 0.96      | 0.99   | 0.97     | 25662   |
|              |           |        |          |         |
| accuracy     |           |        | 0.92     | 29472   |
| macro avg    | 0.55      | 0.42   | 0.45     | 29472   |
| weighted avg | 0.92      | 0.92   | 0.92     | 29472   |