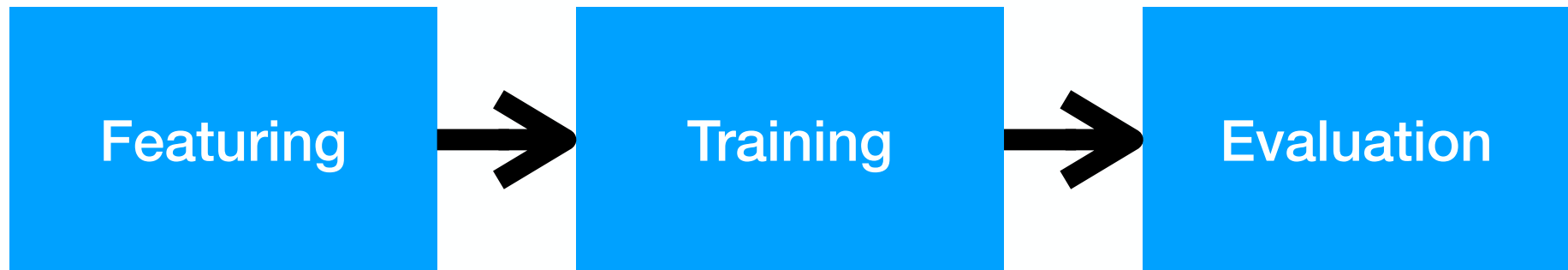


# Outline

- You have to train a Chinese NER model using [sklearn-crfsuite](#) module and evaluate its performance. Character-based Chinese NER dataset in json format will be given.
- Try to improve model performance using suggested methods and compare the results.

# Procedure



# sklearn-crfsuite

## Feature Engineering (Please refer to Featuring suggestions)

```
train_sents[0]
```

```
[('Melbourne', 'NP', 'B-LOC'),  
( '(', 'Fpa', 'O'),  
( 'Australia', 'NP', 'B-LOC'),  
( ')', 'Fpt', 'O'),  
( ',', 'Fc', 'O'),  
( '25', 'Z', 'O'),  
( 'may', 'NC', 'O'),  
( '(', 'Fpa', 'O'),  
( 'EFE', 'NC', 'B-ORG'),  
( ')', 'Fpt', 'O'),  
( '.', 'Fp', 'O')]
```

```
def word2features(sent, i):  
    word = sent[i][0]  
    postag = sent[i][1]  
  
    features = {  
        'bias': 1.0,  
        'word.lower()': word.lower(),  
        'word[-3:]': word[-3:],  
        'word[-2:]': word[-2:],  
        'word.isupper()': word.isupper(),  
        'word.istitle()': word.istitle(),  
        'word.isdigit()': word.isdigit(),  
        'postag': postag,  
        'postag[:2]': postag[:2],  
    }  
  
    if i > 0:  
        word1 = sent[i-1][0]  
        postag1 = sent[i-1][1]  
        features.update({  
            '-1:word.lower()': word1.lower(),  
            '-1:word.istitle()': word1.istitle(),  
            '-1:word.isupper()': word1.isupper(),  
            '-1:postag': postag1,  
            '-1:postag[:2]': postag1[:2],  
        })  
    else:  
        features['BOS'] = True  
  
    if i < len(sent)-1:  
        word1 = sent[i+1][0]  
        postag1 = sent[i+1][1]  
        features.update({  
            '+1:word.lower()': word1.lower(),  
            '+1:word.istitle()': word1.istitle(),  
            '+1:word.isupper()': word1.isupper(),  
            '+1:postag': postag1,  
            '+1:postag[:2]': postag1[:2],  
        })  
    else:  
        features['EOS'] = True  
  
    return features
```

# sklearn-crfsuite

## Assign feature as training data

```
%%time
X_train = [sent2features(s) for s in train_sents]
y_train = [sent2labels(s) for s in train_sents]

X_test = [sent2features(s) for s in test_sents]
y_test = [sent2labels(s) for s in test_sents]
```

## Training

```
%%time
crf = sklearn_crfsuite.CRF(
    algorithm='lbfgs',
    c1=0.1,
    c2=0.1,
    max_iterations=100,
    all_possible_transitions=True
)
crf.fit(X_train, y_train)
```

# sklearn-crfsuite

## Evaluation

```
y_pred = crf.predict(X_test)
metrics.flat_f1_score(y_test, y_pred,
                      average='weighted', labels=labels)
```

```
0.76980231377134023
```

# sklearn-crfsuite

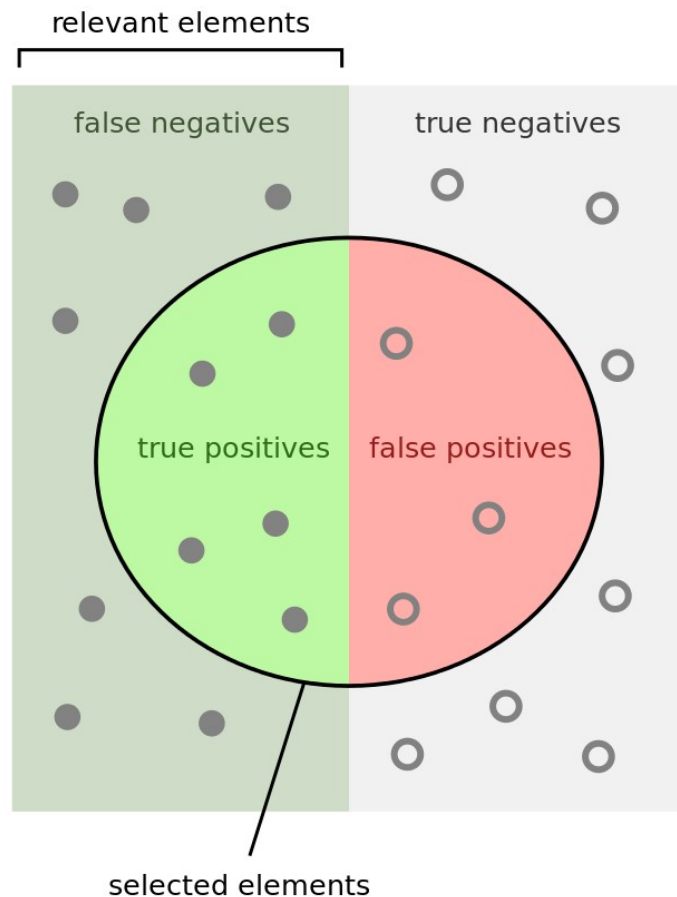
## Evaluation - more detail

```
# group B and I results
sorted_labels = sorted(
    labels,
    key=lambda name: (name[1:], name[0])
)
print(metrics.flat_classification_report(
    y_test, y_pred, labels=sorted_labels, digits=3
))
```

	precision	recall	f1-score	support
B-LOC	0.775	0.757	0.766	1084
I-LOC	0.601	0.631	0.616	325
B-MISC	0.698	0.499	0.582	339
I-MISC	0.644	0.567	0.603	557
B-ORG	0.795	0.801	0.798	1400
I-ORG	0.831	0.773	0.801	1104
B-PER	0.812	0.876	0.843	735
I-PER	0.873	0.931	0.901	634
avg / total	0.779	0.764	0.770	6178

# sklearn-crfsuite

## Evaluation



		Actual Value (as confirmed by experiment)	
		positives	negatives
Predicted Value (predicted by the test)	positives	<b>TP</b> True Positive	<b>FP</b> False Positive
	negatives	<b>FN</b> False Negative	<b>TN</b> True Negative

How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

$$F_1 = \left( \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} \right) = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

## Char-based

```
正 0
如 0
宋 B_Time
代 I_Time
詩 B_Person
人 I_Person
高 I_Person
翥 I_Person
所 0
雲 0
“ 0
南 B_Location
北 I_Location
山 I_Location
頭 I_Location
多 0
墓 B_Thing
田 I_Thing
, 0
清 B_Time
明 I_Time
祭 0
掃 0
各 0
紛 0
然 0
。 0
```

## Word-based

```
正 D
如 P
宋代 Nd
詩人高翥 Nb
所 Nc
雲 Na
“ FW
南北山頭 Nc
多 D
墓田 Nb
, COMMACATEGORY
清明 Nd
祭掃 VC
各 Nes
紛然 VH
。 PERIODCATEGORY
```

Char-based: 27 chars; 14 positive, 13 negative.

Word-based: 16 words; 5 positive, 11 negative.



# dataset

清明是人們祭掃先人，懷念追思的日子。  
正如宋代詩人高翥所雲“南北山頭多墓田，清明祭掃  
紙灰飛作白蝴蝶，淚血染成紅杜鵑”。  
凡清明之時，總是屢屢哀思湧上心頭，對母親懷念的  
母親姓孫名諱秋蘭，所以，我特別喜歡陳毅“幽蘭在

- train\_char.json
- validation\_char.json
- test\_char.json

[token, pos, neLabel]

```
[('正', 'D', '0'),  
( '如', 'P', '0'),  
( '宋', 'Nd', 'B_Time'),  
( '代', 'Nd', 'I_Time'),  
( '詩', 'Na', 'B_Person'),  
( '人', 'Na', 'I_Person'),  
( '高', 'Nb', 'I_Person'),  
( '翥', 'Nb', 'I_Person'),  
( '所', 'Nc', '0'),  
( '雲', 'Na', '0'),  
( '"', 'FW', '0'),  
( '南', 'Ncd', 'B_Location'),  
( '北', 'Ncd', 'I_Location'),  
( '山', 'Nc', 'I_Location'),  
( '頭', 'Nc', 'I_Location'),  
( '多', 'Nb', '0'),  
( '墓', 'Nb', 'B_Thing'),  
( '田', 'Nb', 'I_Thing'),  
( ' ', 'COMMACATEGORY', '0'),  
( '清', 'Nd', 'B_Time'),  
( '明', 'Nd', 'I_Time'),  
( '祭', 'VC', '0'),  
( '掃', 'VC', '0'),  
( '各', 'Nes', '0'),  
( '紛', 'VH', '0'),  
( '然', 'VH', '0'),  
( '。', 'PERIODCATEGORY', '0')]
```

# dataset

清明是人們祭掃先人，懷念追思的日子。  
正如宋代詩人高翥所雲“南北山頭多墓田，清明祭掃各紛然。  
紙灰飛作白蝴蝶，淚血染成紅杜鵑”。  
凡清明之時，總是屢屢哀思湧上心頭，對母親懷念的情愫越發細膩綿長。  
母親姓孫名諱秋蘭，所以，我特別喜歡陳毅“幽蘭在山谷，本自無人識。

- train\_v
- valida
- test\_v

```
[["清明 _ Time", "是 SHI 0", "人們 _ Person", "祭掃 VC 0",  
"先人 _ Person", "， COMMACATEGORY 0", "懷念 VJ 0", "追思  
VK 0", "的 DE 0", "日子 _ Time", "。 PERIODCATEGORY 0"],  
["正 D 0", "如 P 0", "宋代 _ Time", "詩人高翥 _ Person", "  
所 Nc 0", "雲 Na 0", "“ FW 0", "南北山頭 _ Location", "多  
VH 0", "墓田 _ Thing", "， COMMACATEGORY 0", "清明 _ Time"  
， "祭掃 VC 0", "各 Nes 0", "紛然 VH 0", "。 PERIODCATEGORY  
0"], ["紙灰飛 _ Thing", "作 VC 0", "白蝴蝶 _ Thing", "，
```

[token, pos, neLabel]