

自然語言處理

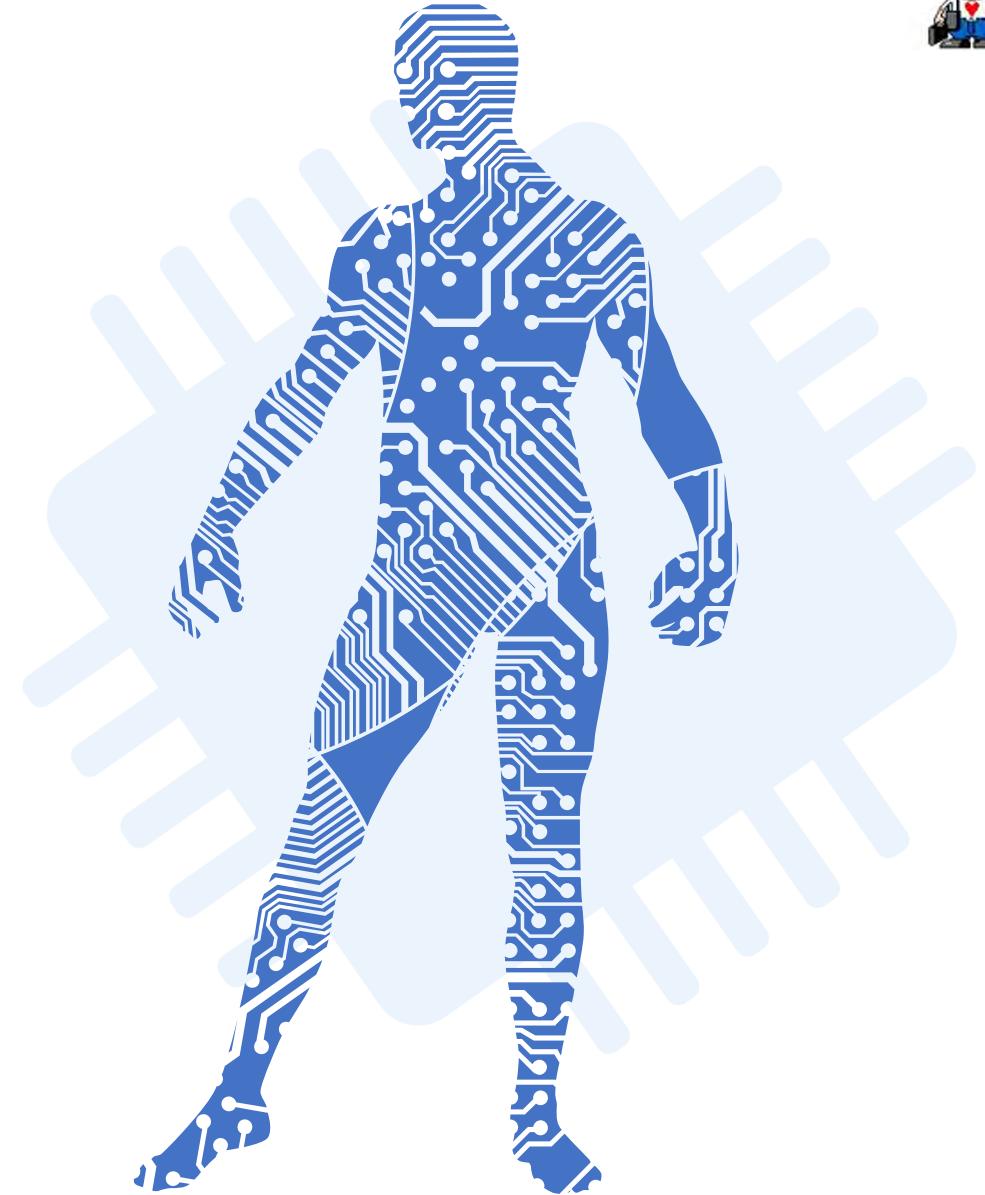
第 7 章 前處理

講師：紀俊男



本章大綱

- 環境設定
- 語系偵測 (Language Detection)
- 斷句 (Sentence Segmentation)
- 顏文字處理 (Emoji Processing)
- 拼寫檢查 (Spelling Check)
- 中文繁簡轉換 (Traditional-Simplified Chinese)
- 斷詞與文本清理 (Tokenization & Text Cleaning)
- 英文文本標準化 (Text Normalization)
- 詞性標註 (Part-of-Speech Tagging) PoS
- 名稱實體識別 (Named Entity Recognition) NER





環境設定

範例完整原始碼：
<https://url.cc/3Jk69b>





原始程式碼



1 # 安裝 NLTK 套件
2 !pip install nltk

3
4 # 安裝 spaCy 套件
5 !pip install spacy

6
7 # 下載中英文語言模型
8 # spaCy 語言模型首頁：<https://github.com/explosion/spacy-models>

9
10 # 下載英文語言模型 (sm = small)
11 !python -m spacy download en_core_web_sm

small

12
13 # 下載中文語言模型 (md = middle)
14 !python -m spacy download zh_core_web_md

medium

15
16 chinese_text = """

17 在這個資訊爆炸的時代，自然語言處理(Natural Language Processing, NLP)技術正在迅速發展，成為人工智慧領域的一個重要分支。NLP的應用範圍廣泛，包括語音識別、機器翻譯、情感分析、文本摘要等。隨著深度學習技術的進步，NLP的研究和應用正在不斷突破傳統的限制，為人們的生活和工作帶來了許多便利。

18
19 自然語言處理涉及到多個步驟，其中包括語言檢測、文本清理、文本標準化、斷句、斷詞、詞性標註等。這些前處理步驟對於後續的分析和模型訓練至關重要。例如，文本清理可以去除無用的標點符號和特殊字符，斷詞則將句子分割成單個的詞彙，便於進一步的處理。

20
21 在NLP的應用中，一個常見的任務是情感分析，即判斷文本中所表達的情感是正面還是負面。這在社交媒體分析、市場研究等領域有著廣泛的應用。另一個重要的應用是機器翻譯，隨著全球化的發展，準確高效的翻譯工具變得越來越重要。

22
23 然而，NLP也面臨著許多挑戰，其中之一是語言的多樣性和複雜性。每種語言都有其獨特的語法結構和詞彙規則，這使得開發通用的NLP模型變得困難。此外，不同語境下同一詞彙的含義可能會有所不同，這就需要模型能夠理解上下文資訊。

24
25 總的來說，自然語言處理是一個充滿挑戰和機遇的領域。隨著技術的不斷進步，我們有理由相信，NLP將在未來發揮更大的作用，為人類社會帶來更多的便利和進步。
"""

Natural Language Tool Kit，能針對英文執行
斷句、斷詞、詞性標註...等處理的優秀套件

老牌英文函式庫

spaCy 需要根據處理的語系，
額外下載語言模型安裝包。



隨堂練習：環境設定

- 請將前一頁的原始碼拷貝貼上，並且執行看看：

```
1 # 安裝 NLTK 套件
2 !pip install nltk
3
4 # 安裝 spaCy 套件
5 !pip install spacy
6
7 # 下載中英文語言模型
8 # spaCy 語言模型首頁： https://github.com/explosion/spacy-models
9
10 # 下載英文語言模型 (sm = small)
11 !python -m spacy download en_core_web_sm
12
13 # 下載中文語言模型 (md = middle)
14 !python -m spacy download zh_core_web_md
15
16 chinese_text = """
17 在這個資訊爆炸的時代，自然語言處理(Natural Language Processing, NLP)技術正在迅速發展，成為人工智慧領域的一個重要分支。NLP的應用範圍廣泛，包括語音識別、機器翻譯、情感分析、文本摘要等。隨著深度學習技術的進步，NLP的研究和應用正在不斷突破傳統的限制，為人們的生活和工作帶來了許多便利。
18
19 自然語言處理涉及到多個步驟，其中包括語言檢測、文本清理、文本標準化、斷句、斷詞、詞性標註等。這些前處理步驟對於後續的分析和模型訓練至關重要。例如，文本清理可以去除無用的標點符號和特殊字符，斷詞則將句子分割成單個的詞彙，便於進一步的處理。
20
21 在NLP的應用中，一個常見的任務是情感分析，即判斷文本中所表達的情感是正面還是負面。這在社交媒體分析、市場研究等領域有著廣泛的應用。另一個重要的應用是機器翻譯，隨著全球化的發展，準確高效的翻譯工具變得越來越重要。
22
23 然而，NLP也面臨著許多挑戰，其中之一是語言的多樣性和複雜性。每種語言都有其獨特的語法結構和詞彙規則，這使得開發通用的NLP模型變得困難。此外，不同語境下同一詞彙的含義可能會有所不同，這就需要模型能夠理解上下文資訊。
24
25 總的來說，自然語言處理是一個充滿挑戰和機遇的領域。隨著技術的不斷進步，我們有理由相信，NLP將在未來發揮更大的作用，為人類社會帶來更多的便利和進步。
26 """
```





語系偵測
Language Detection



何謂「語系偵測」？



- 偵測「文本（Text）」是哪一種語言？

“你好！我叫李傑克。”



中文（zh）

“Hello! My name is Jack Lee.”



英文（en）





為何要做「語系偵測」？



- 不同語系、不同前處理步驟 / 函式庫

英文 (en)

詞語標準化 (Text Normalization)

- 取動詞原型 (Lemmatization)
- 轉成小寫 (Lower cases)

斷詞 (Tokenization)

- import nltk

中文 (zh)

詞語標準化 (Text Normalization)

- 繁簡轉換

斷詞 (Tokenization)

- import jieba





語系偵測實作



```
1 # 套件安裝 (Package Installation)
2 !pip install langid      language id
3
4 # 文本 (Text) 設定
5 chinese_text = "你好！我的名字叫做李傑克！"
6 english_text = "Hello! My name is Jack Lee."
7
8 # 語系偵測 (分數未標準化)
9 import langid
10
11 result = langid.classify(chinese_text)
12 print(result)
```

✖

('zh', -153.1489453315735)

常見語系代碼 (ISO-639 標準)

Language Code	Language Name	Language Code	Language Name
zh	Chinese	ko	Korean
en	English	ru	Russian
es	Spanish	pt	Portuguese
fr	French	it	Italian
de	German	ar	Arabic
ja	Japanese	tr	Turkish
nl	Dutch	pl	Polish
sv	Swedish	da	Danish
fi	Finnish	no	Norwegian
el	Greek	he	Hebrew
th	Thai	id	Indonesian
vi	Vietnamese		

作者自訂之 Confidence Score ($\log_{10} \frac{1}{\text{分數}}$)





隨堂練習：語系偵測



- 請先撰寫好前一頁的原始碼，並且執行看看：

```
1 # 套件安裝 (Package Installation)
2 !pip install langid
3
4 # 文本 (Text) 設定
5 chinese_text = "你好！我的名字叫做李傑克！"
6 english_text = "Hello! My name is Jack Lee."
7
8 # 語系偵測 (分數未標準化)
9 import langid
10
11 result = langid.classify(chinese_text)
12 print(result)
```

('zh', -153.1489453315735)





將分數標準化



載入必要的套件

```
1 # 語系偵測（分數已標準化）  
2 from langid.langid import LanguageIdentifier, model  
3  
        自建語系偵測物件          載入預訓練好的模型    標準化分數為 0~1  
4 detector = LanguageIdentifier.from_modelstring(model, norm_probs=True)  
5  
        偵測語系  
6 result = detector.classify(chinese_text)  
7 print(result)
```

```
('zh', 0.9999999999667954)
```





指定候選語系



```
1 # 語系偵測（指定候選語系 + 分數標準化）
2 import langid          載入必要的套件
3 from langid.langid import LanguageIdentifier, model
4
5 langid.set_languages(['en','zh']) 指定候選語系
6 detector = LanguageIdentifier.from_modelstring(model, norm_probs=True)    自建語系偵測器 + 標準化分數
7
8 result = detector.classify(chinese_text)
9 print(result)
10
11 result = detector.classify(english_text)
12 print(result)
```

('zh', 0.9999999999667954)

('en', 0.9977290649556259)



隨堂練習：語系偵測



- 請先撰寫好前兩頁的原始碼，並且執行看看：

```
1 # 語系偵測（分數已標準化）
2 from langid.langid import LanguageIdentifier, model
3
4 detector = LanguageIdentifier.from_modelstring(model, norm_probs=True)
5
6 result = detector.classify(chinese_text)
7 print(result)
8
9         # 語系偵測（指定候選語系 + 分數標準化）
10        import langid
11        from langid.langid import LanguageIdentifier, model
12
13        langid.set_languages(['en', 'zh'])
14        detector = LanguageIdentifier.from_modelstring(model, norm_probs=True)
15
16        result = detector.classify(chinese_text)
17        print(result)
18
19        result = detector.classify(english_text)
20        print(result)
```





斷句

Sentence Segmentation



英文斷句

```
1 # 載入 NLTK
2 import nltk
3
4 # 下載斷句用模型 (punkt = Punctuation)
5 nltk.download('punkt')
6
7 # 利用 sent_tokenize 斷句
8 from nltk.tokenize import sent_tokenize
9
10 text = "Hello, world! Welcome to the world of NLP. This is an example of sentence tokenization."
11 sentences = sent_tokenize(text)
12 python list
13 print(sentences)      ['Hello, world!', 'Welcome to the world of NLP.', 'This is an example of sentence tokenization.']}
```





隨堂練習：英文斷句

- 請先撰寫好前一頁的原始碼，並且執行看看：

```
1 # 載入 NLTK
2 import nltk
3
4 # 下載斷句用模型 (punkt = Punctuation)
5 nltk.download('punkt')
6
7 # 利用 sent_tokenize 斷句
8 from nltk.tokenize import sent_tokenize
9
10 text = "Hello, world! Welcome to the world of NLP. This is an example of sentence tokenization."
11 sentences = sent_tokenize(text)
12
13 print(sentences)
```

['Hello, world!', 'Welcome to the world of NLP.', 'This is an example of sentence tokenization.']





中文斷句



```
1  1 import spacy  
2  
3  # 載入中文模型  
4  nlp = spacy.load("zh_core_web_md")  
5  
6  # 使用 spaCy 進行斷句  
7  doc = nlp(chinese_text)  
8  
9  # 去除「空白行」與「行末換行符號 \n」  
10 sentences = [sent.text.strip() for sent in doc.sents if sent.text.strip()]  
11  
12 # 輸出結果  
13 print(sentences)
```

[‘在這個資訊爆炸的時代，自然語言處理(Natural Language Processing, NLP)技術正在迅速發展，成為人工智慧領域的一個重要分支。’，
‘NLP的應用範圍廣泛，包括語音識別、機器翻譯、情感分析、文本摘要等。’，
‘隨著深度學習技術的進步，NLP的研究和應用正在不斷突破傳統的限制，為人們的生活和工作帶來了許多便利。’，
‘自然語言處理涉及到多個步驟，其中包括語言檢測、文本清理、文本標準化、斷句、斷詞、詞性標註等。’，
‘這些前處理步驟對於後續的分析和模型訓練至關重要。’，
‘例如，文本清理可以去除無用的標點符號和特殊字符，斷詞則將句子分割成單個的詞彙，便於進一步的處理。’，
‘在NLP的應用中，一個常見的任務是情感分析，即判斷文本中所表達的情感是正面還是負面。’，
‘這在社交媒體分析、市場研究等領域有著廣泛的應用。’，
‘另一個重要的應用是機器翻譯，隨著全球化的發展，準確高效的翻譯工具變得越來越重要。’，
‘然而，NLP也面臨著許多挑戰，其中之一是語言的多樣性和複雜性。’，
‘每種語言都有其獨特的語法結構和詞彙規則，這使得開發通用的NLP模型變得困難。’，
‘此外，不同語境下同一詞彙的含義可能會有所不同，這就需要模型能夠理解上下文資訊。’，
‘總的來說，自然語言處理是一個充滿挑戰和機遇的領域。’，
‘隨著技術的不斷進步，我們有理由相信，NLP將在未來發揮更大的作用，為人類社會帶來更多的便利和進步。’]





隨堂練習：中文斷句



- 請先撰寫好前兩頁的原始碼，並且執行看看：

```
1 import spacy  
2  
3 # 載入中文模型  
4 nlp = spacy.load("zh_core_web_md")  
5  
6 # 使用 spaCy 進行斷句  
7 doc = nlp(chinese_text)  
8  
9 # 去除「空白行」與「行末換行符號 \n」  
10 sentences = [sent.text.strip() for sent in doc.sents if sent.text.strip()]  
11  
12 # 輸出結果  
13 print(sentences)
```

removes leading & trailing whitespaces

- 段 - 句

↑

list comprehension





顏文字處理
Emoji Processing



顏文字 (Emoji) 簡介



• 顏文字是什麼？



• 顏文字的危害與處理

車禍骨折來住院一週 :man: 林希鼎Doctors 技術很棒 幽默:smile:然後骨科
醫師的護理人員團隊..細心:thumbs_up: B1有美容院坐洗髮還滿特別的:heart:

未知字元，造成「無法判斷」或「誤判」

車禍骨折來住院一週 :man: 林希鼎Doctors 技術很棒 幽默:smile:然後骨科
醫師的護理人員團隊..細心:thumbs_up: B1有美容院坐洗髮還滿特別的:heart:

轉成純文字，勉強可以判讀或處理





原始程式碼

```
1 # 安裝 emoji 套件
2 !pip install emoji
3
4 # 引入 emoji 套件
5 import emoji
6
7 # 用 demojize() 將顏文字轉成純文字
8 print(emoji.demojize('Python is ❤'))
```

Python is :red_heart:





隨堂練習：顏文字處理



- 請先**撰寫**好前一頁的**原始碼**，並且**執行**看看：

```
1 # 安裝 emoji 套件
2 !pip install emoji
3
4 # 引入 emoji 套件
5 import emoji
6
7 # 用 demojize() 將顏文字轉成純文字
8 print(emoji.demojize('Python is ❤'))
```

Python is :red_heart:





拼寫檢查
Spelling Check



英文拼寫檢查



```

1  1 # 安裝 language-tool-python 套件
2  2 !pip install language-tool-python
3
4  4 # 載入 language_tool_python 套件
5  5 import language_tool_python
6
7  7 # 指定拼寫檢查之語系
8  8 tool = language_tool_python.LanguageTool('en-US')
9
10 10 # 開始拼寫檢查
11 11 text = 'Speling erors in sentense are anoying.'
12 12 matches = tool.check(text)
13
14 14 # 印出拼寫檢查結果之資料結構
15 15 print(matches)
16
17 17 # 以拼寫檢查結果修正原文並印出
18 18 corrected_text = language_tool_python.utils.correct(text, matches)
19 19 print(corrected_text)

```

Language	Code	Language	Code	Language	Code
English	en	Dutch	nl	French	fr
German	de	Spanish	es	Portuguese	pt
Arabic	ar	Asturian	ast	Belarusian	be
Breton	br	Catalan	ca	Chinese	zh
Danish	da	Esperanto	eo	Galician	gl
Greek	el	Irish	ga	Italian	it
Japanese	ja	Khmer	km	Norwegian	no
Persian	fa	Polish	pl	Romanian	ro
Russian	ru	Slovak	sk	Slovenian	sl
Swedish	sv	Tagalog	tl	Tamil	ta
Ukrainian	uk				

[Match({'ruleId': 'MORFOLOGIK_RULE_EN_US',
 'message': 'Possible spelling mistake found.',
 'replacements': ['Spelling', 'Spewning', 'Spieling'],
 'offsetInContext': 0,
 'context': 'Speling erors in sentense are anoying.',
 'offset': 0, 'errorLength': 7,
 'category': 'TYPO', 'ruleIssueType': 'misspelling',
 'sentence': 'Speling erors in sentense are anoying.'}),
 Match({'ruleId': 'MORFOLOGIK_RULE_EN_US',
 ...
 }]

I
typo

Spelling errors in sentence are annoying.



隨堂練習：英文拼寫檢查



- 請先**撰寫**好前一頁的**原始碼**，並且**執行**看看：

```
1 # 安裝 language-tool-python 套件
2 !pip install language-tool-python
3
4 # 載入 language_tool_python 套件
5 import language_tool_python
6
7 # 指定拼寫檢查之語系
8 tool = language_tool_python.LanguageTool('en-US')
9
10 # 開始拼寫檢查
11 text = 'Speling erors in sentense are anoying.'
12 matches = tool.check(text)
13
14 # 印出拼寫檢查結果之資料結構
15 print(matches)
16
17 # 以拼寫檢查結果修正原文並印出
18 corrected_text = language_tool_python.utils.correct(text, matches)
19 print(corrected_text)
```





中文拼寫檢查



1 # 安裝 PyCorrector 套件
2 !pip install pycorrector

2 # 載入 MLM (Masked Language Model) as Correction using BERT (MacBERT) 模型
3 from pycorrector import MacBertCorrector
4 corrector = MacBertCorrector()

5 # 測試用中文文本
6 error_sentences = [
7 '今天新情很好',
8 '我遇到一位老友跟我療天。',
9 '他們只能有兩個選擇：接受降薪或自動離職。'
10]

11 # 進行拼寫檢查
12 batch_results = corrector.correct_batch(error_sentences)
13
14 # 印出結果
15 for result in batch_results:
16 print(result)
17 print(result['source'])
18 print(result['target'])
19 print("-----")

一次修正多句

(100句左右為佳)

↑ 把掉

X
model

B idirectional 上下文
E ncoder
R epresentations
from
T ransformers

```
{'source': '今天新情很好', 'target': '今天心情很好', 'errors': [('新', '心', 2)]}  
今天新情很好  
今天心情很好  
-----  
'source': '我遇到一位老友跟我療天。', 'target': '我遇到一位老友跟我聊天。', 'errors': [('療', '聊', 9)]}  
我遇到一位老友跟我療天。  
我遇到一位老友跟我聊天。  
-----  
'source': '他們只能有兩個選擇：接受降薪或自動離職。',  
'target': '他們只能有兩個選擇：接受降薪或自動離職。', 'errors': [('薪', '薪', 13)]}  
他們只能有兩個選擇：接受降薪或自動離職。  
他們只能有兩個選擇：接受降薪或自動離職。
```

一個錯 = 一筆記錄



隨堂練習：中文拼寫檢查

- 請先撰寫好前一頁的原始碼，並且執行看看：

```
1 # 安裝 PyCorrector 套件
2 !pip install pycorrector
3
4 # 載入 MLM (Masked Language Model) as Correction using BERT (MacBERT) 模型
5 from pycorrector import MacBertCorrector
6 corrector = MacBertCorrector()
7
8 # 測試用中文文本
9 error_sentences = [
10     '今天新情很好',
11     '我遇到一位老友跟我療天。',
12     '他們只能有兩個選擇：接受降薪或自動離職。'
13 ]
14
15 # 進行拼寫檢查
16 batch_results = corrector.correct_batch(error_sentences)
17
18 # 印出結果
19 for result in batch_results:
20     print(result)
21     print(result['source'])
22     print(result['target'])
23     print("-----")
```

```
{'source': '今天新情很好', 'target': '今天心情很好', 'errors': [('新', '心', 2)]}
今天新情很好
今天心情很好
-----
{'source': '我遇到一位老友跟我療天。', 'target': '我遇到一位老友跟我聊天。', 'errors': [('療', '聊', 9)]}
我遇到一位老友跟我療天。
我遇到一位老友跟我聊天。
-----
{'source': '他們只能有兩個選擇：接受降薪或自動離職。',
 'target': '他們只能有兩個選擇：接受降薪或自動離職。', 'errors': [('薪', '薪', 13)]}
他們只能有兩個選擇：接受降薪或自動離職。
他們只能有兩個選擇：接受降薪或自動離職。
```





中文繁簡轉換

Traditional-Simplified Chinese



中文繁簡轉換簡介

- 何謂「中文繁簡轉換」
 - “文本标准化” → “文本標準化”
- 為何需要「中文繁簡轉換」？
 - 某些 Python 函式庫僅接受簡體中文。
- 中文繁簡轉換的難度
 - 一對多對應： “后” → “后” / “後” ？
 - 用語不同： “信息” → “資訊”
- 解決方案
 - OpenCC 套件





原始程式碼



```
1 # 安裝並引入 OpenCC 套件
2 !pip install opencc
3 from opencc import OpenCC
4
5 # 定義一個專門做「繁簡轉換」的函數
6 def convert_chinese(text, target='simplified', dialect_convert=True):
7     # 根據傳入的 'target' 與 'dialect_convert' 參數，取得對應的 JSON 檔
8     configurations = {
9         ('simplified', True): 'tw2sp.json',
10        ('simplified', False): 'tw2s.json',
11        ('traditional', True): 's2twp.json',
12        ('traditional', False): 's2tw.json'
13    }
14     configuration = configurations.get((target, dialect_convert))
15
16     # 初始化對應的轉換器
17     converter = OpenCC(configuration) configuration = 'tw2sp.json'
18
19     # 轉換並將結果傳回去
20     return converter.convert(text)
21
22     # 呼叫轉換函數，並印出轉換結果
23     print(convert_chinese(chinese_text))
```

在這個資訊爆炸的時代，自然語言處理(Natural Language Processing, NLP)
技術正在迅速發展，成為人工智慧領域的一個重要分支。... (後略) ...

在这个信息爆炸的时代，自然语言处理(Natural Language Processing, NLP)
技术正在迅速发展，成为人工智能领域的一个重要分支。.... (后略) ...



隨堂練習：中文繁簡轉換



- 請先撰寫好前一頁的原始碼，並且執行看看：

```
1 # 安裝並引入 OpenCC 套件
2 !pip install opencc
3 from opencc import OpenCC
4
5 # 定義一個專門做「繁簡轉換」的函數
6 def convert_chinese(text, target='simplified', dialect_convert=True):
7     # 根據傳入的 'target' 與 'dialect_convert' 參數，取得對應的 JSON 檔
8     configurations = {
9         ('simplified', True): 'tw2sp.json',
10        ('simplified', False): 'tw2s.json',
11        ('traditional', True): 's2twp.json',
12        ('traditional', False): 's2tw.json'
13    }
14     configuration = configurations.get((target, dialect_convert))
15
16     # 初始化對應的轉換器
17     converter = OpenCC(configuration)
18
19     # 轉換並將結果傳回去
20     return converter.convert(text)
21
22 # 呼叫轉換函數，並印出轉換結果
23 print(convert_chinese(chinese_text))
```

在这个信息爆炸的时代，自然语言处理(Natural Language Processing, NLP)技术正在迅速发展，成为人工智能领域的一个重要分支。.... (後略) ...





斷詞與文本清理

Tokenization & Text Cleaning



英文斷詞與文本清理



```
1 # 引入必要函式庫
2 import nltk
3 from nltk.tokenize import word_tokenize
4 from nltk.corpus import stopwords
5 import string
6
7 # 下載 NLTK 的 Tokenizer 和 stopwords 資源
8 nltk.download('punkt')
9 nltk.download('stopwords')
10
11 # 獲取英語的停用詞列表
12 stop_words = set(stopwords.words('english'))
13
14 # 定義一個函數來進行斷詞、移除標點符號和停用詞
15 def tokenize_without_punctuation_and_stopwords(text):
16     # 使用 NLTK 的 word_tokenize 進行斷詞
17     tokens = word_tokenize(text)
18     # 移除標點符號和停用詞
19     filtered_tokens = [token for token in tokens
20         if token not in string.punctuation and token.lower() not in stop_words]
21     return filtered_tokens
22
23 # 測試文本
24 text = "Hello, how are you doing today?"
25
26 # 呼叫函數進行斷詞並移除標點和停用詞
27 tokens = tokenize_without_punctuation_and_stopwords(text)
28 print(tokens)
```

['Hello', 'today']

{ tuple
list
dictionary
set



隨堂練習：英文斷詞與文本清理



- 請先撰寫好前一頁的原始碼，並且執行看看：

```
1 # 引入必要函式庫
2 import nltk
3 from nltk.tokenize import word_tokenize
4 from nltk.corpus import stopwords
5 import string
6
7 # 下載 NLTK 的 Tokenizer 和 stopwords 資源
8 nltk.download('punkt')
9 nltk.download('stopwords')
10
11 # 獲取英語的停用詞列表
12 stop_words = set(stopwords.words('english'))
13
14 # 定義一個函數來進行斷詞、移除標點符號和停用詞
15 def tokenize_without_punctuation_and_stopwords(text):
16     # 使用 NLTK 的 word_tokenize 進行斷詞
17     tokens = word_tokenize(text)
18     # 移除標點符號和停用詞
19     filtered_tokens = [token for token in tokens
20         if token not in string.punctuation and token.lower() not in stop_words]
21     return filtered_tokens
22
23 # 測試文本
24 text = "Hello, how are you doing today?"
25
26 # 呼叫函數進行斷詞並移除標點和停用詞
27 tokens = tokenize_without_punctuation_and_stopwords(text)
28 print(tokens)
```

['Hello', 'today']





中文斷詞與文本清理

```
1 # 安裝 Jieba (結巴)
2 !pip install jieba
3
4 # 取得繁體中文自定義辭典 (可自行修改添加, 選用)
5 import os
6 Dictionary_File = 'dict.txt.big' 名字可任取
7
8 if not os.path.isfile(Dictionary_File):
9     os.system('wget https://raw.githubusercontent.com/cnchi/datasets/master/' + Dictionary_File)
10
11 # 取得繁體中文的「停止詞」辭典 (如：的、了、...)
12 StopWords_File = "stopWords_big5.txt"
13
14 if not os.path.isfile(StopWords_File):
15     os.system('wget https://raw.githubusercontent.com/cnchi/datasets/master/' + StopWords_File)
16
17 # 載入 Jieba
18 import jieba
19
20 # 載入繁體中文自定義辭典 (選用)
21 # jieba.set_dictionary(Dictionary_File)
```

紀俊男
心不要斷詞





中文斷詞與文本清理

```
1 23 # 設定標點符號與停止詞清單
2 24 punctuation = set("$/!&%\(\)+-*/_,.。？；‘＼<=>^`|~[]{}'0123456789?“”、《》！、：；？「」（）")
3 25 stopWords = set()
4 26 with open(Stopwords_File, "rt", encoding="utf-8") as f:
5     for line in f:
6         line = line.strip() # Remove trailing \n
7         stopWords.add(line)
8
9
10 31 # 將斷完句的文字，一次一句送入斷詞 & 文本清理
11 32 result = []
12 33 for sentence in sentences:
13     # 針對本句子斷詞
14     seg_sentence = list(jieba.cut(sentence))
15
16     # 移除標點符號
17     seg_sentence_no_punct = [word for word in seg_sentence if word not in punctuation]
18
19     # 移除停止詞
20     seg_sentence_no_stopwords = [word for word in seg_sentence_no_punct if word not in stopWords]
21
22     # 將結果加入 result 串列
23     result.append(seg_sentence_no_stopwords)
24
25 44 # 印出結果
26 45 print(result)
```

```
[['資訊', '爆炸', '時代', '自然', '語言', '處理', 'Natural',
 'Language', 'Processing', 'NLP', '技術', '正在', '迅速',
 ... (後略) ...]
 ]
```





隨堂練習：中文斷詞與文本清理



- 請先撰寫好前兩頁的原始碼，並且執行看看：





英文文本標準化

Text Normalization



英文文本需要哪些標準化處理？



小寫化
(Lower Cases)

The
↓
the

詞幹提取
(Stemming)

leaves
↓
leav

詞形還原
(Lemmatization)

leaves
↓
leaf
e.g. irregular
verbs, plurals





原始程式碼



```
1 # 載入必要套件
2 import nltk
3 from nltk.stem import PorterStemmer, WordNetLemmatizer
4 from nltk.corpus import wordnet
5
6 # 下載 NLTK 資源
7 nltk.download('punkt')
8 nltk.download('wordnet')
9 nltk.download('averaged_perceptron_tagger')
10
11 # 初始化 Stemmer 和 Lemmatizer
12 stemmer = PorterStemmer()
13 lemmatizer = WordNetLemmatizer()
14
15 # 定義一個函數來展示小寫、Stemming 和 Lemmatization
16 def process_text(text):
17     # 轉成小寫
18     text = text.lower()
19     # 使用 NLTK 的 word_tokenize 進行斷詞
20     words = nltk.word_tokenize(text)
21
22     # Stemming
23     stemmed_words = [stemmer.stem(word) for word in words] # 取字根
24
25     # Lemmatization
26     # 為了更精確的 Lemmatization, 我們需要提供正確的詞性 (Part-of-Speech)
27     lemmatized_words = [lemmatizer.lemmatize(word, pos=get_wordnet_pos(word)) for word in words]
28
29 return stemmed_words, lemmatized_words
```

word, sentence

一次取一個詞 # 取字根

POS Tagging.



原始程式碼

Python dictionary get()
prototype:



```
1 31 # NLTK 的詞性與 WordNet 詞性的對應
2 32 def get_wordnet_pos(word):
3 33     """Map NLTK's part-of-speech tags to those used by WordNet Lemmatizer"""
4 34     tag = nltk.pos_tag([word])[0][1][0].upper() ← 根據傳入的 word · 取得 "J" , "N" ,
5 35     tag_dict = {"J": wordnet.ADJ,           "V" , "R" 之一的詞性標註
6 36         "N": wordnet.NOUN,                 ← 定義 NLTK 與 WordNet 辭典的
7 37         "V": wordnet.VERB,                詞性標籤對應表
8 38         "R": wordnet.ADV}
9
10    return tag_dict.get(tag, wordnet.NOUN) ← 將 NLTK 詞性標籤 · 轉為 WordNet 標籤。
11
12 # 測試文本
13 text = "The leaves on the tree have fallen due to strong winds."
14
15 # 處理文本
16 stemmed, lemmatized = process_text(text)
17 print("Stemmed Words:", stemmed)
18 print("Lemmatized Words:", lemmatized)
```





隨堂練習：英文文本標準化



- 請先撰寫好前兩頁的原始碼，並且執行看看：

```
1 # 載入必要套件
2 import nltk
3 from nltk.stem import
4 from nltk.corpus import
5
6 # 下載 NLTK 資源
7 nltk.download('punkt')
8 nltk.download('wordnet')
9 nltk.download('averag
10
11 # 初始化 Stemmer 和 L
12 stemmer = PorterStem
13 lemmatizer = WordNetL
14
15 # 定義一個函數來展示小
16 def process_text(text
17     # 轉成小寫
18     text = text.lower()
19     # 使用 NLTK 的 wo
20     words = nltk.word
21
22     # Stemming
23     stemmed_words = [
24
25     # Lemmatization
26     # 為了更精確的 Len
27     lemmatized_words
28
29     return stemmed_words, lemmatized_words
```

```
31 # NLTK 的詞性與 WordNet 詞性的對應
32 def get_wordnet_pos(word):
33     """Map NLTK's part-of-speech tags to those used by WordNet Lemmatizer"""
34     tag = nltk.pos_tag([word])[0][1][0].upper()
35     tag_dict = {"J": wordnet.ADJ,
36                 "N": wordnet.NOUN,
37                 "V": wordnet.VERB,
38                 "R": wordnet.ADV}
39
40     return tag_dict.get(tag, wordnet.NOUN)
41
42 # 測試文本
43 text = "The leaves on the tree have fallen due to strong winds."
44
45 # 處理文本
46 stemmed, lemmatized = process_text(text)
47 print("Stemmed Words:", stemmed)
48 print("Lemmatized Words:", lemmatized)
```

#字根
并詞型還原





詞性標註

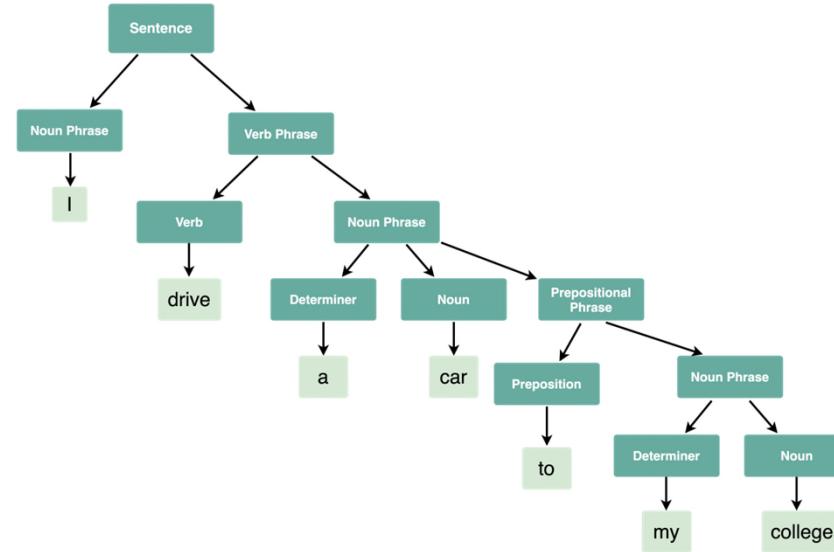
Part-Of-Speech Tagging



詞性標註簡介



- 何謂「詞性標註」(POS Tagging) ?
 - 將詞彙標上「**名詞、動詞、冠詞...**」等標記。
- 為何需要「詞性標註」?
 - 方便未來的「**句型分析**」(Sentence Analysis)，理解「**語意**」。





英文詞性標註



- 1 # 載入必要的函式庫 NLTK

```
1 # 載入必要的函式庫 NLTK
2 import nltk
3 from nltk.tokenize import word_tokenize
```
- 2 # 下載 NLTK 的 Tokenizer 和 POS Tagger 資源

```
5 # 下載 NLTK 的 Tokenizer 和 POS Tagger 資源
6 nltk.download('punkt')
7 nltk.download('averaged_perceptron_tagger')
```
- 3 # 定義一個函數來執行詞性標註

```
9 # 定義一個函數來執行詞性標註
10 def pos_tagging(text):
11     # 斷詞
12     words = word_tokenize(text)
13     # 執行詞性標註
14     pos_tags = nltk.pos_tag(words)
15     return pos_tags
```
- 4 # 測試文本

```
17 # 測試文本
18 text = "The quick brown fox jumps over the lazy dog."
```
- 5 # 獲得詞性標註結果

```
20 # 獲得詞性標註結果
21 tagged_text = pos_tagging(text)
22 print("POS Tagging Result:", tagged_text)
```

Tag	Description	Example	Tag	Description	Example
CC	Coordinating conjunction	and, but, or	CD	Cardinal number	one, two, 123
DT	Determiner	the, a, some	EX	Existential there	there is, there was
FW	Foreign word	d'hoevre	IN	Preposition or subordinating conjunction	in, of, like
JJ	Adjective	big, happy	JJR	Adjective, comparative	bigger, happier
JJS	Adjective, superlative	biggest, happiest	LS	List item marker	1, One:
MD	Modal	can, could, will	NN	Noun, singular or mass	dog, car
NNS	Noun, plural	dogs, cars	NNP	Proper noun, singular	John, London
NNPS	Proper noun, plural	Vikings, Americans	PDT	Predeterminer	all, both
POS	Possessive ending	's	PRP	Personal pronoun	I, you, he
PRP\$	Possessive pronoun	my, your, his	RB	Adverb	very, too, not
RBR	Adverb, comparative	better	RBS	Adverb, superlative	best
RP	Particle	up, off	SYM	Symbol	+, %, &
TO	to	to	UH	Interjection	ah, oops, ouch
VB	Verb, base form	eat, walk	VBD	Verb, past tense	ate, walked
VBG	Verb, gerund or present participle	eating, walking	VBN	Verb, past participle	eaten, walked
VBP	Verb, non-3rd person singular present	eat, walk	VBZ	Verb, 3rd person singular present	eats, walks
WDT	Wh-determiner	which, whatever	WP	Wh-pronoun	who, whoever
WP\$	Possessive wh-pronoun	whose	WRB	Wh-adverb	where, when

NLTK 使用 Penn Treebank 定義之詞性標籤

```
[('The', 'DT'), ('quick', 'JJ'), ('brown', 'NN'),
('fox', 'NN'), ('jumps', 'VBZ'), ('over', 'IN'),
('the', 'DT'), ('lazy', 'JJ'), ('dog', 'NN'), ('.', '.')]
```



隨堂練習：英文詞性標註



- 請先撰寫好前一頁的原始碼，並且執行看看：

```
1 # 載入必要的函式庫 NLTK
2 import nltk
3 from nltk.tokenize import word_tokenize
4
5 # 下載 NLTK 的 Tokenizer 和 POS Tagger 資源
6 nltk.download('punkt')
7 nltk.download('averaged_perceptron_tagger')
8
9 # 定義一個函數來執行詞性標註
10 def pos_tagging(text):
11     # 斷詞
12     words = word_tokenize(text)
13     # 執行詞性標註
14     pos_tags = nltk.pos_tag(words)
15     return pos_tags
16
17 # 測試文本
18 text = "The quick brown fox jumps over the lazy dog."
19
20 # 獲得詞性標註結果
21 tagged_text = pos_tagging(text)
22 print("POS Tagging Result:", tagged_text)
```

```
[('The', 'DT'), ('quick', 'JJ'), ('brown', 'NN'),
('fox', 'NN'), ('jumps', 'VBZ'), ('over', 'IN'),
('the', 'DT'), ('lazy', 'JJ'), ('dog', 'NN'), ('.', '.')]
```





中文詞性標註

```
1 # 載入 spaCy 套件
2 import spacy
3
4 # 載入中文模型
5 nlp = spacy.load("zh_core_web_md")
6
7 # 定義繁體中文的句子
8 text = "我來到位於台北的政治大學校園。"
9
10 # 處理文本
11 doc = nlp(text)
12
13 # 輸出每個詞的文本和詞性標籤
14 for token in doc:
15     print(f"{token.text} ({token.pos_})")
```

Tag	Description	Tag	Description
ADJ	adjective	ADP	adposition
ADV	adverb	AUX	auxiliary
CONJ	conjunction	CCONJ	coordinating conjunction
DET	determiner	INTJ	interjection
NOUN	noun	NUM	numeral
PART	particle	PRON	pronoun
PROPN	proper noun	PUNCT	punctuation
SCONJ	subordinating conjunction	SYM	symbol
VERB	verb	X	other
SPACE	space		

我 (PRON)
來到 (VERB)
位於 (ADP)
台北 (PROPN)
的 (PART)
政治 (NOUN)
大學 (NOUN)
校園 (NOUN)
。 (PUNCT)





隨堂練習：中文詞性標註



- 請先撰寫好前一頁的原始碼，並且執行看看：

```
1 # 載入 spaCy 套件
2 import spacy
3
4 # 載入中文模型
5 nlp = spacy.load("zh_core_web_md")
6
7 # 定義繁體中文的句子
8 text = "我來到位於台北的政治大學校園。"
9
10 # 處理文本
11 doc = nlp(text)
12
13 # 輸出每個詞的文本和詞性標籤
14 for token in doc:
15     print(f'{token.text} ({token.pos_})')
```

我 (PRON)
來到 (VERB)
位於 (ADP)
台北 (PROPN)
的 (PART)
政治 (NOUN)
大學 (NOUN)
校園 (NOUN)
。 (PUNCT)





實體名稱識別

Named Entity Recognition



實體名稱識別簡介



- 何謂「實體名稱」(Named Entity)

- 具有特殊意義的詞彙。
- 如：人名、地名、組織名...etc.

- 為何需要識別「實體名稱」？

- 抽取一篇文章中，具有特殊意義詞彙之關係。
- 如：「**台北101**是**台灣**的一個地標，位於**台北市**。」

↑
建築物名稱
(FAC) ↑
地區或國家
(NORP)

↑
城市名稱
(GPE)





英文實體名稱識別

```
1 # 載入 spaCy 套件
2 import spacy
3
4 # 載入英文的小型 NLP 模型
5 nlp = spacy.load("en_core_web_sm")
6
7 # 定義一個函數進行命名實體識別
8 def named_entity_recognition(text):
9     # 處理文本
10    doc = nlp(text)
11    # 提取命名實體
12    entities = [(ent.text, ent.label_) for ent in doc.ents]
13    return entities
14
15 # 測試文本
16 text = "Apple Inc. announced its new iPhone in San Francisco."
17
18 # 執行命名實體識別
19 entities = named_entity_recognition(text)
20
21 # 印出識別出的實體
22 print("Identified Entities:", entities)
```

Entity Type	Description
PERSON	People, including fictional.
NORP	Nationalities or religious or political groups.
FAC	Buildings, airports, highways, bridges, etc.
ORG	Companies, agencies, institutions, etc.
GPE	Countries, cities, states.
LOC	Non-GPE locations, mountain ranges, bodies of water.
PRODUCT	Objects, vehicles, foods, etc. (Not services.)
EVENT	Named hurricanes, battles, wars, sports events, etc.
WORK_OF_ART	Titles of books, songs, etc.
LAW	Named documents made into laws.
LANGUAGE	Any named language.
DATE	Absolute or relative dates or periods.
TIME	Times smaller than a day.
PERCENT	Percentage, including "%".
MONEY	Monetary values, including unit.
QUANTITY	Measurements, as of weight or distance.
ORDINAL	"first", "second", etc.
CARDINAL	Numerals that do not fall under another type.

```
[('Apple Inc.', 'ORG'),
 ('iPhone', 'PRODUCT'),
 ('San Francisco', 'GPE')]
```





隨堂練習：英文實體名稱識別



- 請先撰寫好前一頁的原始碼，並且執行看看：

```
1 # 載入 spaCy 套件
2 import spacy
3
4 # 載入英文的小型 NLP 模型
5 nlp = spacy.load("en_core_web_sm")
6
7 # 定義一個函數進行命名實體識別
8 def named_entity_recognition(text):
9     # 處理文本
10    doc = nlp(text)
11    # 提取命名實體
12    entities = [(ent.text, ent.label_) for ent in doc.ents]
13    return entities
14
15 # 測試文本
16 text = "Apple Inc. announced its new iPhone in San Francisco."
17
18 # 執行命名實體識別
19 entities = named_entity_recognition(text)
20
21 # 印出識別出的實體
22 print("Identified Entities:", entities)
```

[('Apple Inc.', 'ORG'),
 ('iPhone', 'PRODUCT'),
 ('San Francisco', 'GPE')]





中文實體名稱識別

```
1 # 載入 spaCy 套件
2 import spacy
3
4 # 載入中文的中型 NLP 模型
5 nlp = spacy.load("zh_core_web_md")
6
7 # 定義繁體中文的句子
8 text = "台北101是台灣的一個地標，位於台北市。"
9
10 # 處理文本
11 doc = nlp(text)
12
13 # 印出識別到的命名實體
14 for ent in doc.ents:
15     print(f"{ent.text} ({ent.label_})")
```

台北 (GPE)
101 (CARDINAL)
台灣 (GPE)
地標 (LOC)
台北市 (GPE)

Entity Type	Description
PERSON	People, including fictional.
NORP	Nationalities or religious or political groups.
FAC	Buildings, airports, highways, bridges, etc.
ORG	Companies, agencies, institutions, etc.
GPE	Countries, cities, states.
LOC	Non-GPE locations, mountain ranges, bodies of water.
PRODUCT	Objects, vehicles, foods, etc. (Not services.)
EVENT	Named hurricanes, battles, wars, sports events, etc.
WORK_OF_ART	Titles of books, songs, etc.
LAW	Named documents made into laws.
LANGUAGE	Any named language.
DATE	Absolute or relative dates or periods.
TIME	Times smaller than a day.
PERCENT	Percentage, including "%".
MONEY	Monetary values, including unit.
QUANTITY	Measurements, as of weight or distance.
ORDINAL	"first", "second", etc.
CARDINAL	Numerals that do not fall under another type.





隨堂練習：中文實體名稱識別



- 請先撰寫好前一頁的原始碼，並且執行看看：

```
1 # 載入 spaCy 套件
2 import spacy
3
4 # 載入中文的中型 NLP 模型
5 nlp = spacy.load("zh_core_web_md")
6
7 # 定義繁體中文的句子
8 text = "台北101是台灣的一個地標，位於台北市。"
9
10 # 處理文本
11 doc = nlp(text)
12
13 # 印出識別到的命名實體
14 for ent in doc.ents:
15     print(f"{ent.text} ({ent.label_})")
```

台北 (GPE)
101 (CARDINAL)
台灣 (GPE)
地標 (LOC)
台北市 (GPE)

