

✓ Lab#1, NLP Spring 2025

This is due on 2025/03/10 16:00, commit to your github as a PDF (lab1.pdf) (File>Print>Save as PDF).

IMPORTANT: After copying this notebook to your Google Drive, please paste a link to it below. To get a publicly-accessible link, hit the *Share* button at the top right, then click "Get shareable link" and copy over the result. If you fail to do this, you will receive no credit for this lab!

LINK: *paste your link here*

<https://colab.research.google.com/drive/xxxxxxxxxx>

Student ID:

Name:

✓ Question 1 (100 points)

Let's switch over to coding! Write some code in this cell to compute the number of unique word **tokens** in this paragraph (5 steps of Text Normalisation: 1. Lowercase Conversion, 2. Remove punctuations, 3. Stemming, 4. Lemmatisation, 5. Stopword Removal). Use a whitespace tokenizer to separate words (i.e., split the string by white space). Be sure that the cell's output is visible in the PDF file you turn in on Github.

按兩下(或按 Enter 鍵) 即可編輯

```
1 paragraph = '''Last night I dreamed I went to Manderley again. It seemed to me
2 that I was passing through the iron gates that led to the driveway.
3 The drive was just a narrow track now, its stony surface covered
4 with grass and weeds. Sometimes, when I thought I had lost it, it
5 would appear again, beneath a fallen tree or beyond a muddy pool
6 formed by the winter rains. The trees had thrown out new
7 low branches which stretched across my way. I came to the house
8 suddenly, and stood there with my heart beating fast and tears
9 filling my eyes.'''
10
11 # DO NOT MODIFY THE VARIABLES
12 tokens = 0
13 word_tokens = []
14
15 import nltk
16 from nltk.stem import PorterStemmer, LancasterStemmer, SnowballStemmer, WordNetLemmatizer
17 from nltk.corpus import stopwords
18 from nltk.tokenize import word_tokenize
19
20
21 nltk.download("stopwords")
22 nltk.download("punkt")
23 nltk.download("wordnet")
24
25 # 1. 小寫轉換 (Lowercase Conversion)
26 word_tokens = word_tokenize(paragraph)
27 tokens_lower = [token.lower() for token in word_tokens]
28
29 # 2. 移除標點符號 (Remove punctuations)
30 def remove_punct(tokens):
31     return [word for word in tokens if word.isalpha()]
32
33 clean_tokens = remove_punct(tokens_lower)
34
35 # 3. 語幹提取 (Stemming)
36 porter = PorterStemmer()
37 stemmed_port = [porter.stem(token) for token in clean_tokens]
38
39 lanc = LancasterStemmer()
40 stemmed_lanc = [lanc.stem(token) for token in clean_tokens]
41
42 snow = SnowballStemmer("english")
43 stemmed_snow = [snow.stem(token) for token in clean_tokens]
44
45
46 # 4. 詞形還原 (Lemmatisation)
47 lemmatiser = WordNetLemmatizer()
48 lemmatised = [lemmatiser.lemmatize(token) for token in stemmed_snow]
49
```

```
50 # 5. 停用詞去除 (Stopword Removal)
51 stop_words = set(stopwords.words("english"))
52 words_no_stop = [word for word in lemmatised if word not in stop_words]
53
54 word_tokens = set(words_no_stop)
55
56 tokens = len(word_tokens)
57
58
59 # DO NOT MODIFY THE BELOW LINE!
60 print('Number of word tokens: %d' % (tokens))
61 print("printing lists separated by commas")
62 print(*word_tokens, sep=", ")
63
```

↗ Number of word tokens: 51
printing lists separated by commas
sometim, new, way, thrown, night, pool, beyond, driveway, rain, surfac, fallen, grass, stoni, cover, form, across, sudden, heart, came, winter, l
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Package wordnet is already up-to-date!