

✓ 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/12KA2Pt7Kf36Em8vknY42Sjn2FU3Hj8Me?usp=sharing>

Student ID: B1129028

Name: 陳奕劭

✓ 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 from enum import unique
2 paragraph = '''Last night I dreamed I went to Manderley again. It seemed to me
3 that I was passing through the iron gates that led to the driveway.
4 The drive was just a narrow track now, its stony surface covered
5 with grass and weeds. Sometimes, when I thought I had lost it, it
6 would appear again, beneath a fallen tree or beyond a muddy pool
7 formed by the winter rains. The trees had thrown out new
8 low branches which stretched across my way. I came to the house
9 suddenly, and stood there with my heart beating fast and tears
10 filling my eyes.'''
11
12 # DO NOT MODIFY THE VARIABLES
13 tokens = 0
14 word_tokens = []
15
16 # YOUR CODE HERE! POPULATE THE tokens and word_tokens VARIABLES WITH THE CORRECT VALUES!
17
18 import nltk
19 from nltk.tokenize import word_tokenize
20 from nltk.corpus import stopwords
21 from nltk.stem import PorterStemmer, LancasterStemmer, SnowballStemmer, WordNetLemmatizer
22
23 nltk.download('punkt')
24 nltk.download('punkt_tab')
25 nltk.download('stopwords')
26 nltk.download('wordnet')
27
28 # Step 1: Lowercase Conversion
29 paragraph = paragraph.lower()
30
31 # Step 2: Tokenization
32 tokens = word_tokenize(paragraph)
33
34 # Step 3: Remove Punctuation
35 def remove_punctuation(text):
36     return [word for word in text if word.isalpha()]
37
38 tokens = remove_punctuation(tokens)
39
40 # Step 4: Stemming using PorterStemmer
41 stemmer = SnowballStemmer("english")
42 stemmed_tokens = [stemmer.stem(token) for token in tokens]
43
44 # Step 5: Lemmatization using WordNetLemmatizer
```

```
45 lemmatizer = WordNetLemmatizer()
46 lemmatized_tokens = [lemmatizer.lemmatize(token) for token in stemmed_tokens]
47
48 # Step 6: Stopword Removal
49 stop_words = set(stopwords.words("english"))
50 word_tokens = [word for word in lemmatized_tokens if word not in stop_words]
51
52 # Compute number of unique tokens
53 tokens = len(set(word_tokens))
54
55
56 # DO NOT MODIFY THE BELOW LINE!
57 print('Number of word tokens: %d' % (tokens))
58 print("printing lists separated by commas")
59 print(*word_tokens, sep = ", ")
```

```
↗ Number of word tokens: 51
printing lists separated by commas
last, night, dream, went, manderley, seem, wa, pas, iron, gate, led, driveway, drive, wa, narrow, track, stoni, surfac, cover, grass, weed, some
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]   Package punkt_tab is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
```