

✓ 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/12RhWFKoKnB5cyf9t8NgXtkekd9hv6tVe#scrollTo=SgNZTjrhcHa0>

Student ID:B1129014

Name:黎秉鑫

✓ 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 # YOUR CODE HERE! POPULATE THE tokens and word_tokens VARIABLES WITH THE CORRECT VALUES!
16
17 import nltk
18 from nltk.stem import PorterStemmer, LancasterStemmer, SnowballStemmer, WordNetLemmatizer
19 from nltk.tokenize import word_tokenize
20 from nltk.corpus import stopwords, wordnet
21 nltk.download('punkt', quiet=True)
22 nltk.download('punkt_tab', quiet=True)
23 nltk.download('stopwords', quiet=True)
24 nltk.download('wordnet', quiet=True)
25
26 # 1. Lowercase Conversion
27 paragraph_lower = paragraph.lower()
28
29 # 2. Remove punctuations
30 word_tokens = word_tokenize(paragraph_lower)
31 def remove_punct(token):
32     return [word for word in token if word.isalpha()]
33 punct_removed = remove_punct(word_tokens)
34
35 # print(*punct_removed, sep = ", ") # for debugging
36
37 # 3. Stemming
38 port = PorterStemmer()
39 stemmed_port = [port.stem(token) for token in punct_removed]
40 # lanc = LancasterStemmer()
41 # stemmed_lanc = [lanc.stem(token) for token in stemmed_port]
42 # snow = SnowballStemmer("english")
43 # stemmed_snow = [snow.stem(token) for token in stemmed_lanc]
44 stemmed = stemmed_port
45
46 # 4. Lemmatisation
47 from nltk.stem import WordNetLemmatizer
48 lemmatizer = WordNetLemmatizer()
49 lemmatised = [lemmatizer.lemmatize(token) for token in stemmed]
```

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

↻ Number of word tokens: 51
printing lists separated by commas
beat, narrow, would, cover, beneath, pas, weed, fast, went, stretch, appear, came, thrown, stood, stoni, driveway, iron, drive, surfac, wa, some

