# ▾ Lab#1, NLP Spring 2023

This is due on 2023/03/06 15:30, 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/10V35CX7Bw-JZ6KIVEw6cWBMmpwZBAE2x?usp=share_link

---

**Student ID**:B0928018

**Name**:呂哲睿

# ▾ 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

```
paragraph = '''Last  night  I  dreamed  I  went  to  Manderley  again.  It  seemed  to  me
that  I  was  passing  through  the  iron  gates  that  led  to  the  driveway.
The  drive  was  just  a  narrow  track  now,  its  stony  surface  covered
with  grass  and  weeds.  Sometimes,  when  I  thought  I  had  lost  it,  it
would  appear  again,  beneath  a  fallen  tree  or  beyond  a  muddy  pool
formed  by  the  winter  rains.  The  trees  had  thrown  out  new
low  branches  which  stretched  across  my  way.  I  came  to  the  house
suddenly,  and  stood  there  with  my  heart  beating  fast  and  tears
filling  my  eyes.'''

# DO NOT MODIFY THE VARIABLES
tokens = 0
word_tokens = []
```

```python
# YOUR CODE HERE! POPULATE THE tokens and word_tokens VARIABLES WITH THE CORRECT VALUE

#1
paragraph_lower = paragraph.lower()


#2
import nltk
nltk.download("punk")
def remove_punct(paragraph_lower):
    return [word for word in paragraph_lower if word.isalpha]
sent = remove_punct(paragraph_lower)

#3
from nltk.stem import PorterStemmer, LancasterStemmer, SnowballStemmer


port = PorterStemmer()
stemmed_port = [port.stem(token) for token in sent]

lanc = LancasterStemmer()
stemmed_lanc = [lanc.stem(token) for token in sent]

snow = SnowballStemmer("english")
stemmed_snow = [snow.stem(token) for token in sent]

#4
from nltk.stem import WordNetLemmatizer
nltk.download("wordnet")
nltk.download('omw-1.4')

lemmatiser = WordNetLemmatizer()
lemmatised = [lemmatiser.lemmatize(token) for token in stemmed_snow ]

#5
from nltk.corpus import stopwords
nltk.download("stopwords")

stop_words = set(stopwords.words("english"))

words_no_stop = [word for word in lemmatised if word not in stop_words]

tokens = len(words_no_stop)
word_tokens = words_no_stop

# DO NOT MODIFY THE BELOW LINE!


print('Number of word tokens: %d' % (tokens))
print("printing lists separated by commas")
print(*word_tokens, sep = ", ")

    Number of word tokens: 316
    printing lists separated by commas
```

```
l, , n, g, h, , , r, e, e, , , w, e, n, , , n, e, r, l, e, , g, n, ., , , , e, e, e,
, h, , , w, , p, n, g, , h, r, u, g, h, , h, e, , r, n, , g, e, , h, , l, e, , , h
, h, e, , r, v, e, , w, , j, u, , , n, r, r, w, , r, c, k, , n, w, ,, , , n, , u, r
, w, h, , g, r, , n, , w, e, e, ., , e, e, ,, , w, h, e, n, , , h, u, g, h, , , h,
, w, u, l, , p, p, e, r, , g, n, ,, , b, e, n, e, h, , , f, l, l, e, n, , r, e, e, , r
, f, r, e, , b, , h, e, , w, n, e, r, , r, n, ., , h, e, , r, e, e, , h, , h, r, w, n
, l, w, , b, r, n, c, h, e, , w, h, c, h, , r, e, c, h, e, , c, r, , , w, ., , , c, e
, u, e, n, l, ,, , n, , , h, e, r, e, , w, h, , , h, e, r, , b, e, n, g, , f, , n,
, f, l, l, n, g, , , e, e, .
[nltk_data] Error loading punk: Package 'punk' not found in index
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]    Package wordnet is already up-to-date!
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data]    Package omw-1.4 is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]    Package stopwords is already up-to-date!
```