

## R12725049 徐尚淵 作業二

1. 執行環境：Jupyter Notebook
2. 程式語言：Python (版本 3.11.4)

```
In [1]: from platform import python_version
        print(python_version())

3.11.4
```

3. 執行方式：

直接使用 Jupyter Notebook GUI Run code



4. 作業處理邏輯說明

1. 先做與作業一相同的前處理 (簡化截圖)

```
In [2]: import os
import math

# 資料夾路徑
folder_path = r'./data'

# 存儲文檔內容的列表
document_contents = []

# 列出資料夾中的所有文件
for filename in os.listdir(folder_path):
    # 文件路徑
    file_path = os.path.join(folder_path, filename)

    if os.path.isfile(file_path):
        with open(file_path, "r", encoding="utf-8") as file:
            document_contents.append(file.read())

In [3]: print(document_contents[0])

the white house is also keeping a close watch on yugoslavia, where
opposition forces are about to step up the pressure on president slobodan
milosevic. but will it work? nbc's jim maceda is in belgrade tonight.
```

```
In [7]: # Empty List for storing processed document
processed_texts = []

for document_content in document_contents:
    result = []
    # Lower casting
    document_content = document_content.lower()
    # Tokenized
    tokenized_content = tokenize_text(document_content)
    # Stopwords removal
    filtered_tokens = [token for token in tokenized_content if token not in stop]
    # Stemming
    for t in filtered_tokens:
        result.append(ps.stem(t))
    # 存入List
    processed_texts.append(result)

# 以文件一做測試
print(processed_texts[0])

['white', 'hous', 'also', 'keep', 'close', 'watch', 'yugoslavia', 'opposit', 'f',
orc', 'step', 'pressur', 'presid', 'slobodan', 'milosev', 'work', 'nbc', 'jim',
'maceda', 'belgrad', 'tonight', 'serbia', 'eve', 'gener', 'strike', 'two-hour',
'roadblock', 'tast', 'come', 'tomorrow', 'say', 'opposit', 'nationwid', 'work',
```



## 5. 將每個文件的 tf-idf 儲存

```
output_folder = "output"
os.makedirs(output_folder, exist_ok=True)

# 將每個文件的 tf-idf 儲存
for document_index, words in enumerate(processed_texts, start=1):
    # 文件名為 Docid.txt
    document_filename = f"{document_index}.txt"

    with open(os.path.join(output_folder, document_filename), "w", encoding="utf-8"):
        # Header
        file.write("t_index\ttf-idf\n")

        for tfidf_info in tfidf_vectors[document_index - 1]:
            line = "{:8}\t{}\n".format(tfidf_info['index'], tfidf_info['tf-idf'])
            file.write(line)
```

## 6. Cosine Similarity Function

```
In [13]: def cosine(vector_x, vector_y):
    # 計算向量內積
    dot_product = sum(x['tf-idf'] * y['tf-idf'] for x, y in zip(vector_x, vector_y))

    # 計算向量大小
    magnitude_x = math.sqrt(sum(x['tf-idf']**2 for x in vector_x))
    magnitude_y = math.sqrt(sum(y['tf-idf']**2 for y in vector_y))

    # Cosine similarity
    similarity = dot_product / (magnitude_x * magnitude_y)

    return similarity

def unit_vector(vector):
    # 計算向量大小
    magnitude = math.sqrt(sum(x['tf-idf']**2 for x in vector))

    # 計算單位向量
    unit_vector = [{'tf-idf': x['tf-idf'] / magnitude} for x in vector]

    return unit_vector
```

## 7. 先計算單位向量再計算 Cosine Similarity

```
In [14]: # 計算兩文件各自的單位向量
unit_vector_x = unit_vector(tfidf_vectors[154])
unit_vector_y = unit_vector(tfidf_vectors[958])

# 計算 cosine similarity
similarity = cosine(unit_vector_x, unit_vector_y)

print(f"文件 x 和文件 y 的 cosine similarity: {similarity}")

文件 x 和文件 y 的 cosine similarity: 0.6642675123400708
```