

▼ Lab#4, NLP@CGU Spring 2023

This is due on 2023/04/20 16:00, commit to your github as a PDF (lab4.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/xxxxxxx>

此內容會顯示為程式碼

Student ID:B0928006

Name: 廖采葳

▼ Word Embeddings for text classification

請訓練一個使用 Word2Vec 分類器來和 Google's Universal Sentence Encoder (a fixed-length 512-
儲存成功! 分類結果比較

```
!wget -O Dcard.db https://github.com/cjwu/cjwu.github.io/raw/master/courses/nlp2023/lab4-Dcard-Dataset.db

--2023-04-24 07:32:07-- https://github.com/cjwu/cjwu.github.io/raw/master/courses/nlp2023/lab4-Dcard-Dataset.db
Resolving github.com (github.com)... 140.82.112.4
Connecting to github.com (github.com)|140.82.112.4|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://raw.githubusercontent.com/cjwu/cjwu.github.io/master/courses/nlp2023/lab4-Dcard-Dataset.db [following]
--2023-04-24 07:32:07-- https://raw.githubusercontent.com/cjwu/cjwu.github.io/master/courses/nlp2023/lab4-Dcard-Dataset.db
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.110.133, 185.199.111.133, 185.199.109.133, ...
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.110.133|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 151552 (148K) [application/octet-stream]
Saving to: 'Dcard.db'

Dcard.db          100%[=====>] 148.00K  --.-KB/s    in 0.02s

2023-04-24 07:32:07 (6.67 MB/s) - 'Dcard.db' saved [151552/151552]

import sqlite3
import pandas as pd

conn = sqlite3.connect("Dcard.db")
df = pd.read_sql("SELECT * FROM Posts;", conn)
df
```

```
df.loc[I.flatten(), cols_to_show]
```

```
precision = 0.8
```

- ▼ Implement Your kNN or SVM classifier Here!

請比較分類結果中選出 topk 相近的筆數，並計算 forum_zh 是否都有在 query text 的 forum_zh 中

[開了新頻道] [Youtuber | 頻道 | 有趣 | 日常 | 搞笑]

330

YouTuber

```
precision = 0
topk = 10
```

```
# kNN
from sklearn.neighbors import KNeighborsClassifier
```

```
#train
X_train = embeddings
y_train = df["forum zh"]
```

```
knn = KNeighborsClassifier(n_neighbors=5)
knn.fit(X_train, y_train)
```

```
predicted_labels = knn.predict(embeddings[I.flatten()])
```

```
for i, label in enumerate(predicted_labels):
    if label == plabel:
```

```
precision += 1
```

```
h())["forum_zh"].values.tolist()
```

```
# if label in searched_labels:
#     print(f"{label} is in the searched labels.")
# else:
#     print(f"{label} is not in the searched labels.")
```

```
# # DO NOT MODIFY THE BELOW LINE!
print("precision = ", precision/topk)
```

儲存成功!