

▼ 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/1rqn1x6NgtvDSMQBAkDYGMWIHZXkWxtId?usp=share_link

Student ID: B0928026

Name: 洪詩晴

▼ Word Embeddings for text classification

請訓練一個 kNN或是SVM 分類器來和 Google's Universal Sentence Encoder (a fixed-length 512-dimension embedding) 的分類結果比較

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

--2023-04-24 15:03:11-- https://github.com/cjwu/cjwu.github.io/raw/master/courses/nlp2023/lab4-Dcard-Dataset.db
Resolving github.com (github.com)... 192.30.255.113
Connecting to github.com (github.com)|192.30.255.113|: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 15:03:11-- 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 15:03:11 (9.39 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
```

```

    createdAt      title      excerpt  categories
0      2022-03-04T07:54:19.886Z  專題需要數據🥺🥺 幫填~      希望各位能花個20秒幫我填一下
1      2022-03-04T07:42:59.512Z      #詢問 找衣服🥺      想找這套衣服🥺，但發現不知道該用什麼關鍵字找，（圖是草屯团仔的校園演唱會截圖）      詢問      衣服 | 鞋子 | 衣物 |
2      2022-03-      #黑特 網購50% FIFTY      因為文會有點長，先說結論是，50%是目前網購過的平台退      黑特 | 網購 | 三思 |

!pip3 install -q tensorflow_text
!pip3 install -q faiss-cpu

6.0/6.0 MB 63.7 MB/s eta 0:00:00
17.6/17.6 MB 19.7 MB/s eta 0:00:00

import tensorflow_hub as hub
import numpy as np
import tensorflow_text
import faiss

embed_model = hub.load("https://tfhub.dev/google/universal-sentence-encoder-multilingual/3")

docid = 355
texts = "[" + df['title'] + ']' + df['topics'] + ']' + df['excerpt']
texts[docid]

'[開了新頻道] [Youtuber | 頻道 | 有趣 | 日常 | 搞笑] 昨天上了第一支影片，之前有發過沒有線條的動畫影片，新的頻道改成了有線條的，感覺大家好像比
內容主要是分享自己遇到的小故事，不知道這樣的頻道大家是否會想要看呢？喜歡的話也'

embeddings = embed_model(texts)
embed_arrays = np.array(embeddings)
index_arrays = df.index.values
topk = 10
# Step 1: Change data type
embeddings = embed_arrays.astype("float32")

# Step 2: Instantiate the index using a type of distance, which is L2 here
index = faiss.IndexFlatL2(embeddings.shape[1])

# Step 3: Pass the index to IndexIDMap
index = faiss.IndexIDMap(index)

# Step 4: Add vectors and their IDs
index.add_with_ids(embeddings, index_arrays)

D, I = index.search(np.array([embeddings[docid]]), topk)

plabel = df.iloc[docid]['forum_zh']

cols_to_show = ['title', 'excerpt', 'forum_zh']
plist = df.loc[I.flatten(), cols_to_show]

precision = 0
for index, row in plist.iterrows():
    if plabel == row["forum_zh"]:
        precision += 1

print("precision = ", precision/topk)
precision = 0

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]
```

```
docid = 355
texts = "[" + df['title'] + ' ' + df['topics'] + ' '
texts[docid]
```

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

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

```
# YOUR CODE HERE!
# IMPLEMENTIG TRIE IN PYTHON
from sklearn.svm import SVC
```

```
X_train = embeddings
y_train = df["forum_zh"]
```

```
svm = SVC(kernel="linear")
svm.fit(X_train, y_train)
```

```
predicted_labels = svm.predict(embeddings[I.flatten()])
precision = 0
for i, label in enumerate(predicted_labels):
    if label == plabel:
        precision += 1
```

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

```
precision = 0.8
```

[Colab 付費產品](#) - [按這裡取消合約](#)

✓ 0 秒 完成時間：晚上11:09

