### ▼ 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/1izrKoiMML17tmWUpuPzVm7P\_PmVJnBLk?usp=sharing

#### Student ID: B0928003

Name:郭玉俊

df

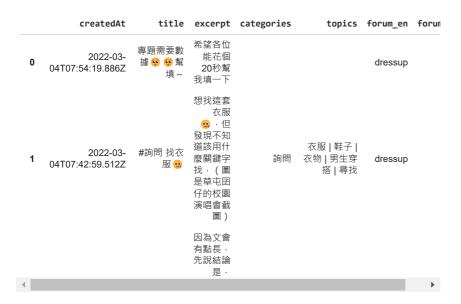
## Word Embeddings for text classification

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

### 按兩下 (或按 Enter 鍵) 即可編輯

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

```
--2023-04-24 06:53:21-- https://github.com/cjwu/cjwu.github.io/raw/master/courses/nlp2023/lab4-Dcard-Dataset.db
     Resolving github.com (github.com)... 20.27.177.113
     Connecting to github.com (github.com) | 20.27.177.113 | :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 06:53:21-- https://raw.githubusercontent.com/cjwu/cjwu.github.io/master/courses/nlp2023/lab4-Dcard-Dataset.db
     Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.108.133, 185.199.109.133, 185.199.111.133, ...
     Connecting to raw.githubusercontent.com (raw.githubusercontent.com) | 185.199.108.133 | :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 0.03s
     2023-04-24 06:53:22 (5.25 MB/s) - 'Dcard.db' saved [151552/151552]
import sqlite3
import pandas as pd
conn = sqlite3.connect("Dcard.db")
   = pd.read_sq1("SELECT * FROM Posts;", conn)
```



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

```
import tensorflow hub as hub
import numpy as np
import tensorflow_text
embed_model = hub.load("https://tfhub.dev/google/universal-sentence-encoder-multilingual/3")
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,新的頻道改成有線條的,感覺大家好像比較喜歡這種風格,試試看新的風格,影片
    內容主要是分享自己遇到的小故事,不知道清楼的艏道大家是否會想要看呢? 查數的話也!
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
```

forum_zh	excerpt	title	
YouTuber	昨天上了第一支影片·之前有發過沒有線條的動畫影片· 新的頻道改成有線條的·感覺大家好像比較喜歡	開了新頻道	355
YouTuber	哈哈哈哈·沒錯我就是親友團來介紹一個我覺得很北七的 頻道·現在觀看真的低的可憐·也沒事啦·就多	一個隨性系YouTube 頻道	359
YouTuber	又來跟大家分享新的作品了~·頻道常常分享 {縫紉} {服裝 製作} 等相關教學·大家對服裝製	《庫洛魔法使》(迷 你)服裝製作	330
YouTuber	勾惡幫主在自己頻道簡介跟每部影片的下方都已經說明了·要分會會長以上才能看全部影片·這個說明已	自己沒搞清楚狀況就 不要亂黑勾惡	342
YouTuber	友人傳了這篇文給我·我一看·十大廚師系YouTuber·就 猜一定有MASA·果不其然·榜上有	廚師系YouTuber	338
有趣	小時候都很喜歡看真珠美人魚和守護甜心·但是!!·每 次晚餐看電視的時候·只要有播映到這種場景	毁我童年的家人	243
YouTuber		喜歡看寵物頻道的有	349

# Implemement 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]

```
precision = 0
topk = 10
# YOUR CODE HERE!
# IMPLEMENTIG TRIE IN PYTHON
# 載入 Universal Sentence Encoder Multilingual 模型
embed_model = hub.load("https://tfhub.dev/google/universal-sentence-encoder-multilingual/3")
# 讀取資料並處理成所需格式
texts = "[" + df['title'] + '] [' + df['topics'] + '] ' + df['excerpt'] embeddings = embed_model(texts)
embed arrays = np.array(embeddings)
# 設置 k 值和查詢文本的 ID
k = 10
docid = 355
# 創建 kNN 模型
index = faiss.IndexFlatL2(embeddings.shape[1])
index.add(embed arrays)
# 進行搜索
D, I = index.search(np.array([embed_arrays[docid]]), k)
# 選出 topk 筆相似的資料
topk_indices = I.flatten()
topk_labels = df.loc[topk_indices]['forum_zh']
# 計算精確度
plabel = df.iloc[docid]['forum_zh']
\label{eq:precision} \mbox{ = (topk\_labels == plabel).sum() / k}
# # DO NOT MODIFY THE BELOW LINE!
print("precision = ", precision)
     precision = 0.8
```

Colah 付费產品 - 按這裡取消合約

✓ 9秒 完成時間: 下午3:35