

▼ 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

####

Student ID:B0928002

Name:林力

▼ 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
```

```
--2023-04-24 07:59:07-- https://github.com/cjwu/cjwu.github.io/raw/master/cov
Resolving github.com (github.com)... 140.82.121.3
Connecting to github.com (github.com)|140.82.121.3|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://raw.githubusercontent.com/cjwu/cjwu.github.io/master/courses
--2023-04-24 07:59:08-- https://raw.githubusercontent.com/cjwu/cjwu.github.io
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.108
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.108
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 0.003s
```

```
2023-04-24 07:59:08 (50.4 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	topics	forum_en	fc
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			dressup	
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，（圖是草屯囡仔的校園演唱會截圖） 因為文會有點長，先說結論是，50%	詢問	衣服 鞋子 衣物 男生穿搭 尋找	dressup	

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

6.0/6.0 MB 43.0 MB/s eta 0:00:00

17.0/17.0 MB 16.6 MB/s eta 0:00:

```
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-multili

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，不知道這樣的頻道大家早不會想要看呢？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
```

	title	excerpt	forum_zh
355	開了新頻道	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...	YouTuber
359	一個隨性系 YouTube 頻道	哈哈哈哈哈，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...	YouTuber
330	《庫洛魔法使》(迷你) 服裝製作	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...	YouTuber
342	自己沒搞清楚狀況就不要亂黑勾惡	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...	YouTuber
338	廚師系 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...	YouTuber
243	毀我童年的家人	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...	有趣
349	喜歡看寵物頻道的有嗎？🐶🐱		YouTuber

▼ 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]

```
precision = 0
topk = 10
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split
import pandas as pd
import numpy as np
# YOUR CODE HERE!
# IMPLEMENTIG TRIE IN PYTHON
# Define the query text
X = df['title'] + ' ' + df['topics'] + ' ' + df['excerpt']
y = df['forum_zh']
vectorizer = TfidfVectorizer()
X = vectorizer.fit_transform(X)

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_sta

# Train the classifier
n_neighbors = 5
clf = KNeighborsClassifier(n_neighbors=n_neighbors)
clf.fit(X_train, y_train)

# Predict the labels of the testing set
y_pred = clf.predict(X_test)

# Compute the precision of the top k retrieved neighbors
k = 10
precisions = []
for i, y_true in enumerate(y_test):
    neighbors = clf.kneighbors(X_test[i], n_neighbors=k, return_distance=False)[0]
    if y_true in y[neighbors]:
        precisions.append(1)
    else:
        precisions.append(0)
precision = np.mean(precisions)
# # DO NOT MODIFY THE BELOW LINE!
print("precision = ", precision/topk)
```

precision = 0.0

[Colab paid products](#) - [Cancel contracts here](#)

✓ 0s completed at 4:02 PM

● ×