

NCKU CSIE 112 Senior Project

# 以文會友

Find Your Best Study Partner on Facebook

分析個人在 Facebook 發佈的貼文，來媒合有共同興趣的朋友，並推薦其適合的揪課課程

Team 4-2

指導教授

專題生

李強

連思涵 F74099017



# Outline

00 Introduction

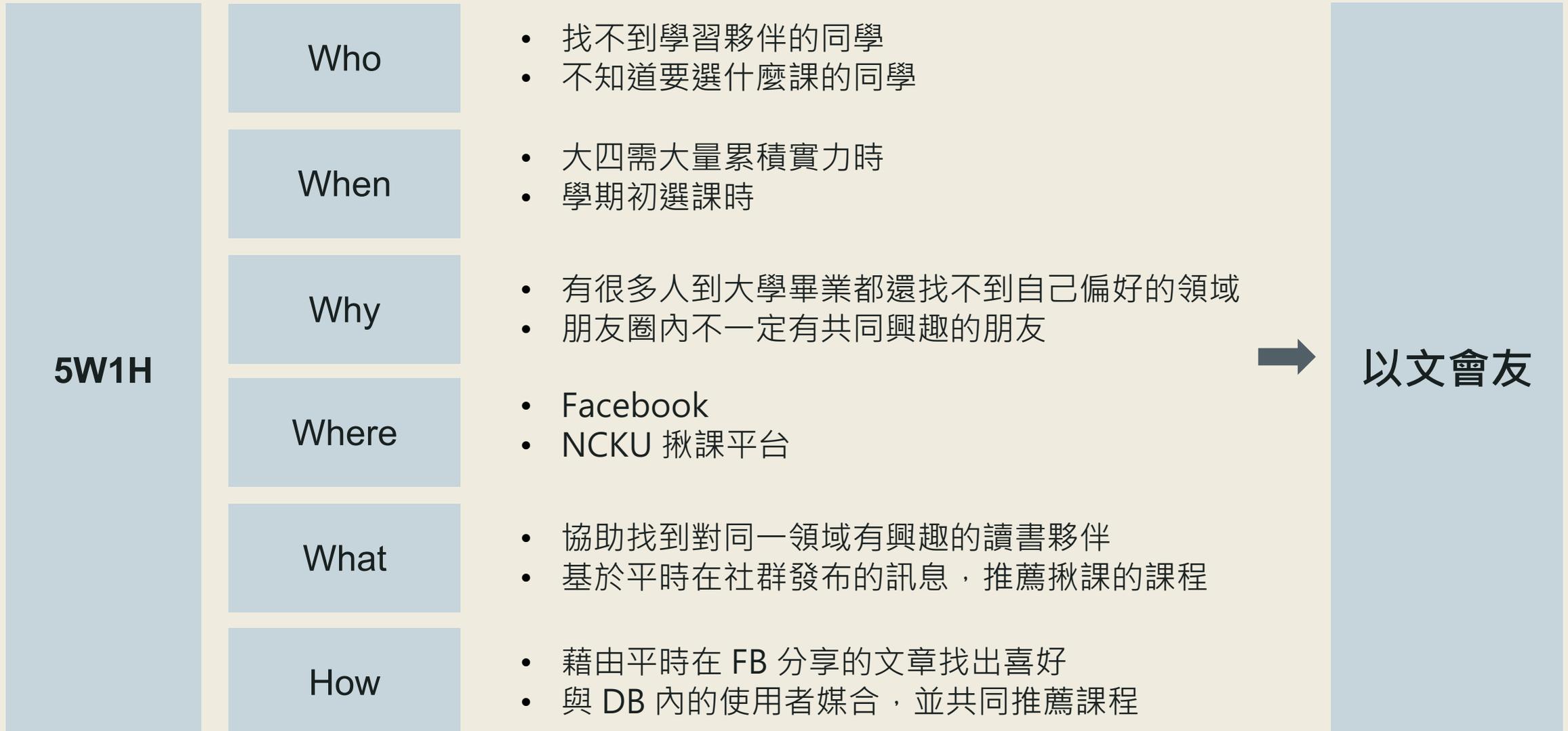
01 Text Compare and Recommendation

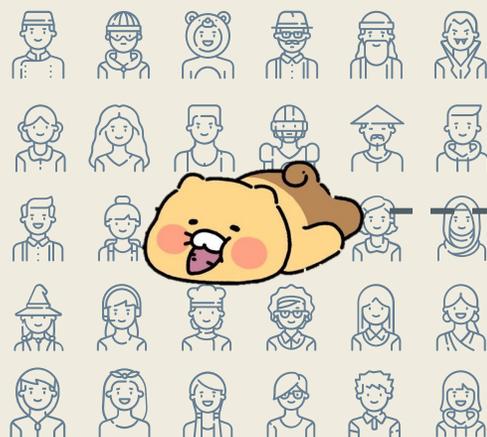
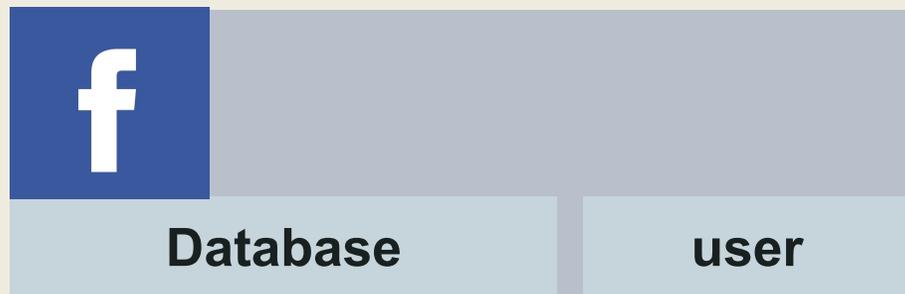
02 Post Analysis and Recommendation

03 Experience and Review

# 00 Introduction

## 研究動機的 5W1H





# 以文會友

Text Compare

Post Analysis

Recommendation with Collaborative Filtering

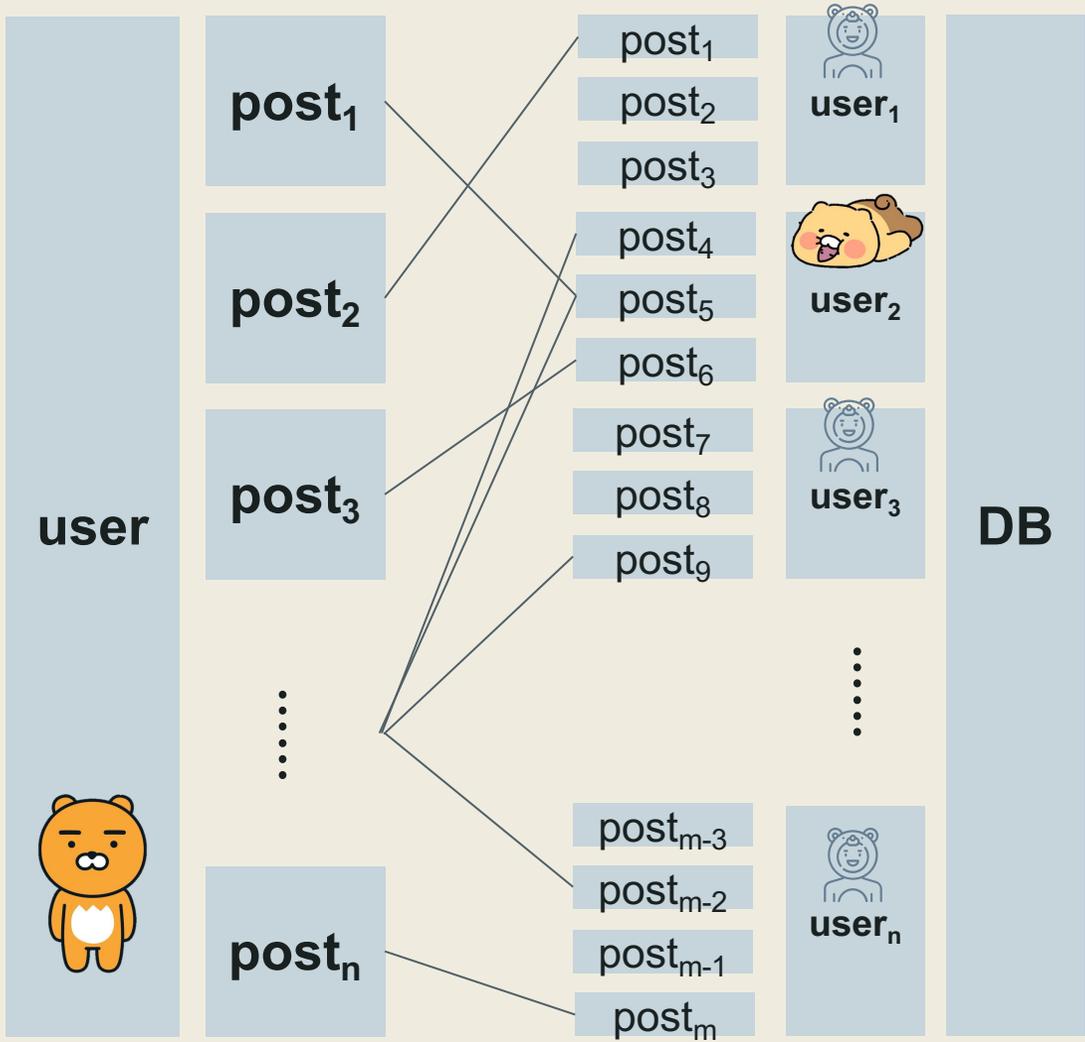
Content-based Recommendation



# **01 Text Compare and Recommendation**

# 01 Text Compare and Recommendation

## 將使用者的 n 篇貼文與 DB 內的資料做 TF-IDF 文本比對，找出最佳讀書夥伴



	PG 財經筆記=0	PG 財經筆記=1	PG 財經筆記=2	PG 財經筆記=3	PG 財經筆記=4	PG 財經筆記=5	PG 財經筆記=6
砂谷阿雅 Anya Cheng=0	0.0	0.21061544120311700	0.0	0.015361979603767400	0.07673442363739910	0.0011455945204943400	0.07391513884067540
砂谷阿雅 Anya Cheng=1	0.0	0.08701366931200030	0.04551013559103010	0.0	0.0	0.0642264363634758	0.0
砂谷阿雅 Anya Cheng=2	0.01428665965795520	0.07027796655893330	0.036756981164217000	0.0	0.04267445579171180	0.051873501390218700	0.009457391686737540
砂谷阿雅 Anya Cheng=3	0.05020962655544280	0.0	0.0	0.02014680579304700	0.0	0.01452986802905800	0.04460490867495540
砂谷阿雅 Anya Cheng=4	0.0	0.0	0.0	0.04368483158607480	0.0	0.0	0.0
砂谷阿雅 Anya Cheng=5	0.01930840127170090	0.0	0.0	0.0659237653017044	0.0	0.0	0.0
砂谷阿雅 Anya Cheng=6	0.008547567762434480	0.04473619535565380	0.001196285360492770	0.02158987894654270	0.05155201256275180	0.08111820369959880	0.1028934121131900
砂谷阿雅 Anya Cheng=7	0.034765198826789900	0.0	0.0	0.006786812096834180	0.0	0.030891088768839800	0.09519081562757490
砂谷阿雅 Anya Cheng=8	0.017523692920804000	0.13220977783203100	0.009789655764102940	0.00615262845531106	0.03211237117648180	0.011916029267013100	0.06276460736989980
砂谷阿雅 Anya Cheng=9	0.0	0.10345801711082500	0.03482940047979360	0.02193448506295680	0.0	0.024576593190431600	0.008961433544754980
砂谷阿雅 Anya Cheng=10	0.0	0.12083148956298800	0.0	0.0	0.1467435210943220	0.0	0.03252088278532030
砂谷阿雅 Anya Cheng=11	0.0	0.0	0.0	0.0	0.0	0.0	0.0
砂谷阿雅 Anya Cheng=12	0.040348924696445500	0.0	0.0	0.018096202984452200	0.060261406004428900	0.013354961760342100	0.0
砂谷阿雅 Anya Cheng=13	0.018389329314231900	0.0014181813457980800	0.0007417412125505510	0.008587183430790900	0.03196415677666680	0.04125290364027020	0.035349804908037200
砂谷阿雅 Anya Cheng=14	0.0476633757352829	0.05116431042551990	0.011860466562211500	0.0272291749717580500	0.05348987504839900	0.03646244559879303	0.034092675894498800

	半路出家軟體工程師在砂谷=1	半路出家軟體工程師在砂谷=2	半路出家軟體工程師在砂谷=3	半路出家軟體工程師在砂谷=4	半路出家軟體工程師在砂谷=5	半路出家軟體工程師在砂谷=6
砂谷阿雅 Anya Cheng=0	0.03230162337422370	0.09413079917430880	0.16493353247642500	0.09936847537755970	0.0706122219625310	0.0742359310388565
砂谷阿雅 Anya Cheng=1	0.020086079835891700	0.05602323263883590	0.06831202656030660	0.03114149160683160	0.08724986880414960	0.06453309208154680
砂谷阿雅 Anya Cheng=2	0.020322639495134400	0.08869554102420810	0.06959442049264910	0.10377360880374900	0.0735383098520810	0.07031203806400300
砂谷阿雅 Anya Cheng=3	0.050991594791412400	0.00588524155318737	0.10716690123081200	0.006763797253370290	0.030549844726920100	0.02786768227815630
砂谷阿雅 Anya Cheng=4	0.007944576442241670	0.019083848764896400	0.115008294864896400	0.04882487654685970	0.0	0.0135298045352100
砂谷阿雅 Anya Cheng=5	0.004155758302658800	0.02840055523028990	0.07302099466323850	0.05544073134680720	0.015372841618955100	0.010579117573797700
砂谷阿雅 Anya Cheng=6	0.0484774003932000	0.09301892668008800	0.17500029504299200	0.15268391370773300	0.0387902557849884	0.06955569492412570
砂谷阿雅 Anya Cheng=7	0.06076361984014510	0.09727168828248980	0.15979471802711500	0.0406737373576160	0.0035688746568262600	0.07203321903944020
砂谷阿雅 Anya Cheng=8	0.020340092490182600	0.10245874524116500	0.2702205771768270	0.05756952613592150	0.16627046465873700	0.060548584908247000
砂谷阿雅 Anya Cheng=9	0.020630884915590300	0.04924929141998290	0.13180261850357100	0.09415735304355620	0.14674808084964800	0.057659655809402500
砂谷阿雅 Anya Cheng=10	0.0	0.023886656388640400	0.1284452577614600	0.08582510054111840	0.0	0.016902871429920200
砂谷阿雅 Anya Cheng=11	0.0	0.0	0.020670926198363300	0.0	0.13742424547672300	0.0
砂谷阿雅 Anya Cheng=12	0.00883545447140932	0.047613758593797700	0.08018850535154340	0.09152454882860180	0.1786891222000120	0.05572190508246420
砂谷阿雅 Anya Cheng=13	0.026880767196416900	0.08219444006681440	0.2609350383281710	0.10437951982021300	0.06579418480396270	0.052596960216760600
砂谷阿雅 Anya Cheng=14	0.03547798469662670	0.10454515367746400	0.14531861245632200	0.14667677879333500	0.06827839463949200	0.06818445771932600

	Min 的投資說書小檔=0	Min 的投資說書小檔=1	Min 的投資說書小檔=2	Min 的投資說書小檔=3	Min 的投資說書小檔=4	Min 的投資說書小檔=5	Min 的投資說書小檔=6
砂谷阿雅 Anya Cheng=0	0.0	0.03144350647926330	0.038628607988357500	0.021221471950411800	0.0679646059870200	0.012101402506232300	0.1079131896363030
砂谷阿雅 Anya Cheng=1	0.0	0.00521071208640933	0.02520398423075680	0.026119481772184400	0.0010995753109455100	0.07187148747825620	0.020624930039048200
砂谷阿雅 Anya Cheng=2	0.0	0.030006371438503300	0.02385772578418260	0.025626037269830700	0.004727283958345650	0.09201629459885794	0.016241729259491000
砂谷阿雅 Anya Cheng=3	0.0	0.0	0.00012107313377782700	0.00029072435572743400	0.028933946043252900	0.11216455698013300	0.0281015015757431980
砂谷阿雅 Anya Cheng=4	0.0	0.18392598628997800	0.019330322742462200	0.009730505757033830	0.010304443538188900	0.0	0.04341162368655210
砂谷阿雅 Anya Cheng=5	0.0	0.1422995628926420	0.01495545357465740	0.012663706205785300	0.023916976526379600	0.027891401201486600	0.03358663618564610
砂谷阿雅 Anya Cheng=6	0.02829722873866560	0.04344220831990240	0.044765204191207900	0.026664482429623600	0.025461072102189100	0.12088080495595900	0.04573161154985430
砂谷阿雅 Anya Cheng=7	0.0	0.06269685924053190	0.06692685931921010	0.028788655996322600	0.0391208790242672	0.13510340452194200	0.03511188551783560
砂谷阿雅 Anya Cheng=8	0.027852758765220600	0.028659239411354100	0.050835173577070200	0.0437321811914444	0.02817085012793540	0.058897923678159700	0.06546416878700260
砂谷阿雅 Anya Cheng=9	0.0	0.036091312766075100	0.012270506471395500	0.0387757308781147	0.005685921758413320	0.026327423751354200	0.03381908684698500
砂谷阿雅 Anya Cheng=10	0.0	0.008460070066549360	0.03932003304362300	0.03644002601504330	0.001785261556506160	0.11361657828092600	0.030484752202034
砂谷阿雅 Anya Cheng=11	0.0	0.0	0.020924221724271800	0.01297870185226200	0.05704332888126370	0.0	0.01308429326417400
砂谷阿雅 Anya Cheng=12	0.0	0.08867335319519040	0.000944624887779350	0.03508240729570390	0.0257823895663023	0.011516711674630600	0.04939771290779110
砂谷阿雅 Anya Cheng=13	0.008939091116189960	0.00584833137691021	0.0365571264564991	0.05703437700867650	0.051629888590765	0.0559527613222599	0.0607214495396650
砂谷阿雅 Anya Cheng=14	0.025669759139418600	0.04104606434702870	0.01988590766671000	0.05829806625843050	0.03972093388438230	0.08640435338020330	0.05514468625187870

## 將使用者的 $n$ 篇貼文與 DB 內的資料做 TF-IDF 文本比對，找出最佳讀書夥伴

user

post<sub>n</sub>



### 1. 資料清洗

- 去除標點符號
- 去除停用詞
- 去除數字
- 以 **jieba** 斷字

我是一個不知道要做什麼工作的人！

我/是/一個/不/知道/要/做/什麼/工作/的/人

「一個」、「知道」、「做」、「工作」

### 2. 建立詞袋

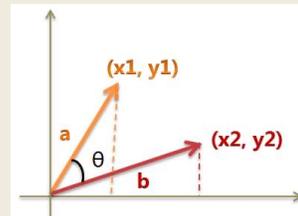
( **bag of words** )

- 對所有詞編號  
ex : 「你好」, 4
- 將貼文都轉成二元組的向量  
ex : 「你好」出現兩次  
→ ( 4, 2 )



### 3. 取得兩兩貼文相似度

- 得出詞袋中每個詞的 **TF-IDF** 值
- 分析 post<sub>n</sub> 和 DB 裡每篇貼文的餘弦相似度  
( **cosine similarity** )



某字詞在某貼文中出現的頻率

ex : 十萬 / 青年 / 十萬 / 肝  
→  $TF = 2/3 = 0.66666$

TF

X

文件數除以某特定字詞有被多少貼文所提及的數量取 log

ex : DB 總共有  $m$  篇貼文  
其中有 5 個文件提及「十萬」  
→  $IDF = \log(10/5) = 0.3$

IDF

post<sub>1</sub>

post<sub>2</sub>

post<sub>3</sub>

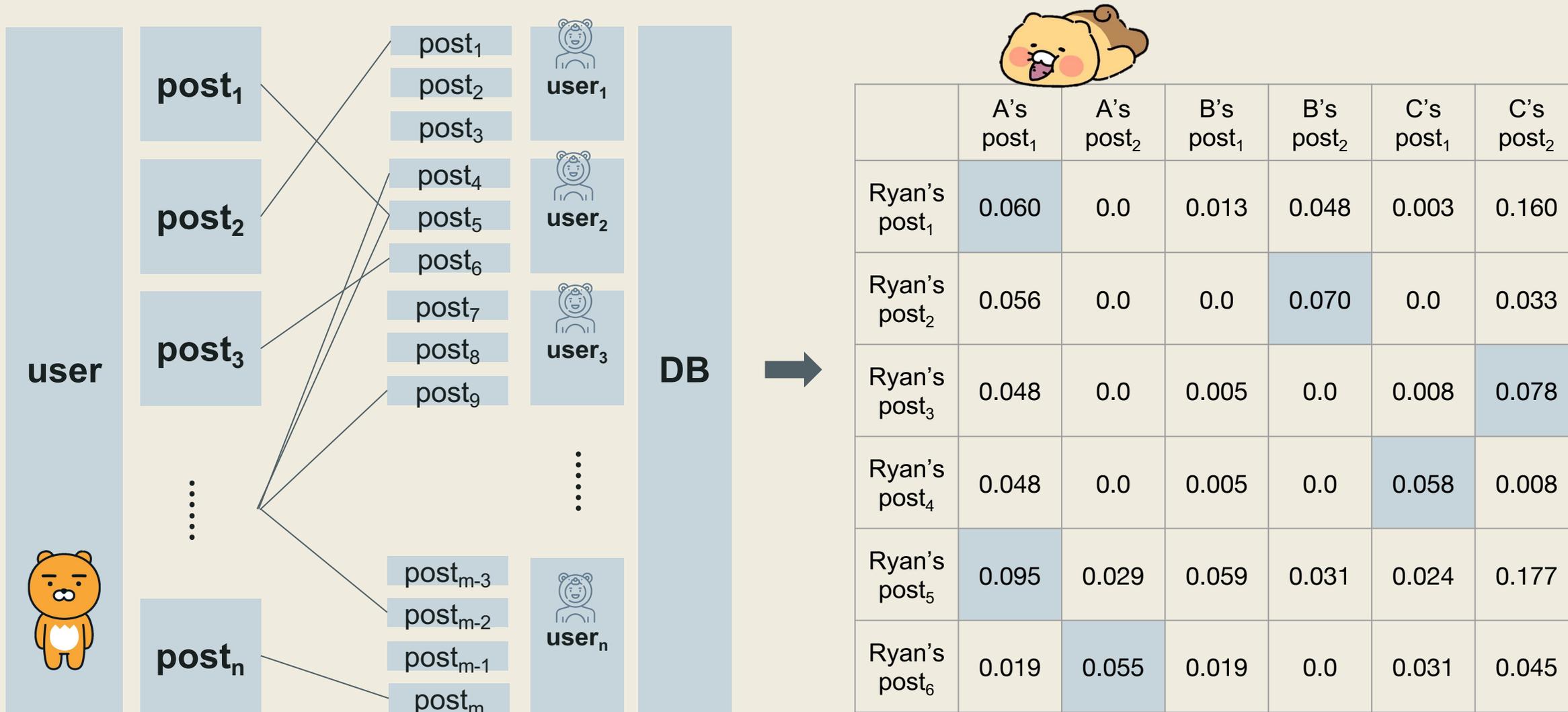
⋮

post<sub>m</sub>

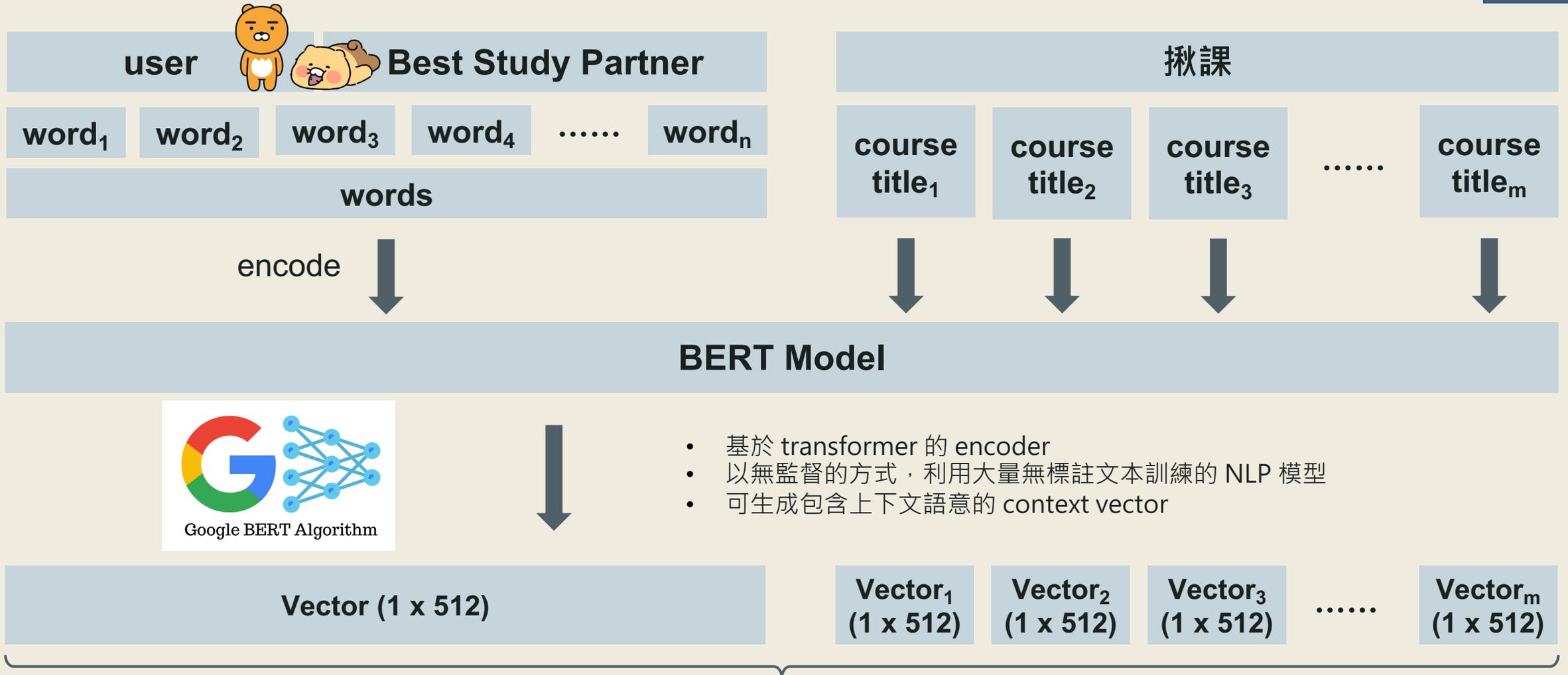
Someone from DB



將使用者的  $n$  篇貼文與 DB 內的資料做 TF-IDF 文本比對，找出最佳讀書夥伴



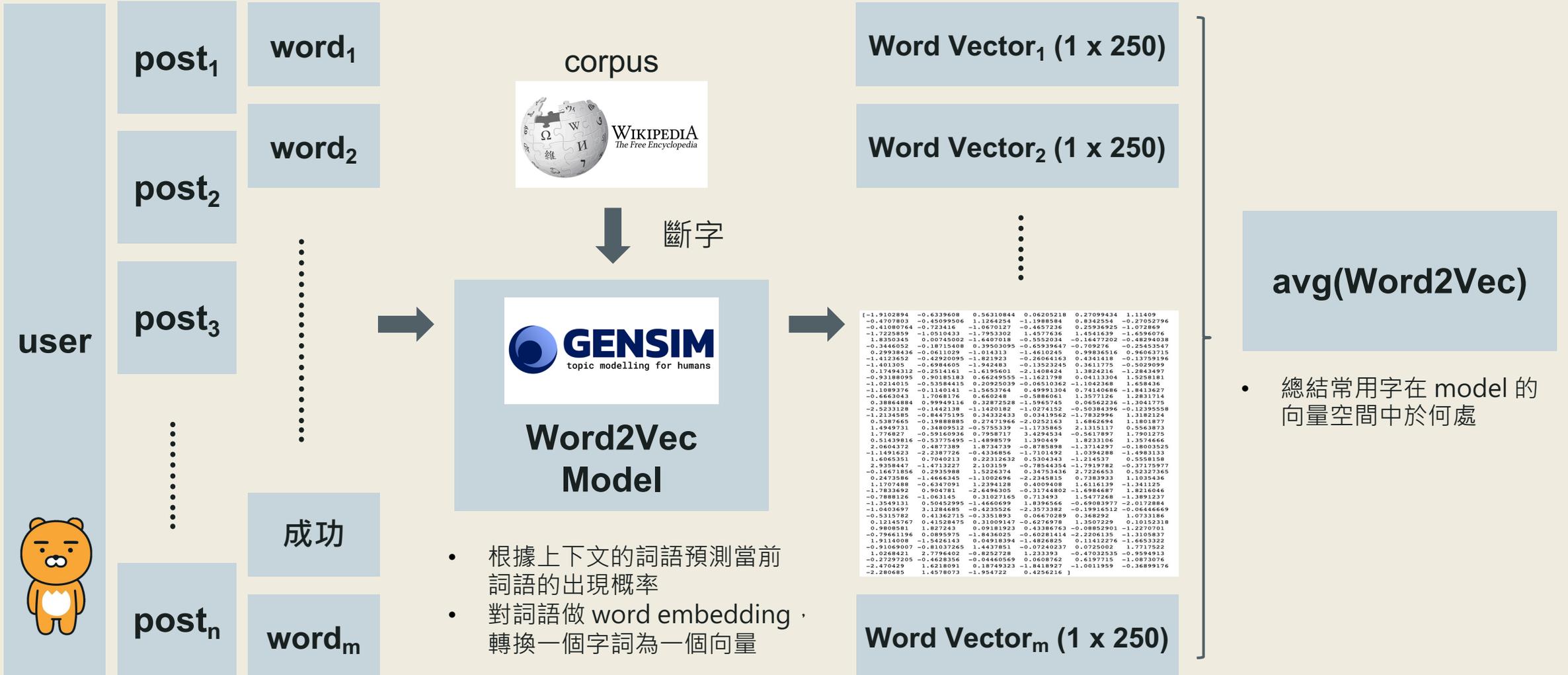
# 將使用者與讀書夥伴 TF-IDF 值都很高的字、揪課標題透過 BERT model 轉詞向量



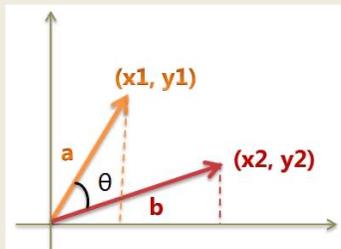
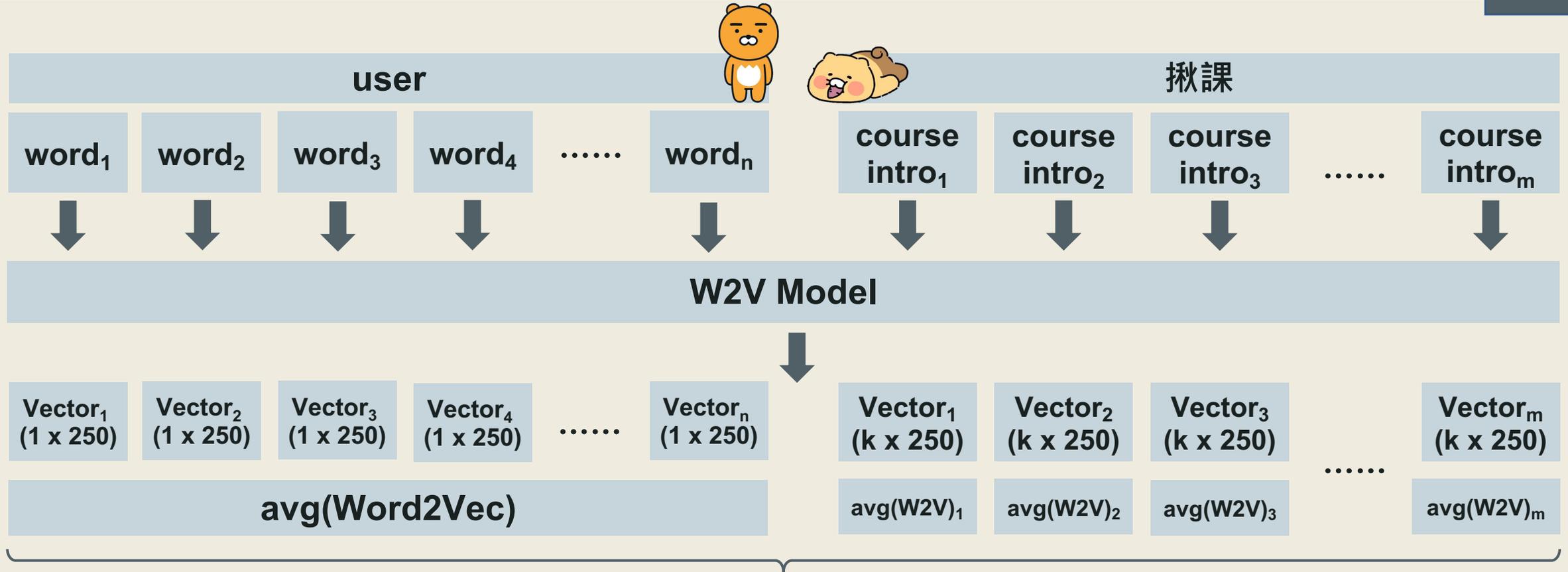
cosine similarity → 找出最相近的向量即為最推薦的課程

## **02 Post Analysis and Recommendation**

## 找出使用者所有貼文中 TF-IDF 高的字輸入進 w2v model 後，取出所有詞向量的平均



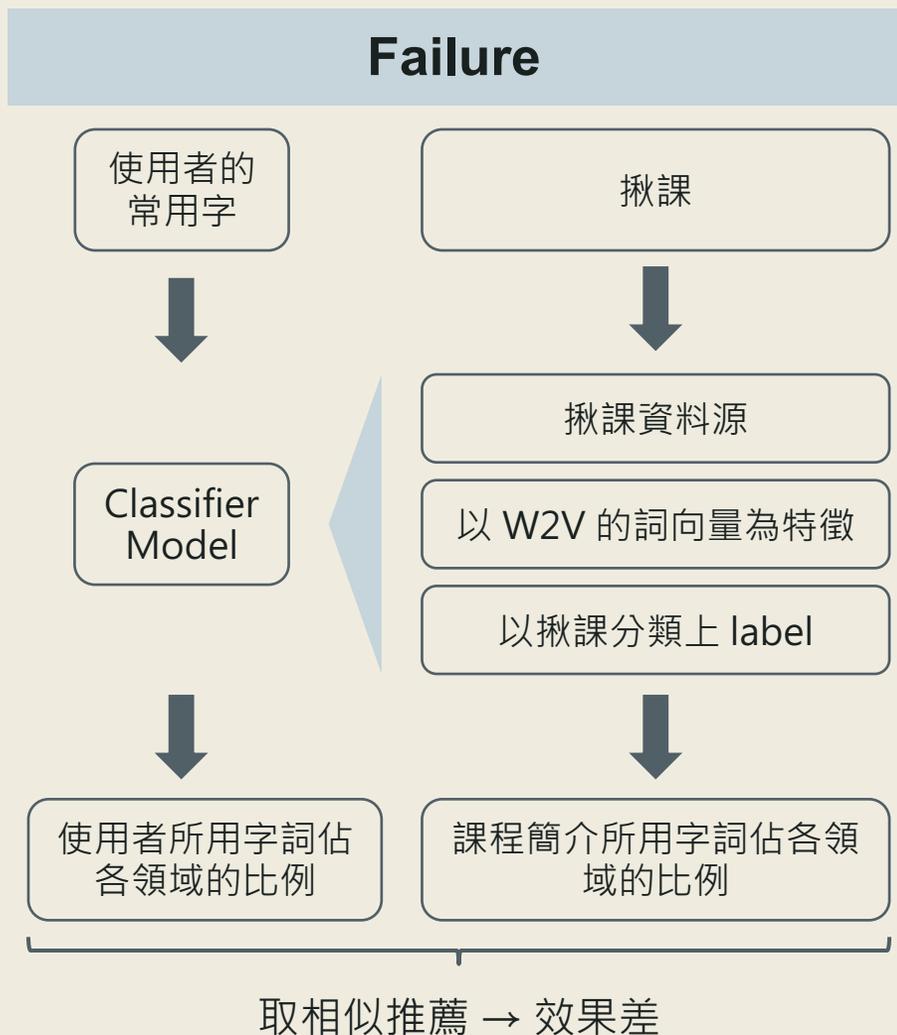
同理到揪課上所有課程的簡介後，比對相似度以推薦課程



cosine similarity → 找出最相近的向量即為最推薦的課程

## **03 Experience and Review**

## 實作失敗的經歷 & 本次專題需克服之處



### Flaw

#### W2V model

- 在處理 corpus 時，由於斷字無法處理複合字詞，會造成 model 的 key error (ex: 成功大學 → 「成功」、「大學」)
- 以維基百科的資料為 corpus，使得 model 預測的結果與想像有落差，若爬取社群媒體的文章可能可以改善 (ex: 「獅子」最相關為「兩隻」)

```
print(model.wv.most_similar('成功大學', topn=10), '\n')
```

```
KeyError                                Traceback (most recent call last)  
<ipython-input-135-fde0a1741ddd> in <module>  
----> 1 print(model.wv.most_similar('成功大學', topn=10), '\n')
```

```
print(model.wv.most_similar('獅子', topn=20), '\n')
```

```
[('兩隻', 0.5327986478805542), ('獨角獸', 0.519851565361023), ('公雞', 0.5154509544372559), ('烏鴉', 0.5109748840332031), ('鷺', 0.5069186687469482), ('鬃毛', 0.4994606673717499), ('獅鬚', 0.4973223805427551), ('牙吠', 0.49647188186645)
```

### Reality

#### Facebook

- 大部分同學的貼文大多與生活日常有關，較少分享學術知識、想法觀點
- 分享的內容因網頁結構複雜，尚未包含在分析的貼文裡

#### 揪課

- 資料數不夠多，分類器 precision 不到 40%
- 平台上的多分類造成字詞語意的分類不精確

**Thank you!**