

文字分析於校務研究的應用

淡江大學資訊處 曹乃龍

為何要作文字分析？

不要只依靠冰冷的數字

校園之外的資訊也很重要

LEVELS OF ANALYTICS

Analytics

文字雲

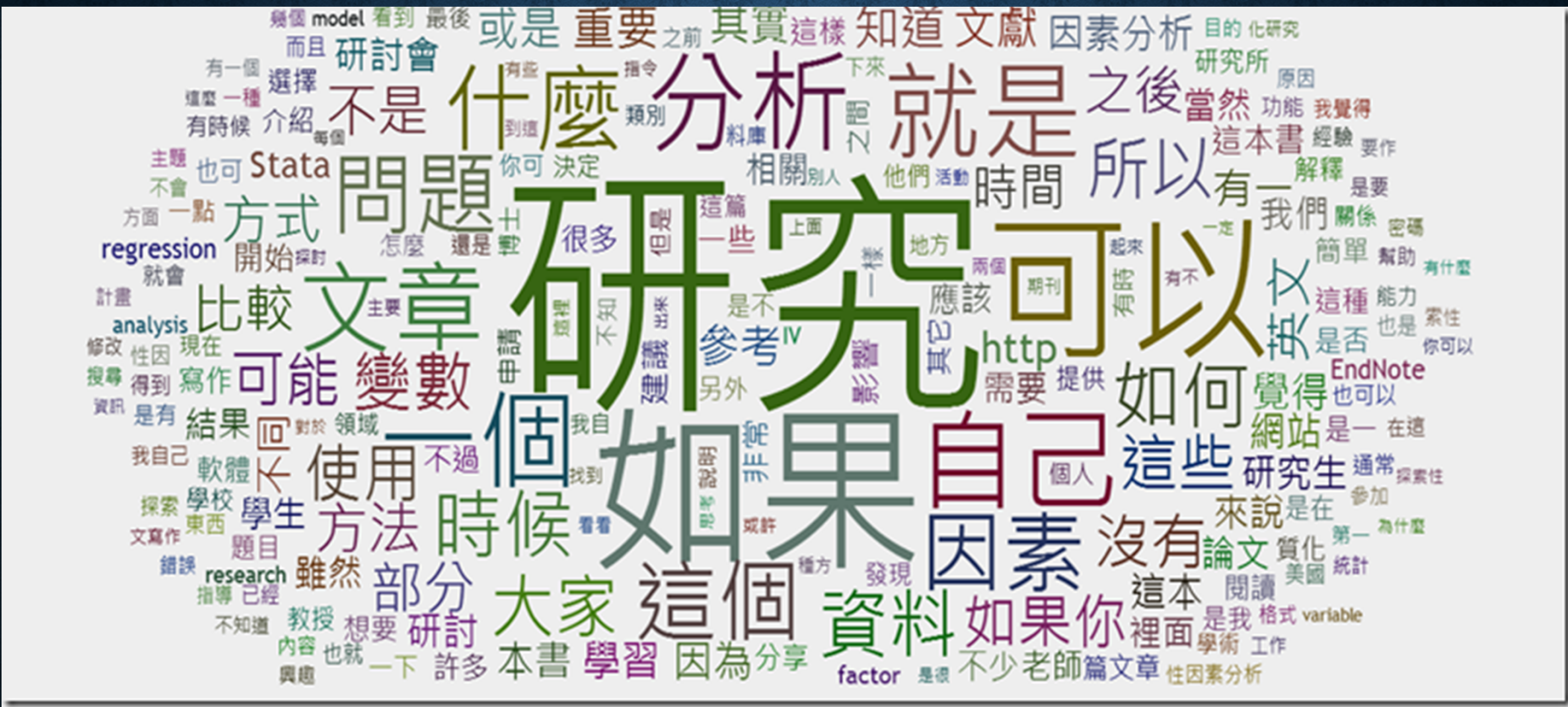
Handy Analytics

討論區評分

Actionable Analytics

課程設計輔助

ANALYTICS



HANDY ANALYTICS

- 定義：進行後可即時實際使用的分析結果
- 範例：討論區評分
 - 合作夥伴：滁州学院计算机与信息工程学院 程曦資訊
 - 資料來源：滁州学院計算機概論MOOC平台討論內容，共有4,868主題及42,135回應。
 - 外部資料：Google Search

使用討論區評分 (Before)

- 目的：利用課程討論區了解學生參與程度並給予評價
- 方法：利用發表主題及回應的次數及字數作為評分標準
- 問題：
 - 主題：哪个杀毒软件好？
 - 回應一： 360卫士
 - 回應二： 我不知道！

使用討論區評分 (After)

- 目的：分析回應內容判斷是否符合主題，以利給與正確評價
- 方法：利用Google Search 對回應進行內容擴充
- 問題：
 - 主題：哪个杀毒软件好？
 - 回應一： 360卫士，擴充後包含杀毒软件
 - 回應二： 我不知道！，擴充後不知所云

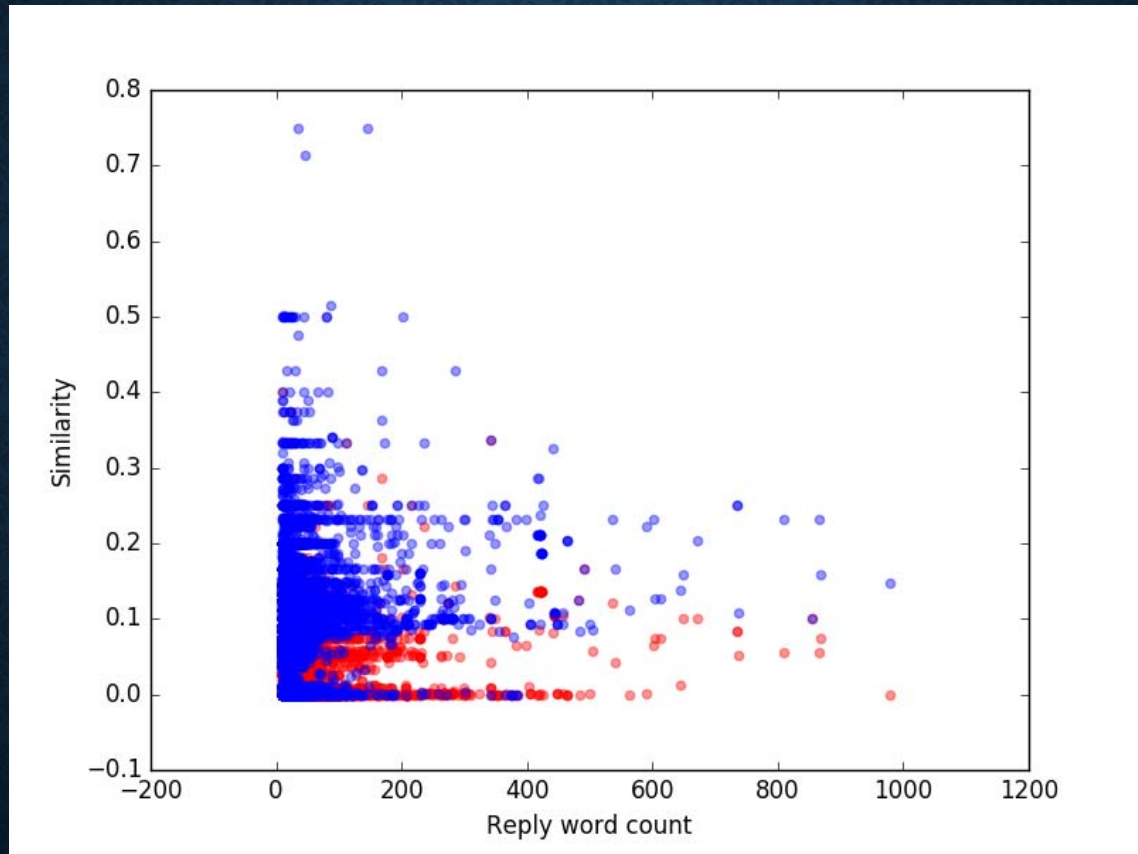


Text Similarity

Google Search

Text Analytics Engine

School Insight



Red: Before expansion, Blue: After expansion

SCHOOL INSIGHT INTEGRATION

- Data import
- Create dashboard menu item
- Pick attributes for visualization

DATA IMPORT

cx_topic_reply_sim	cx_topic_reply_sim	14	1	674319	1001281623	357820
cx_topic_reply_sim	cx_topic_reply_sim	19	1	327710	1001280939	816933
cx_topic_reply_sim	cx_topic_reply_sim	22	1	660827	1001280589	157471
cx_topic_reply_sim	cx_topic_reply_sim	24	1	846383	1001280665	357598
cx_topic_reply_sim	cx_topic_reply_sim	26	1	327728	1001280752	157471
cx_topic_reply_sim	cx_topic_reply_sim	26	1	327728	1001280978	816933
cx_topic_reply_sim	cx_topic_reply_sim	19	1	327710	1001280667	275322
cx_topic_reply_sim	cx_topic_reply_sim	22	1	660827	1001280611	816855
cx_topic_reply_sim	cx_topic_reply_sim	24	1	846383	1001509469	816933
cx_topic_reply_sim	cx_topic_reply_sim	26	1	327728	1001281036	443025
cx_topic_reply_sim	cx_topic_reply_sim	19	1	327710	1001239606	184964
cx_topic_reply_sim	cx_topic_reply_sim	22	1	660827	1001282160	792160
cx_topic_reply_sim	cx_topic_reply_sim	24	1	846383	1001280978	816834
cx_topic_reply_sim	cx_topic_reply_sim	26	1	327728	1001280667	364590
cx_topic_reply_sim	cx_topic_reply_sim	19	1	327710	1001280695	392228
cx_topic_reply_sim	cx_topic_reply_sim	22	1	660827	1001280970	151123
cx_topic_reply_sim	cx_topic_reply_sim	24	1	846383	1001280978	816876
cx_topic_reply_sim	cx_topic_reply_sim	26	1	327728	1001283636	157516
cx_topic_reply_sim	cx_topic_reply_sim	19	1	327710	1001240271	357667
cx_topic_reply_sim	cx_topic_reply_sim	22	1	660827	1001281692	816855
cx_topic_reply_sim	cx_topic_reply_sim	24	1	846383	1001282871	357598
cx_topic_reply_sim	cx_topic_reply_sim	26	1	327728	1001280939	143206
cx_topic_reply_sim	cx_topic_reply_sim	19	1	327710	1001281902	185162
cx_topic_reply_sim	cx_topic_reply_sim	22	1	660827	1001280970	148901
cx_topic_reply_sim	cx_topic_reply_sim	24	1	846383	1001281091	357751
cx_topic_reply_sim	cx_topic_reply_sim	26	1	327728	1001281288	157479
cx_topic_reply_sim	cx_topic_reply_sim	19	1	327710	1001280978	357598
cx_topic_reply_sim	cx_topic_reply_sim	22	1	660827	1001280390	377781
cx_topic_reply_sim	cx_topic_reply_sim	24	1	846383	1001239915	816906
cx_topic_reply_sim	cx_topic_reply_sim	26	1	327728		

原始数据

```

{
  "index": "cx_topic_reply_sim",
  "type": "cx_topic_reply_sim",
  "id": "14",
  "version": 1,
  "score": 1,
  "source": {
    "nst_sim": 1,
    "reply_id": 674319,
    "reply_user_id": 1001281623,
    "topic_id": 357820,
    "sim_char_TR": 0.000018,
    "sim_char_TE": 0.071516,
    "sim_word_TR": 0.000018,
    "sim_word_TE": 0.125031,
    "reply_word_count": 12
  }
}

```

CREATE MENU ITEM

Create [Home](#) > [次選單](#) > Create

新增

主選單	回應與主題相似度
次選單名稱	回應與主題相似度
啟用	<input checked="" type="checkbox"/> ON
排序	
URL	c2481002
Controller	c2481002
Index	Index
參數	
圖示	fa-bar-chart-o fa-fw

PICK THE ATTRIBUTES FOR VISUALIZATION

欄位設定

TYPE: cx topic replv sim

報表欄位

報表樣式

- sim_char_TE
- sim_word_TR
- topic_id
- reply_id
- sim_char_TR
- reply_user_id

柱狀圖

細項欄位

sim char TF

欄位名稱: 欄位名稱

已選擇欄位

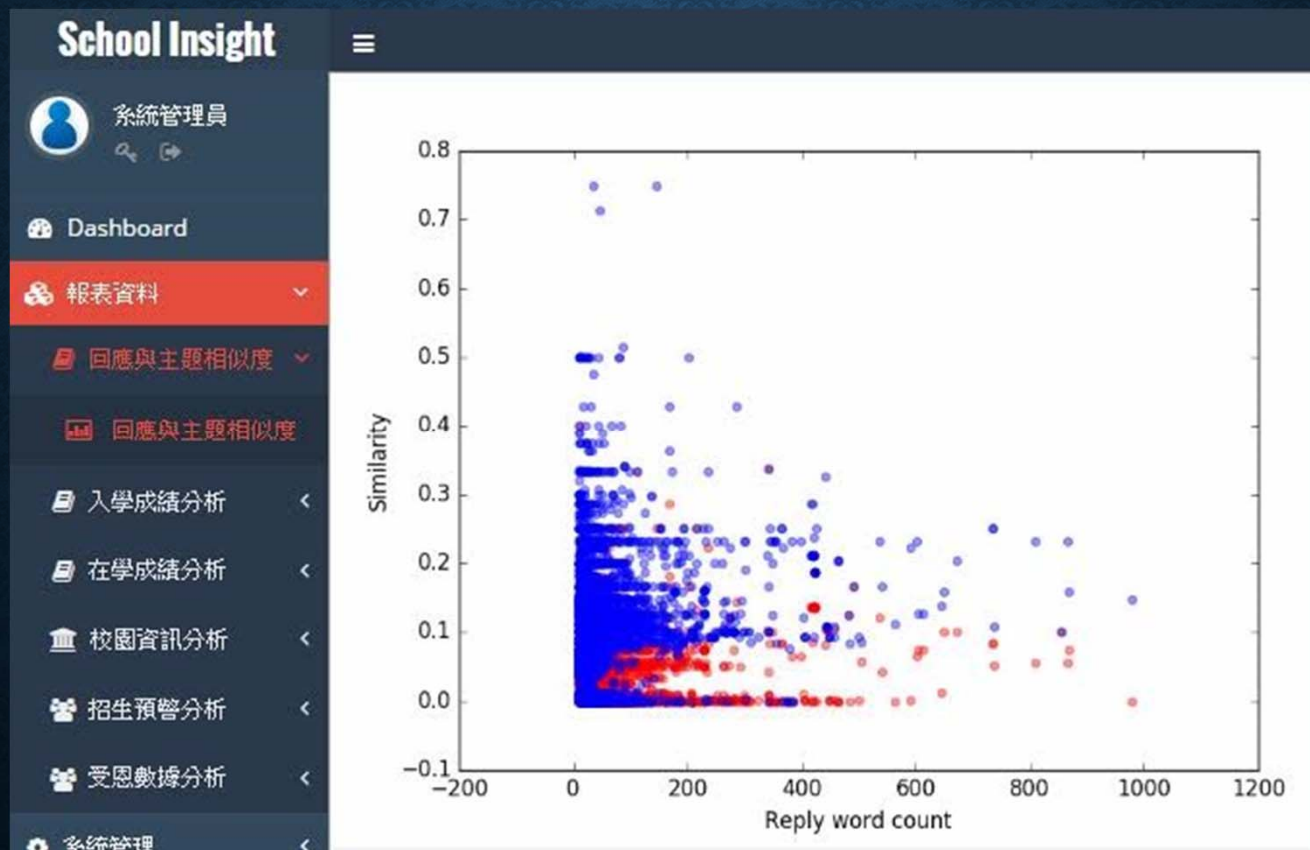
- before expansion

報表樣式: 散點圖

細項欄位: sim word TR

欄位名稱: before expansion

VOILÀ



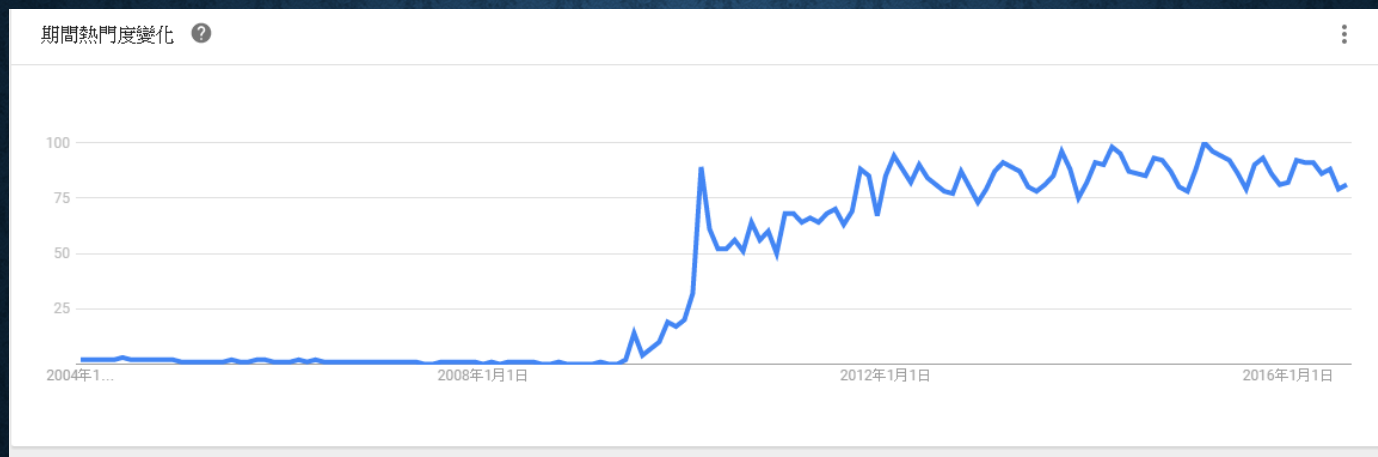
WHY SCHOOL INSIGHT?

1. Reusable modules
2. Easy to integrate web visualization libraries, such as Highchart, matplotlib, D3.js and Google Chart.
3. Easy to integrate analytics engine designed by common development tools, such as Python, Java and R.

ACTIONABLE ANALYTICS

- 定義：進行後可供實際決策的分析結果
- 範例：課程設計輔助系統
 - 最終目的：縮短學用落差
 - 方法：擷取企業所需且正在高成長的技術
 - 資料來源：104人力銀行、Wikipedia、Google Trend及教學計畫表

NOSQL



	數量	備註
104.com.tw	303職缺	2016.8.31資料
淡江課程	3門課	102-104, 同一教授

APPROACHES

1. Get job description
2. Extract key term candidates
3. Use Wikipedia to filter key terms
4. Extract Google Trend for these key terms
5. Count occurrence in teaching plans

GET JOB DESCRIPTION

- 方法：使用104 API
- 目標：職缺類別為軟體設計工程師
- 資料：共有4468個職缺
- 技術字彙來源：職稱、職缺描述及其他條件

EXTRACT KEY TERM CANDIDATES

- 使用斷詞器(Jieba)
- 以小單位為優先，例如「平行處理經驗」會被斷為「平行」「處理」「經驗」
- 技術字彙擷取時除了考慮單字，也考慮bigram

Nosql	Asp.net MVC	Machine learning	Stored Procedure	Programming skills
○	○	○	○	○

USE WIKIPEDIA TO FILTER KEY TERMS

- Only keep those key terms which are also [Wikipedia article title](#)
- Future work: Chinese Wikipedia

Nosql	Asp.net MVC	Machine learning	Stored Procedure	Programming skills
○	○	○	○	×

EXTRACT GOOGLE TREND FOR THESE KEY TERMS

- Use the difference between trends of last year and first year

Nosql	Asp.net MVC	Machine learning	Stored Procedure	Programming skills
○	○	○	×	×

COUNT OCCURRENCE IN TEACHING PLANS

- 22,898 teaching plans from 102-1 to 104-1

Nosql	Asp.net MVC	Machine learning	Stored Procedure	Programming skills
○	○	×	×	×

Term	Freq	MI/RP	Trend	Plan
asp.net mvc	68	44.86	83	0
android studio	43	39.7	91	4
entity framework	37	38.7	89	0
deep learning	48	36.88	78	0
mvc	451	17.75	33	2
node.js	131	17.55	90	0
json	125	17.48	92	4
angularjs	114	17.35	95	0
nosql	104	17.21	85	3
mongodb	73	16.7	95	3
github	67	16.58	89	2
laravel	61	16.45	85	0
xcode	59	16.4	38	3
openstack	52	16.21	85	0
opencv	47	16.07	54	5
devops	46	16.04	83	0
redis	45	16.01	84	3
nginx	44	15.97	92	0

WHY IT'S ACTIONABLE?

- 反例：學生畢業高中分佈
 - 雖然能知道分佈情形，但由於只是數字多寡，而且逐年不同，管理者該如何利用這個分析並無明確做法。
- 為什麼這個分析可執行？
 - 管理階層可明確了解業界所需及高成長、高需求的技術內容
 - 系所及教師可根據分析結果加強課程內容或訂定研究方向

THANK YOU!