

交大宅宅森77 靠北交大貼文分析

Group 7
0316213 蒲郁文 0316238 蔡孟軒 0316248 何鳳雯

Table of Contents

- 1. Introduction
- 2. Problem Definition
- 3. Related Works
- 4. Challenges
- 5. Dataset Description
- 6. Methods
- 7. Experimental Results
- 8. Issues Faced
- 9. Other Efforts
- 10. Conclusions & Future Works
- 11. Job Description
- 12. References

Introduction

靠北交大，交大人不可或缺的佈告欄。

學校發的公告沒人理，靠交上的資訊卻能迅速地傳達給交大人。

Introduction

10	教師送交本學期非應屆畢業 轉系申請截止	11
17	下午1:50 預約整骨	18

的佈告欄。

上的資訊卻能迅速地傳

Introduction



靠北交大
7月13日 · *

#靠交40179
"提醒各位 7/18 記得要去整骨
揪咪*.<
<http://i.imgur.com/fgZLmj5.jpg>"

教師送交本學期非應屆畢業生成績截止
7月 3日 (星期一)
暑期班開始
7月 10日 (星期一)
教師送交本學期非應屆畢業生成績截止
轉系申請截止
7月 18日 (星期二)
下午1:50 預約整骨

是誰○_○

i.imgur.com
I.IMGUR.COM

讚 留言 分享

你、[REDACTED] 和其他 1,630 人

最相關留言 ▾

Problem Definition

1. 靠北文有什麼常見的主題？
2. 靠北文主題和時間的關係？

Problem Definition

Input :

- 靠北交大貼文資料集

Output :

- 各貼文的主題
- 各主題的時間分佈

Related Works

A Comparison of Document Clustering Techniques
(KDD 2000) Steinbach, M., Karypis, G., & Kumar, V.

- K-means
 - better
- agglomerative hierarchical clustering
 - performs poorly
 - since, in many cases, the nearest neighbors of a document are of different classes

Challenges

- Posts on social networking sites are dirty.
- Have a lot of grammar errors, spelling errors, etc.

Dataset Description

	ID	Link	Message	Created_time	Story
Data Type	粉專編號_文章編號	連結：分享文、來源文章連結	貼文內容	貼文發布時間	動態描述：分享文、來源粉專資訊
Example	55_94	https://www.facebook.com/NCTULIB/photos/a.688796024509631.10737	#靠交49561 圖書館用我的照片 都不用講一下 的? https://www.facebook.com/NCTULIB/posts/1659934950729062	2017-12-05T04:34:12+0000	靠北交大 shared 國立交通大學圖書館's post.

Dataset D

靠北交大分享了國立交通大學圖書館National Chiao Tung University Library 的貼文。

2017年12月5日 · *

#靠交49561

"圖書館用我的照片都不用講一下的？

<https://www.facebook.com/NCTULIB/posts/1659934950729062>



國立交通大學圖書館National Chiao Tung University Library
2017年12月4日 · ●

歐巴就是要這麼大看才過癮阿！
圖書館新玩具！超大平板，挑戰眼珠可移動的範圍

讚 留言 分享

19

最相關留言 ▾

	ID	Link	reated_me	Story
Data Type	粉專編號_文章編號	連結來源	貼文發布時間	動態描述：分享文、來源粉專資訊
Example	55_94	http://com.pho602.737	2017-12-05T4:34:12+000	靠北交大 shared 國立交通大學圖書館's post.

Methods

Data Collection

Data Preprocessing

Clustering

K-means

Latent Dirichlet
allocation(LDA)

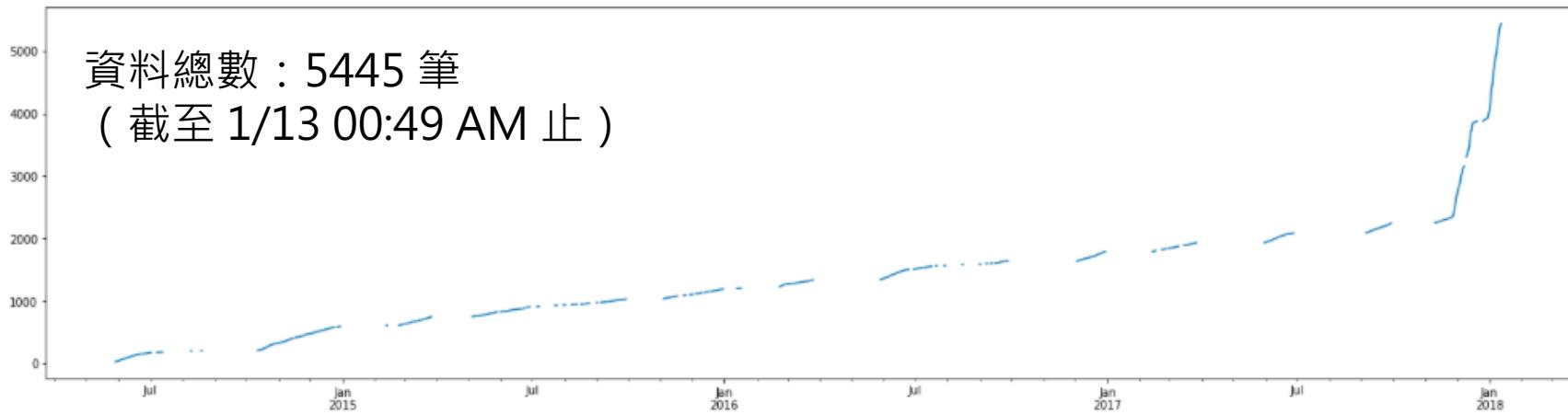
Non-negative matrix
factorization (NMF)

Data Collection

不斷透過 Facebook Graph API 蔊集資料

(最後更新時間 : 2018/01/03 00:49)

Data Observations



累計貼文數量對時間的關係圖

Data Preprocessing

1. Remove posts that have no content.
2. Word Segmentation: Jieba + human intelligence
(e.g. 梅竹, 13舍)
3. Removing Stopwords: SnowNLP stopword list +
human intelligence (e.g. 😂, 哈哈哈)
4. Build TF-IDF / TF document vectors.
(using both unigrams and bigrams; 14275 features in total)

Clustering Methods

1. Hard-Clustering

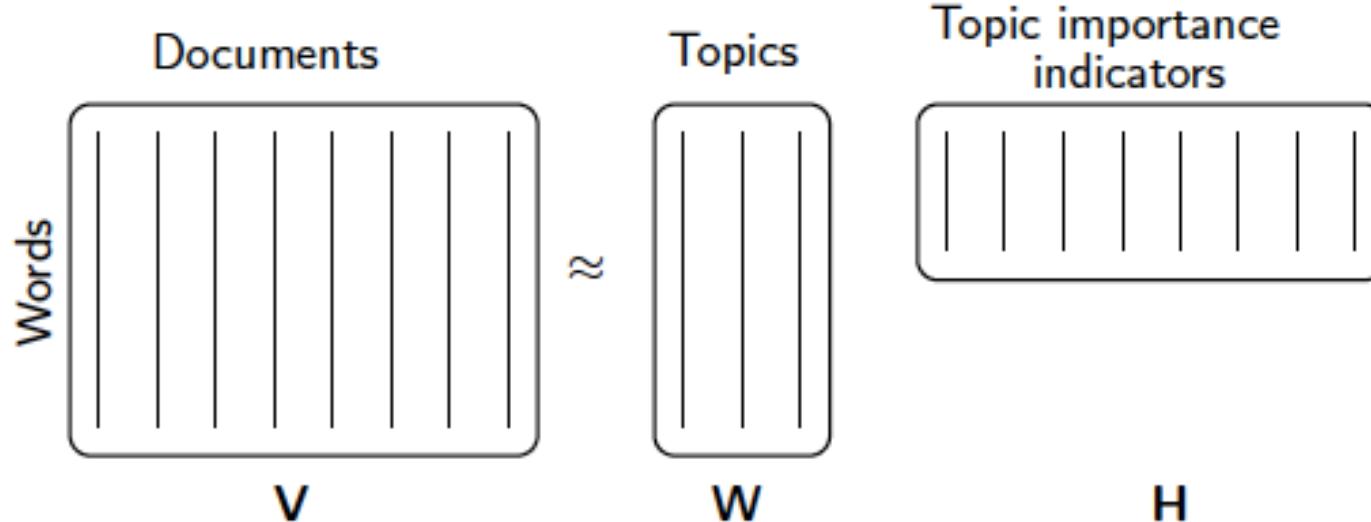
a. K-means

2. Soft-Clustering

a. Latent Dirichlet Allocation (LDA)

b. Non-negative Matrix Factorization (NMF)

Clustering (NMF)



source: <http://www.cnblogs.com/pinard/p/6812011.html> 17

Clustering (LDA)

Topics

gene	0.04
dna	0.02
genetic	0.01
...	

life	0.02
evolve	0.01
organism	0.01
...	

brain	0.04
neuron	0.02
nerve	0.01
...	

data	0.02
number	0.02
computer	0.01
...	

Documents

Topic proportions and assignments

Seeking Life's Bare (Genetic) Necessities

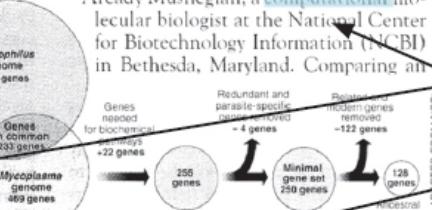
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains

Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.



Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

source: <http://deliveryimages.acm.org/10.1145/2140000/2133826/figs/f1.jpg>

Experimental Results

Experimental Results (NMF)

Topic 0: 2592

別人 / 講 / 尊重 / 事 / 只會 / 討厭 / 懂得 / 大聲 / 讀書 / 別人 讀書 / 感到 / 整天 / 衣服

Topic 6: 102

助教 / 學店 / 分數 / 課 / 學生 / 教授 / 作弊 / 考卷 / 可愛 助教 / 題目 / 成績 / 安安 / 跪求

Topic 8: 156

吃 / 吃 吃 / 滷味 / 蔬菜 / 路 / 二餐 / 好吃 / 滷味 吃 / 我加 / 面 / 一餐 / 夜市 / 姐妹 / 女二

Topic 11: 120

梅竹 / 清 / 梅竹 清 / 必勝 / 倒數 / 比賽 / 梅竹賽 / 女籃 / 宣傳 / 清大 / 梅竹 倒數 / 清 女籃

Topic 12: 178

考試 / 作業 / 寫 / 寫 作業 / 交 作業 / 期末考 / 交 / 考試 作業 / 期中考 / 上課 / code / e3

Experimental Results (LDA)

Topic 3: 135

臺灣 / 電影 / 路上 / 幸福 / 中 / 幸福路上 / 社團 / 故事 / 希望 / 手機 / 透過 / 學生 / 領袖

Topic 4: 210

辦 / 活動 / 營隊 / 系學會 / 學妹 / 梅竹 / 參加 / 社團 / 妳 / 女生 / 找 / 系 / 明明 / 學校

Topic 7: 180

吃 / 狗 / 我會 / 留言 / 宿舍 / 感謝 / 同學 / 檢 / 咬 / 私訊 / 怕 / 留言我會 / 希望 / 樓

Topic 10: 147

倒數 / 助教 / 題目 / 學生 / 同學 / 時間 / 2de / ICLAB / 期末 / pattern / project / 梅竹 / 學校

Topic 16: 184

系統 / 妳 / Windows / 內容 / 目標 / 環境 / 型別 / 配置 / 攻擊 / 程式設計 / 方法 / 原理

Experimental Results (K-means)

Topic 7: 137

文 / 聽說 / 告白 / 搬運 / 搬 / 帥哥 / 搬運工 / 哥 / 貼文 / 系 / index / 搬文 / 屁孩 / ☺

Topic 14: 148

找 / 系學會 / 明明 / 學妹 / 朋友 / 辦 / 時間 / 社團 / 幫忙 / 活動 / 還要 / 學長 / 營隊 / 逆

Topic 15: 72

梅竹 / 清 / 比賽 / 梅竹賽 / 倒數 / 必勝 / 梅後 / 梅竹倒數 / 宣傳 / 梅竹清 / 明年 / 清大 / 號

Topic 19: 106

助教 / 微積分 / 寫 / 掛 / 考試 / 作業 / 物理 / 微積分大會考 / 物理微積分 / 題目 / 分數

Topic 24: 90

八卦 / 畢業 / 實驗室 / 門檻 / 畢業門檻 / 畢業找 / swap 八卦 / swap / 女友 / 找 / 聽說

Results Visualization (NMF)

男朋友 可愛 櫃檯 運管
啦啦隊 蛋包飯 機 樂隊 男朋友 聽說
妹子 你好 女孩 #怪 耶舞 力 阿西 鑽石 名字
可愛 姊 電資 機械 奈米 阿西 可愛 騷擾
水 戴眼鏡 學期結束 阿西 師 可愛 tag 二餐 笨股 七
可愛 男友 可愛 魔術 阿西 師 可愛 二餐 笨股 迷人
你好 可愛 室友 可愛 樣子 可愛 可愛 笨股 學姊
付 告白 土木系 名字 李 師 認識一下 進場 笨股 皮卡丘
應 條子 土木系 名字 李 教 倒數 倒數 可愛 可愛 可愛
你 死會 打球 追 可愛 機 可愛 可愛 跳高
死會 打球 跳舞 小天使 可愛 可愛 助教 我家 高妹
號樣子 倒數 跳舞 可愛 可愛 可愛 電物 馬尾 好奇
教 梅竹 倒數 小天使 可愛 可愛 可愛 電物 馬尾 好奇
樣子 男友 可愛 外文 發 復活 想要 男朋友 戴 可惜
你 可愛 女生 可愛 奶廚 店員 恩

直播 瞧不起 肥宅 不見
明年 天硬幣 賴 亂
女生 tag 交清 掛 不少
帥 隔壁 棒球隊 布條
打爆 收拾 清 不見
吹 梅竹 找
工具 梅竹 倒數 热舞
換 男排 號 梅竹 後援會 綠色
商學院 忘記 班
兩校 銅板 火力
吹 梅竹 賽
大 梅竹 母 梅後
全體 文章
後援會 清交 攻陷
臺下 梅竹 傳
梅竹 竹狐 是 梅後
太帥 竹狐 是 梅後
有大 竹狐 是 梅後
model 梅竹 播 轉
夜市 热情 昨天
球館 球票
直播 必勝 成功 號
記者 喇 聲
研究生 橋藝 簡
梅 吃掉 顯眼
男排 男排 男排
梅竹 男排 男排
跑跑 跑跑 跑跑
加油 加油 加油
掰掰 女朋友
啦啦 大學
支援 體育館

外面 演唱會 公共空間 不可思議 草地 幾個 周到

施工 LOL 衛生紙 力點 爛 討論 路

幾點 用氣音 很慢書 疊 情侶 男 絲 路

整間 五樓 滑鼠耳聲 六樓 浩然 那區 生 咳想 我要 降

擦屁股 靜音 故意 講話 同 打字 還 拜託 咯

洗澡 放閃 講 隨機 機打字 過期 煙味 麻煩 邀

旁邊 一把 祝 為啥 學妹 圖書館 公共 次走

安靜 過夜 一團 活動 一群 時大聲 週末 讀書

拉基 求解 吵 圖書館 照片 一首 聊天 大聲 聊天

拜託 路

滾出去 五樓 叫死 一團 聊天 滑鼠 那區

圖書館 大聲 聊天



朋友
Tag
身邊
Tag
身邊

比賽
 男生 朋友 睡不著
 朋友 他會
 笑臉
 有時候
 整天
 感覺
 大學
 光粉
 臭
 陪
 男
 有個朋友
 留言
 關係
 聊
 離
 黑人
 傻
 開門
 這科
 大學生
 作業
 49286
 興趣

缺
 宅宅
 難
 DT
 小聲
 句
 女朋友
 身上
 多麼
 背地裡
 課
 生活
 順便
 半夜
 睡不著
 朋友
 關係
 XXX
 睡
 tag
 姓黃
 半夜
 遊
 輕鬆
 社會
 看不出
 又帥獻世
 欠
 看不
 希望

紹
 家
 搭車
 49879
 Ans
 熟
 照鏡子
 課堂
 人生
 重修夥伴
 偷偷
 房間
 走

朋友
Tag
身邊
Tag
身邊

比賽
 男生 朋友 睡不著
 朋友 他會
 笑臉
 有時候
 整天
 感覺
 大學
 光粉
 臭
 陪
 男
 有個朋友
 留言
 關係
 聊
 離
 黑人
 傻
 開門
 這科
 大學生
 作業
 49286
 興趣

難
 句
 女朋友
 身上
 多麼
 背地裡
 課
 生活
 順便
 半夜
 睡不著
 朋友
 關係
 XXX
 睡
 tag
 姓黃
 半夜
 遊
 輕鬆
 社會
 看不出
 又帥獻世
 欠
 看不
 希望

紹
 家
 搭車
 49879
 Ans
 熟
 照鏡子
 課堂
 人生
 重修夥伴
 偷偷
 房間
 走

陶 滾 長 拒吃 很潮 土木 寫個

一首 麥當勞 優質 bbs NCTU 信徒

團購 成效 NCTU 發了 SENIORHIGH 置頂 級

bbs SENIORHIGH 國家 3853 禮拜 置頂 中午

詩 澱 友 文 學弟 113 支援 TALK 奬學金

酸 問卦 NCTU TALK E6 よ 新聞 買 便宜 教主

發現 質答 明天 中午 sex 明天 繼 性

ptt SENIORHIGH 19204 一條

52795 交上

GOSSIPING

bbs Gossiping

This image is a word cloud centered around the theme of dogs (狗). The size of each character varies based on its frequency or importance in the context. Key characters include:

- 狗** (Dog): The central and largest character.
- 单身** (Single): A large blue character.
- 女二** (Female No. 2): A large purple character.
- 汪汪** (Wang Wang): A red character.
- 愛** (Love): A yellow character.
- 累** (Tired): A yellow character.
- 忙** (Busy): A green character.
- 叫** (Call): A green character.

Other prominent characters include:

- 邊** (Edge), **舍** (Dormitory), **人** (Person), **士** (Person), **夜** (Night), **半** (Half), **單** (Single), **身** (Body), **邊** (Edge), **舍** (Dormitory).
- 方位词: 上 (Up), 下 (Down), 左 (Left), 右 (Right).
- 状态词: 忙 (Busy), 累 (Tired), 快 (Fast), 痛 (Pain).
- 行为词: 跑 (Run), 走 (Walk), 咬 (Bite), 抓 (Catch), 打 (Fight), 攻擊 (Attack).
- 地点词: 宿舍区 (Dormitory Area), 校园 (Campus), 图书馆 (Library), 学校 (School), 街市 (Street Market), 市场 (Market), 地下室 (Basement).
- 人物词: 女生 (Female Student), 女二 (Female No. 2), 女三 (Female No. 3), 工人 (Worker), 阿福 (AFU).
- 动物词: 狗 (Dog), 猫 (Cat), 鸡 (Chicken), 猪 (Pig), 蟑螂 (Cockroach).
- 生活词: 食物 (Food), 粮食 (Grain), 猪油 (Pork Fat), 酱油 (Soy Sauce), 老虎 (Tiger), 蜈蚣 (Centipede).
- 文化词: 传统 (Traditional), 现代 (Modern), 未来 (Future).
- 时间词: 昨天 (Yesterday), 今天 (Today), 明天 (Tomorrow), 未来 (Future).
- 其他词: 喵 (Meow), 汪 (Wang), 呼籲 (Appeal), 呼喊 (Shout), 挑逗 (Tease), 趾高气扬 (Proud), 可悲 (Sad).

The background features a large watermark of a dog's face with the text "NIHON NO SHIBAKEN" and the date "2018年1月".

哥哥 直播 正義影片 嘻嘻 大喊 考 那位 那篇 微積分 痘癰 憂 大會考 留言 求學 幾篇 中提到 宿舍 死忠 支援 中毒 手槍 救 光頭 哥哥 哥哥 哥哥 哥哥 光頭 母湯 欧 哥哥 哥哥 哥哥 哥哥 光頭 哥哥 母湯 TAG 明明出門 看著

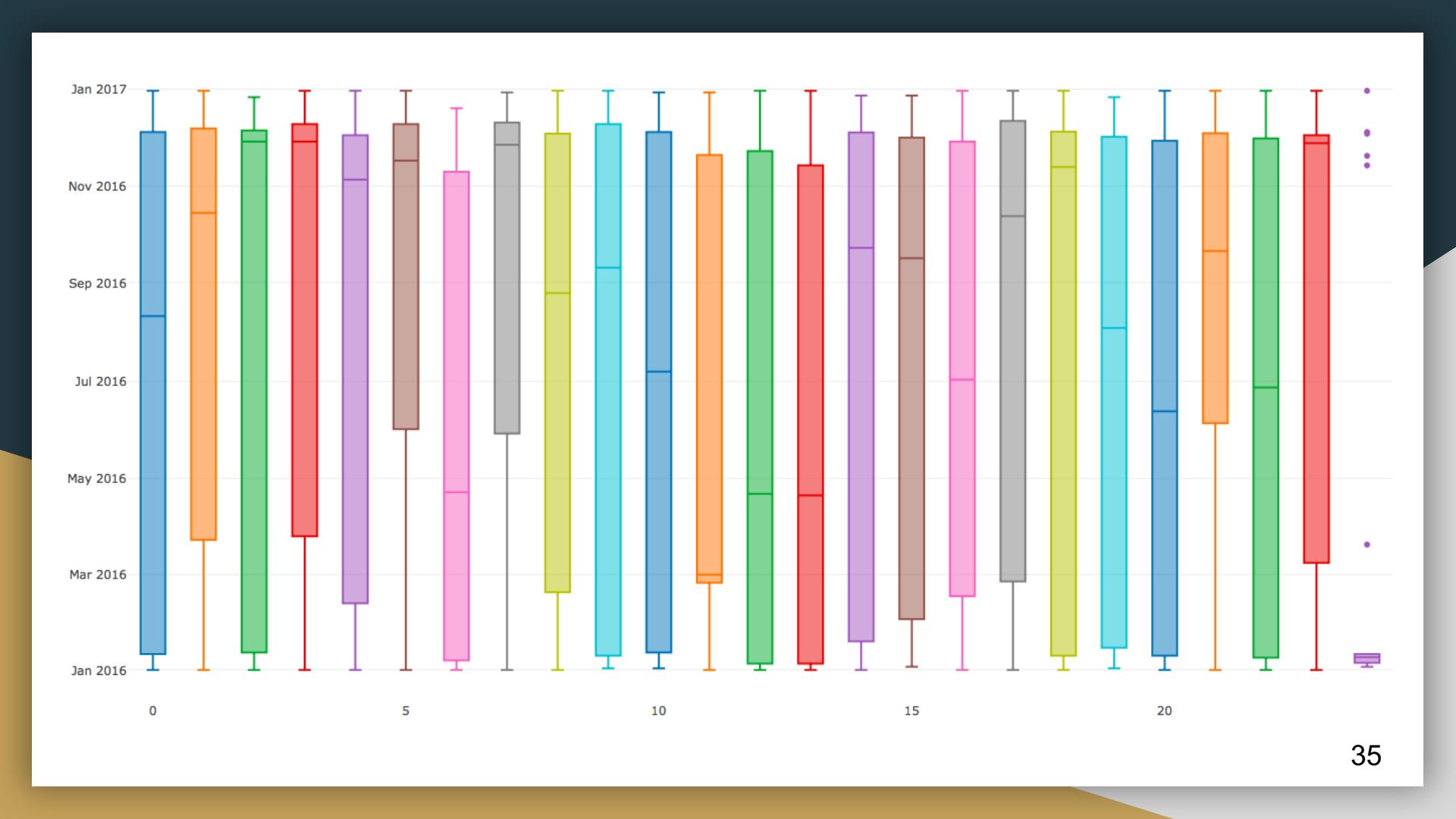
希望 光頭 加油 滿分 粉 好帥 看個 中 七舍 七舍 半夜 微積分 痘癰 憂 大會考 留言 求學 幾篇 中提到 宿舍 死忠 支援 中毒 手槍 救 光頭 哥哥 哥哥 哥哥 哥哥 光頭 母湯 欧 哥哥 哥哥 哥哥 哥哥 光頭 哥哥 母湯 TAG 明明出門 看著

對付 希望 同學 那位 聽說 大學畢業 笑話 大學畢業 考 滿分 學 期末考 分 竟然 大會考 哥哥 哥哥 哥哥 哥哥 光頭 哥哥 母湯 TAG 明明出門 看著

底下 欧 提到 試試看 版本 哥哥 哥哥 哥哥 哥哥 光頭 哥哥 母湯 TAG 明明出門 看著

多人 娛樂 哥哥 哥哥 哥哥 哥哥 光頭 哥哥 母湯 TAG 明明出門 看著

哥哥 直播 正義影片 嘻嘻 大喊 考 那位 那篇 微積分 痘癰 憂 大會考 留言 求學 幾篇 中提到 宿舍 死忠 支援 中毒 手槍 救 光頭 哥哥 哥哥 哥哥 哥哥 光頭 母湯 欧 哥哥 哥哥 哥哥 哥哥 光頭 哥哥 母湯 TAG 明明出門 看著



Issues Faced

- Data Collection
 - Facebook 對資料蒐集的限制
- 自然語言處理
 - 文本分群演算法的選擇
 - 最佳化各方法的參數
- 成果呈現
 - 如何視覺化

Other Efforts

- Agglomerative Clustering
 - 效果不佳
 - 與文獻結論一致

Conclusions

- 自動化找出貼文的主題是可行的
- 分群效果：NMF > k-means > LDA
- 分群時間：LDA (10.1 s) > k-means (8.08 s) >
NMF (997 ms)

Future Works

- 結合貼文讚數、留言
 - 議題熱門程度分析
- 結合天氣資料
 - 下雨、寒冷是否會讓人更容易 7pupu
- 將本專案整理成開源套件
 - 目前處理繁體中文社群網站文本的套件仍十分稀少

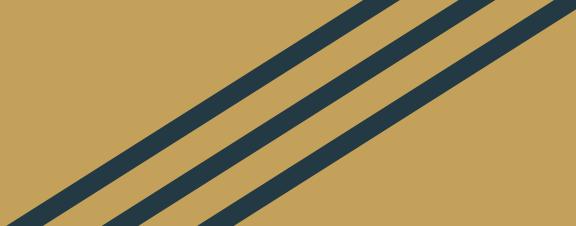
Job Description

- 蒲郁文(34%) - data preprocessing; clustering
- 蔡孟軒(33%) - clustering; visualization
- 何鳳雯(33%) - data collection; clustering

References

1. Xu, W., Liu, X., & Gong, Y. (2003, July). Document clustering based on non-negative matrix factorization. In Proceedings of the 26th annual international ACM SIGIR conference on Research and development in information retrieval (pp. 267-273). ACM.
2. Steinbach, M., Karypis, G., & Kumar, V. (2000, August). A comparison of document clustering techniques. In KDD workshop on text mining (Vol. 400, No. 1, pp. 525-526).
3. Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. Journal of Machine Learning Research, 3(Jan), 993-1022.

Q & A



Thank You!