

# 社群媒介輿論分析111

## Announcement

如欲修課的同學，請先填寫修課意願表 (<https://forms.gle/XCBU6tB15zEKRrQWA>)。並請於開學第一週前將近500字研究提案以E-mail寄送給開課教師(MUST)，研究提案內容主要為你感興趣的研究議題是什麼？亦可參考範本撰寫([Social Opinion Proposal Template](#))。如果已經考過Proposal或報題者，請直接繳交Proposal或報題表單即可。

## Description

因為在社群平台上，人人都可以是資訊生產者，因此，巨量、去中心化、扁平化、流動、匿名、即興等社群平台的特色亦影響著訊息的特性、傳散方式與效果，與過往的集中式、權威式的訊息大不相同。尤其是Twitter或Youtube這樣的社群平台，留過什麼言、Follow過誰、Retweet過誰都被記錄下來。近期Twitter甚至開放一則Tweet被閱讀過幾次的功能。透過這樣的平台，我們可以研究例謠言的流傳、用語的流變、意識形態或政治立場的極化、泛政治化的現象、數位公眾外交的互動網絡、網軍偵測、對性別議題的看法、民粹主義、負面黨性的言談等。

本課程採用的方法會是大數據方法、量化方法、文本探勘方法。甫接觸計算機方法在社會議題的應用時，往往會過度粗暴地運用計算工具對文本資料直接做分析。也因此在講課的同時，本課程也希望帶同學也必須思考，相較於過往社會科學的研究方法，無論是質性或量化，計算機科學是一個新方法嗎？還是他只是用以克服大數據的挑戰？他和傳統方法又有何異同？使用上要特別注意哪些環節呢？如何適切地用這些演算法工具？而過去內容分析方法必定有許多常見的理論，又有哪些理論所關係到的研究，常常被計算機方法所運用呢？

## Purpose

本門課有兩個重點，第一個是社群輿論的研究方法與寫作。本課程將帶學生從事社群媒介上的輿論分析研究。常見的社群媒介分析對象包含臉書、PTT、推特、YOUTUBE等，但由於資料取用限制，本課程將以YOUTUBE與推特為主。授課過程將從選讀文獻、尋找研究議題、撰寫研究目的與問題、規劃文獻、實作研究方法並進行信效度評估、到結果分析完成一個小論文。

第二個重點是介紹文本探勘的各種方法，包含詞嵌入、文本分群、文本分類、網絡方法、主題模型等。除了方法的使用介紹外，將著重在閱讀前人研究的應用來了解方法的適用時機與信效度的驗證。學生會需要在老師引導下，閱讀並報告論文。

## Requirement

本學期課程將帶學生撰寫個人Essay，需以社群輿論為分析對象，以傳播為議題如政治傳播、科學傳播、災害傳播、計算社會科學等。Essay內容應包含問題意識、研究問題、文獻探討架構與核心論文摘要、研究方法、和初步研究成果。大學生應以大專生研究計畫、中華傳播學會年會、資訊社會學年會、計算傳播年會投稿為目標；碩士生應以寫出嚴謹的研究計畫書為目標，並提供初步的分析結果。

本課程並非程式語言教學課程，修課生必須至少修過一學期的程式課程，且需有能夠自行撰寫程式的自信。學生將被要求依照所介紹的論文與分析方法，就自己的研究議題自行實作。程式部分會有線上課程可以看，亦列計在授課時數中。

## Textbook & Resources

- PDF of Reference <https://paperpile.com/shared/HzGKFN>
- Stefanowitsch, A. (2020). **Corpus linguistics: A guide to the methodology.** Language Science Press. [Corpus linguistics](#)
- **Text Mining: A Guidebook for the Social Sciences** (sagepub.com): This book can be accessed for free through NTU VPN. It was written by Gabe Ignatow and Rada Mihalcea. The former is a sociologist from North Texas and the latter is a computer scientist. In addition to introducing the common techniques of text mining, the book also introduces three traditional social science methods related to text mining, including thematic analysis, narrative analysis, and metaphor analysis.
- Jungherr, A. (2015). **Analyzing Political Communication with Digital Trace Data:** The Role of Twitter Messages in Social Science Research. Springer International Publishing. <https://doi.org/10.1007/978-3-319-20319-5>
- Desagulier, G. (2017). *Corpus Linguistics and Statistics with R*. Springer International Publishing. <https://doi.org/10.1007/978-3-319-64572-8>
- Jockers, M. L., & Thalken, R. (n.d.). *Text Analysis with R*. Springer International Publishing. <https://doi.org/10.1007/978-3-030-39643-5>
- Baker, P. (2006). Using Corpora in Discourse Analysis. A&C Black. <https://play.google.com/store/books/details?id=QMsSpoNX2EMC>
- [quanteda tutorials :: Tutorials for quanteda](#)

## Grading Policy

### Essay (45%)

(All presentations are peer-graded by classmates)

1. Pitching and research questions 5% (Standup presentation + doc)
2. Literature review & Methods 5%
3. Exploratory data analysis 5% (Presentation + doc)
4. Final essay and presentation 30%: (Presentation 5%, Essay 25% at least 10 pages in APA format with correct reference style, and figure/table captions)
  - a. Completeness (15%), Writing format and APA (5%), Logic (5%)

### A ssignment & Practice (55%)

- Note-taking Readings:15%: Summarize paper with bolded title
- Reporting Selected topics 15%: Select one topic to report and lead discussion
- Implementation 25%: Keyness, Collocation and POS, word embeddings, and topic modeling. Required to implement skills or algorithms for specific topics such as political communication, health communication, or user modeling (Examples of tech collocation with PMI, topic modeling and visualization, crawling wiki DB for improving tokenization, using the tech word-embedding for doc classification). Each assignment has two weeks to work on, but there will be no delay in due dates

## Calendar

W	Date	Tutorials	Students
---	------	-----------	----------

1	2/20	Social Opinion Research overview: 1. Twitter related researches overview 2. Journals, papers, scholars, and issues (Stereotypes, polarizations, echo chambers, and filter bubble)	Note Taking
2	2/27	228 Break	
3	3/06	Computational linguistics and text mining for Political Communication	
4	3/13	Process of information retrievals & Research Designs ● Gathering corpus for Youth representation on news ● Populist and Populism: Big data methods vs. Survey	Pitching I Note Taking
5	3/20	Keyness ● Anomaly detection, event detection, and Hashtag activism ● R Student report by 李宜恬	Note Taking
6	3/27	Key Influencer ● Netizen, spammer, cyber warrior, and key influences in digital diplomacy and Information operation ● R Student report by 張聖雯	Note Taking
7	4/03	Spring break	
8	4/10	N-gram, Collocation & Network method applications ● How collocation network vis assists text understanding ● R Student report by 何家慈	Literature Review & Method A Keyness & Collocation
9	4/17	Topic Modeling: Reliability and Interpretability	Note Taking
10	4/24	Topic modeling, Agenda-setting and Framing ● R Student report by 洪欣愉	A Topic Modeling
11	5/01	Word embeddings and stereotypes ● Ideology scaling of US's News Media	E Exploratory data analysis
12	5/08	Ideology Scaling & ENA ● Wordfish, WordScore, Filter bubble, echo chamber and segregation ● R Student report by 巫敏慧	A Word Embeddings Note Taking
13	5/15	Sentiment analysis, POS, and Emoji ● Dictionary methods: Using E-hownet, Snownlp, and NTUSD ● R Student report by 曹惟純	Note Taking
14	5/22	Stance-based sentiment analysis & Classification: Support vs. against ● Support or against Hillary? Stance toward climate and feminist issues.	
15	5/29	Concluding	
16	6/05	Final	
17	6/12	E Final Presentation (In-Class)	

## Weekly

## W1 Overview

- TM00\_OVERVIEW

### Reference

- Jungherr, A. (2015). *Analyzing Political Communication with Digital Trace Data: The Role of Twitter Messages in Social Science Research*. Springer International Publishing.  
<https://doi.org/10.1007/978-3-319-20319-5>

### Take-home Practice I: Twitter as corpus

Twitter由於資料容易取得，且有作者名、時間、群體、轉推、回應等相當豐富的社群行為資訊，所以可以找到相當多類型的研究如政治科學、社會學、傳播學、地理、環境、文學、語言學等。這個練習是為了讓你找到本學期要從事的研究題目，或讓你更習慣查找文獻以發掘新題目。

請從Google Scholar或Jungherr, A. (2015). *Analyzing Political Communication with Digital Trace Data: The Role of Twitter Messages in Social Science Research*. Springer International Publishing. 挑選一篇論文，並記錄在[SMOA-PaperDB](#)中。**請勿和他人重複**。

### Reading for next next week

從下週的文獻中挑選一篇，嘗試歸納用文本探勘、計算語言學可以做什麼樣的研究？或者什麼樣的研究對你而言是有趣的？本次閱讀是希望你從上述文獻中，摘要出與社群輿論分析、文本探勘或計算語言學在傳播科學或政治傳播上的應用。將你閱讀的篇目貼至[SMOA-PaperDB](#)。

通常一篇摘要會希望你摘要近一千字或更多，在這門課並沒有限制你不能使用Google Translate或任何你認為可以用的工具，但摘要文字必須是通順的、經過整理的。且要能夠擷取重要的圖、文進入摘要中。你也可以將你做完註記後的文獻電子檔直接貼至該Notion Page中，來提供你確實好好閱讀過該篇文章的證明。比方說，你可以非常提綱挈領地摘要1000字，但其他全部用註記PDF的方式來達成。

Table

### SMOA-PaperDB

Week	Title	Citation	Year	Note by	Notes	Tags	APA
W1	Stance detection for online public opinion awareness	Cao et al. (2022)	2022	Hsieh	A review of all stance-based opinion analysis	Stance	Cao, R., Luo, X.,

+ New

## W3 Computational Linguistics and text mining for Political Comm.

- TM02-Ling\_Comm\_PoliSci\_Computing

### Reference

- Wilkerson, J. D., & Casas, A. (2016). Large-scale Computerized Text Analysis in Political Science: Opportunities and Challenges. *Annual Review of Political Science*, 20, 529–544.
  - 嘗試舉例說明：「政治學（或社會科學）會比較注重變項間的交互作用，例如多控制一個變項，可能為模型帶來的影響，而不是如機器學習般停止於預測的準確率」

- Molina, M., & Garip, F. (2019). Machine Learning for Sociology. *Annual Review of Sociology*. <https://doi.org/10.1146/annurev-soc-073117-041106>
- Theocharis, Y., & Jungherr, A. (2021). Computational social science and the study of political communication. *Political Communication*, 38(1-2), 1-22.
- Nguyen, D., Değruöz, A. S., Rosé, C. P., & de Jong, F. (n.d.). Computational Sociolinguistics: A Survey. *Computational Linguistics*. [https://doi.org/10.1162/COLI\\_a\\_00258](https://doi.org/10.1162/COLI_a_00258)
- Lucas, C., Nielsen, R. A., Roberts, M. E., Stewart, B. M., Storer, A., & Tingley, D. (2015). Computer-Assisted Text Analysis for Comparative Politics. *Political Analysis: An Annual Publication of the Methodology Section of the American Political Science Association*, 23(2), 254-277. <https://doi.org/10.1093/pan/mpu019>

## W4-1 Corpus

- TM00 Rundown

Barberá. (2021) 的研究中提到了建立語料庫的可能方法和考量。就你個人的研究而言，你的語料庫是什麼？如何說明你所選用的語料庫範圍和抽樣方法是相關(relevant)且具有代表性的(Representative)？

- 何謂「語料庫」？
- 文中的Keyword search和Subject-based兩種方法所指為何？
- 文中所指的「相關(Relevant)」與「具代表性(Representative)」是什麼意思？對你的研究而言，相關文本為何？是否具代表性？
- Example:
  - [Fathering keywords](#): 如果要做BabyMother版的父職論述，抽老公、丈夫、隊友是否相關且具代表性？
  - 如果要做年輕人的新聞評述，抽取「年輕人」是否相關且具代表性？

研究者應就其研究問題界定適當範圍的語料庫。如果是自行爬取資料來建立語料庫的話，必須要思考是用關鍵字搜尋或依賴於既定的資料範圍。如果是用關鍵字搜尋的話，尤其是要呈現議題趨勢變化的話，必須嘗試證明所下的關鍵字能夠準確地收集到相關的資料，既不濫收也不少收。在研究方法中必須要說明所獲得(或預計獲得)的資料範圍，來自哪個媒體或平台？原始資料總筆數？時間範圍為何？

- □ TM02-Introduction to NLP
- [Preparing Corpus \(Notion\)](#)
- <https://cclabtw.notion.site/The-degree-evaluative-construction-Grammaticalization-in-constructionalization>
- [國教院--華語文語料庫應用系統整合入口 \(naer.edu.tw\)](#)

## Reference

- Barberá, P., Boydston, A. E., Linn, S., McMahon, R., & Nagler, J. (2021). **Automated Text Classification of News Articles: A Practical Guide**. *Political Analysis*, 29(1), 19–42. <https://doi.org/10.1017/pan.2020.8>
- Hyland, K., & Tse, P. (2005). **Hooking the reader: a corpus study of evaluative that in abstracts. English for Specific Purposes**, 24(2), 123–139. <https://doi.org/10.1016/j.esp.2004.02.002>
- Liu, M. C., & Chang, C. (2012). The degree-evaluative construction: Grammaticalization in constructionalization. Newest trends in the study of grammaticalization and lexicalization in Chinese, ed. Janet Zhiqun Xing, 115-148.
- 吳雙羽, & 賴惠玲. (2022). 從範例模型及詞彙構式互動的觀點探討新型「被+X」構式. *Taiwan Journal of Linguistics*, 20(1), 115–150. [https://doi.org/10.6519/TJL.202201\\_20\(1\).0003](https://doi.org/10.6519/TJL.202201_20(1).0003)

- 鍾曉芳, & 曾瑋庭. (2022). 類別隱喻構式研究. 中國語文通訊, 101(2), 193–212.

## W4-2 Research designs

- [Research designs](#)
- Examples
  - Gathering corpus for Youth representation on news
  - Populist and Populism: Big data methods vs. Survey
- Writing guidelines
  - How to seek and organize papers with citation managers
  - How to read papers by notions
  - How to quote a paper

這三本都是相當高分的碩論，但除卻議題本身，他們在理論與方法的運用與結合上，大相徑庭，和傳統方法的綜合運用和比較上、就方法而言、和理論對話上，你覺得他們有什麼樣的差異？請各簡要用100字做個摘要。<https://paperpile.com/shared/HzGKFN>

- 江玟. (2021). *Populism and Political Communication on Facebook\_ 2020 Taiwan Presidential Election*. [https://doi.org/10.1016/S1541-4612\(21\)00035-5](https://doi.org/10.1016/S1541-4612(21)00035-5)
- 黃子晞. (2019). 從媒體微博之限制評論設置看中國網路輿論限縮 [國立臺灣大學]. <https://doi.org/10.6342/NTU201901458>
- 齊若堯. (n.d.). *Media representation and bias against Southeast Asia: Analysis of Taiwan's media coverage of Filipinos*.
- (Good) Grimmer, J., & Stewart, B. M. (2013). Text as data: The promise and pitfalls of automatic content analysis methods for political texts. *Political analysis*, 21(3), 267-297.

## W5 Keyness & Hashtags & POS

- What does **Zipf's law** say? i.e., Zipf's law describes the distribution of term frequency.
- What is **keyness**? Is it different from keywords, keyterms or **term frequency**?
- What is **tf-idf**? What is its relationship with term frequency? What problem does it aim to solve?
- What is term **dispersion**? How is it different from term frequency and tf-idf?

- [AS06 Tokenization and TF.ipynb \(Colab\)](#)
- [TM04 Chinese Tokenization.ipynb \(Colab\)](#)
- <https://cclab-tm.notion.site/Segmentation-Tokenization>
- <https://cclab-tm.notion.site/Keyness>
- [Text analysis of Trump's tweets confirms he writes only the \(angrier\) Android half – Variance Explained](#)

## Reference

- Baker, P. (2004). **Querying keywords: Questions of difference, frequency, and sense in keywords analysis.** Journal of English linguistics, 32(4), 346-359.
- Gabrielatos, C. (2018). **Keyness analysis: Nature, metrics and techniques.** InC. Taylor & A. Marchi (Eds.), *Corpus Approaches to Discourse: A Critical Review* (pp. 225–258). Routledge, 10, 9781315179346-11.
- 張莉萍. (2000). 華語教材中的語法用語. 第六屆世界華語文教學研討會. [https://web.ntnu.edu.tw/~lchang/yufayongyu\\_2000.pdf](https://web.ntnu.edu.tw/~lchang/yufayongyu_2000.pdf)
- 朱蘊兒, & 陳百齡. (2016). 眇裡尋他千百度：新聞文本中的關鍵詞篩選與類目建構. 中華傳播學年會.

## Tech Readings

- [7 Case study: comparing Twitter archives | Text Mining with R](#)

- [Who wrote the anti-Trump New York Times op-ed? Using tidytext to find document similarity – Variance Explained](#)
- [Blog | Julia Silge](#)
- [Named Entity Recognition with Spacy and the Mighty roBERTa | by Zoumana Keita](#)
- [Guide to Named Entity Recognition with spaCy and NLTK \(analyticsindiamag.com\)](#)
- [初學者 | 今天掌握SnowNLP好不好 - 每日頭條 \(kknews.cc\)](#)

## W6 Key Influencer: Digital Diplomacy & Cyber Warriors

- [\[Instruction\] Literature Review Guideline](#)

- 關鍵影響人要如何定義？
- 關鍵影響人採用的定義應該是廣域的(比方說整個網絡)，或者僅是局部的(比方說，單一個點的周遭)？
- 網絡指標似乎非常適合用來定義關鍵影響人，Degree, closeness, betweenness centrality, eigencentrality, page rank所定義的關鍵影響人有什麼樣的特性？

## Reference

- Ingenhoff, D., Calamai, G., & Sevin, E. (2021). **Key Influencers in Public Diplomacy 2.0: A Country-Based Social Network Analysis**. *Social Media + Society*, 7(1), 2056305120981053.
- Uren, T., Thomas, E., & Wallis, J. (2019). **Tweeting through the great firewall: preliminary analysis of PRC-linked information operations on the Hong Kong protests**. Australian Strategic Policy Institute.
- [Social network analysis: Understanding centrality measures \(cambridge-intelligence.com\)](#)
  - [Social Network Analysis - Cambridge Intelligence \(cambridge-intelligence.com\)](#)
- [Identifying Influencers on Social Media: A Guide to Social Network Analysis Using Python | by Dina Bavli | Towards Data Science](#)
- [Analyzing Twitter User Network. Analyzing the Twitter user network... | by Mananai Saengsuwan | Towards Data Science](#)

## W8 N-gram and Collocation Network

- TM00 Rundown

Preview: Find and Answer

- Terms和N-gram之間的關係是什麼？
- N-gram和Collocation之間的差異在哪裡？兩者在應用上又會有哪些不同之處？能舉出一個使用N-gram的應用案例和一個使用Collocation的應用案例嗎？何時該使用n-gram何時又該使用Collocation呢？
- Collocation的指標有哪幾種？
- SentencePiece的概念和n-gram異同？用途為何？

Activity: Using Corpus to play Collocation

使用以下兩種不同類型的系統，舉你感興趣的例子各一個，貼上並分享你感興趣的範例。

- [國教院--華語文語料庫應用系統整合入口 \(naer.edu.tw\)](#)
- [Netspeak](#)

- Tutorials of Co-occurrence

- R: [4 Relationships between words: n-grams and correlations | Text Mining with R \(tidytextmining.com\)](#)
  - R: [Analyzing Co-Occurrences and Collocations in R \(slcladal.github.io\)](#)
- Tutorials of Network Visualization
  - [Network visualization with R \(kateto.net\)](#)
  - [Katya Ognyanova: Rutgers Prof., Network Researcher, Data Scientist \(kateto.net\)](#)
  - [Chapter 2 igraph package | Introduction to Network Analysis Using R \(yunranchen.github.io\)](#)
- Sample codes
  - [TM02 Collocation — Programming for Social Scientists \(p4css.github.io\)](#)
  - [TM03 POS Tagging — Programming for Social Scientists \(p4css.github.io\)](#)
    - Python: [TM02\\_collocation.ipynb](#)
  - R: [TM07 Collocation Speech](#), Download data from [toChinaSpeech.rda](#)
  - [Ptt\\_fathering \(Dropbox\)](#)
  - <https://cclab-tm.notion.site/Collocation-0eef6be2850044a68951307004c7f142>
- Slides
  - [TM02 Collocation\\_Py](#) (Actually bigram with distance)
  - [TM02-Collocation\\_R](#)

## Reference

- [4 Relationships between words: n-grams and correlations | Text Mining with R \(tidytextmining.com\)](#)
- Stefanowitsch, A. (2020). **Corpus linguistics: A guide to the methodology.** Language Science Press. <https://library.oapen.org/handle/20.500.12657/43768>
  - Chapter 6 Significant testing
  - **Chapter 7 Collocation**
  - **Chpater 8 Grammar and collocation**
- Manning, C., & Schütze, H. (1999). **Chapter 5. Collocation. Foundations of statistical natural language processing.** MIT press. See <https://paperpile.com/shared/HzGKFN>.
- Chen, A. C.-H. (2022). **Words, constructions and corpora: Network representations of constructional semantics for Mandarin space particles.** *Corpus Linguistics and Linguistic Theory*, 18(2), 209–235. <https://doi.org/10.1515/cllt-2020-0012>

## W9 Topic modeling

- [Talk20230315-Computational Political Science?](#)
- [TM03 Topic modeling with R](#), Coherence and summary part
  - a. <https://www.notion.so/cclab-tm/Topic-modeling-Review-6648ee68fb52420f81528d07e636fb36?pvs=4>
- Sample Code: [RStudio Cloud](#)
  - a. <https://www.dropbox.com/sh/ohq6qzxvpkoidc/AAAU-cikPg80ROtuSYfu3khUa?dl=0>
- Good Tutorials
  - a. [6 Topic modeling | Text Mining with R \(tidytextmining.com\)](#)
  - b. [stm: An R Package for Structural Topic Models \(harvard.edu\)](#) ← Using STM for political science and most of Social Science Applications
    - i. [Structural Topic Modeling with R — Part I | by Jovan Trajceski | Medium](#)
    - ii. [Structural Topic Modeling with R — Part II | by Jovan Trajceski | Medium](#)

- c. [GitHub - MaartenGr/BERTopic: Leveraging BERT and c-TF-IDF to create easily interpretable topics.](#) ← Using BERTTopic in your research if you don't need to include independent factor in your research. With executable Colab examples. Recommended!!
  - i. [About Coherence of topic models · Issue #90 · MaartenGr/BERTopic \(github.com\)](#)
  - ii. [suggestion: incorporating document-level covariates · Issue #360 · MaartenGr/BERTopic \(github.com\)](#)
  - iii. [Model Reproducibility · Issue #987 · MaartenGr/BERTopic \(github.com\)](#)
- d. Top2vec [How to perform topic modeling with Top2Vec. | Towards Data Science](#)  
Without sufficient information.
  - i. [HKML S4E4 - Top2Vec: Distributed Representations of Topics, with application on 2020 10-K - YouTube](#): Video tutorial!!
- e. Classical methods vs. BERTTopic: [Topic Modeling with LSA, pLSA, LDA, NMF · BERTTopic, Top2Vec: a Comparison | by Nicolo Cosimo Albanese | Towards Data Science](#)

### Highlight

1. Implementing Topic modeling by Python/R,
2. Topic modeling for short text documents
3. Explaining topic-document distribution and word-topic distribution correctly
4. How-to evaluate the result of topic, e.g., by coherence score or ...
5. Applications

### Take-home Quiz

- 請簡述Topic modeling之Documents-Topics-Words間的關係。Topic modeling比較適用於哪一種文本？嘗試閱讀論文歸納之。Topic modeling有不同的模型，過去針對不同的文本會如何選用模型？Topic modeling的結果要如何評估？

### Reference

- Egger, R., & Yu, J. (2022). [A topic modeling comparison between LDA, NMF, Top2Vec, and BERTopic to demystify Twitter posts](#). *Frontiers in Sociology*, 7, 886498. <https://doi.org/10.3389/fsoc.2022.886498>
- Maier, D., Waldherr, A., Miltner, P., Wiedemann, G., Niekler, A., Keinert, A., Pfetsch, B., Heyer, G., Reber, U., Häussler, T., Schmid-Petri, H., & Adam, S. (2018). [Applying LDA Topic Modeling in Communication Research: Toward a Valid and Reliable Methodology](#). *Communication Methods and Measures*, 12(2-3), 93–118.
- Lucas, C., Nielsen, R. A., Roberts, M. E., Stewart, B. M., Storer, A., & Tingley, D. (2015). [Computer-Assisted Text Analysis for Comparative Politics](#). *Political Analysis: An Annual Publication of the Methodology Section of the American Political Science Association*, 23(2), 254–277. <https://doi.org/10.1093/pan/mpu019>
- Wang, P. Y. A., & Hsieh, S. K. (2023). Incorporating structural topic modeling into short text analysis. *Concentric*, 49(1), 96-138. [Incorporating structural topic modeling into short text analysis | John Benjamins \(jbe-platform.com\)](#)
- Rauchfleisch, A. (2017). The public sphere as an essentially contested concept: A co-citation analysis of the last 20 years of public sphere research. *Communication and the Public*, 2(1), 3-18.
- Roberts, M. E., Stewart, B. M., Tingley, D., Lucas, C., Leder-Luis, J., Gadarian, S. K., ... & Rand, D. G. (2014). Structural topic models for open-ended survey responses. *American journal of political science*, 58 (4), 1064-1082.
- Roberts, M. E., Stewart, B. M., & Tingley, D. (2019). Stm: An R package for structural topic models. *Journal of Statistical Software*, 91, 1-40.

- Osnabrügge, M., Ash, E., & Morelli, M. (2021). Cross-Domain Topic Classification for Political Texts. *Political Analysis*, 1-22.
  - Abstract. ... \*\*assess the use of supervised learning in cross-domain topic classification. ... classify topics in a labeled source corpus and then predicts topics in an unlabeled target corpus from another domain. .... We demonstrate the method using the case of labeled party platforms (source corpus) and unlabeled parliamentary speeches (target corpus). .... To illustrate the usefulness of the method, we present two case studies on how electoral rules and the gender of parliamentarians influence the choice of speech topics
  - Cross-domain topic classification: Using supervised learning, training labeled party platforms text to predict unlabeled parliamentary speeches
- [When Coherence Score is Good or Bad in Topic Modeling? | Baeldung on Computer Science](#)
- [Evaluate Topic Models: Latent Dirichlet Allocation \(LDA\) | by Shashank Kapadia | Towards Data Science](#)

## W10 Topic modeling Applications: Agenda-setting and Framing?

- Best R Practice: [Training, evaluating, and interpreting topic models | Julia Silge](#)
- Good Python Practice using gensim: [Evaluate Topic Models: Latent Dirichlet Allocation \(LDA\) | by Shashank Kapadia | Towards Data Science](#)
- [2. Topic Modeling: A Naive Example — ENC2045 Computational Linguistics](#) (The coherence computation is implemented in gensim. Tmtoolkit is designed to apply the coherence computation to a sklearn-trained LDA.
- [NLTK :: Sample usage for framenet](#) FrameNet of NLTK
- [FrameGrapher | fndrupal \(berkeley.edu\)](#)
- Applications of Topic modeling:  
<https://cclab-tm.notion.site/Topic-Modeling-Applications-139ae6d7edfb47b990fe272cf2a6159a>

## Reference

- Sturdza, M. D. (2018). **Automated Framing Analysis: A Rule Based System for News Media Text**. Journal of Media Research - Revista de Studii Media, 11(32), 94–110. <https://www.ceeol.com/search/article-detail?id=715200>
- Caswell, D., & Dörr, K. (2018). **Automated Journalism 2.0: Event-driven narratives**. Journalism Practice, 12(4), 477–496. <https://doi.org/10.1080/17512786.2017.1320773>
- 張錦華. (2011). 從 van Dijk 操控論述觀點分析中國大陸省市採購團的新聞置入及報導框架：以台灣四家報紙為例. 中華傳播學刊, 20, 65–93.
- 鄧朝元, 廖達琪, & 黃韋豪. (2022). 國產新冠肺炎疫苗報導之新聞框架:《ETtoday 新聞雲》,《Yahoo! 奇摩新聞》,《中時新聞網》及《自由時報》之 Facebook 粉絲專頁之比較.

## W11 Word embedding

### Resource

- <https://p4css.github.io/jour7088/representation.html>
- [Word-Embeddings \(Notion\)](#)
- [SOMA-Word Embedding](#)
- [Opinion Analysis 2022-05-02-10-22-26](#) (Live Recording)
- [Opinion Analysis 2022-05-09-10-23-35](#) (Live Recording)
- [TM04 Word2vec](#)

- Comparing word2vec to wordfish and wordscore,  
[https://github.com/2048lab/SMA/blob/main/5\\_embedding\\_vs\\_wordfish\\_word\\_score.Rmd](https://github.com/2048lab/SMA/blob/main/5_embedding_vs_wordfish_word_score.Rmd)

## Questions

1. 就課程介紹詞嵌入的方法後，似乎他是一個被用來探查「偏見」或「刻板印象」的好工具。除了種族、國家、或性別外，還有可能有什麼樣的偏見或刻板印象？請構思並舉例說明看看。
2. 固然詞嵌入的方法被大量運用來查找新聞或社群文本中的「偏見」或「刻板印象」，但本身他的概念就是若有兩個詞在句法結構上有相似的特徵的話，那他們兩個詞向量會非常接近。那麼，除了查找「偏見」或「刻板印象」外，除了我們說可以拿來做詞彙擴展外，還可以用來探究什麼樣的議題？

## Reference

1. Garg, N., Schiebinger, L., Jurafsky, D., & Zou, J. (2018). **Word embeddings quantify 100 years of gender and ethnic stereotypes**. Proceedings of the National Academy of Sciences of the United States of America, 115(16), E3635–E3644.  
<https://doi.org/10.1073/pnas.1720347115>
2. Rheault. (2020). **Word Embeddings for the Analysis of Ideological Placement in Parliamentary Corpora**. Political Analysis: An Annual Publication of the Methodology Section of the American Political Science Association.  
<https://doi.org/10.1017/pan.2019.26>
3. Rodman. (2020). A Timely Intervention: Tracking the Changing Meanings of Political Concepts with Word Vectors. Political Analysis: An Annual Publication of the Methodology Section of the American Political Science Association.  
<https://doi.org/10.1017/pan.2019.23>
4. Dyvre, A. (2021). The promise and pitfall of automated text-scaling techniques for the analysis of jurisprudential change. *Artificial Intelligence and Law*, 29(2), 239–269.  
<https://doi.org/10.1007/s10506-020-09274-0>
5. Nissim, M., van Noord, R., & van der Goot, R. (2020). Fair is better than sensational: Man is to doctor as woman is to doctor. Computational Linguistics, 46(2), 487-497.

## W12 Ideology Scaling

- <https://ntucc.webex.com/meet/jerryhsieh>
- Topic over network  
<https://gamma.app/public/SMOALRTopic-over-network-fli2g4qv7pwqkr6>
  - [State-aligned trolling in Iran and the double-edged affordances of Instagram \(notion.so\)](#)
  - [The German Far-right on YouTube: An Analysis of User Overlap and User Comments \(notion.so\)](#)
  - [Fighting Zika with honey: An analysis of YouTube's video recommendations on Brazilian YouTube \(notion.so\)](#)
- Talk20230315-Computational Political Science? → TM06-Ideology-PCA

## Reference

- Bond, R., & Messing, S. (2015). **Quantifying Social Media's Political Space: Estimating Ideology from Publicly Revealed Preferences on Facebook**. *The American Political Science Review*, 109(1), 62–78.  
<https://doi.org/10.1017/S0003055414000525>
  - 張耕齊. (2017). *Ideology Estimation , Media Slant , and Opinion Segregation*

- Barberá, P., Jost, J. T., Nagler, J., Tucker, J. A., & Bonneau, R. (2015). **Tweeting From Left to Right: Is Online Political Communication More Than an Echo Chamber?** *Psychological Science*, 26(10), 1531–1542.  
<https://doi.org/10.1177/0956797615594620>
- Bakshy, E., Messing, S., & Adamic, L. A. (2015). **Exposure to ideologically diverse news and opinion on Facebook.** *Science*, 348(6239), 1130–1132.  
<https://doi.org/10.1126/science.aaa1160>
- Dyevre, A. (2021). The promise and pitfall of automated text-scaling techniques for the analysis of jurisprudential change. *Artificial Intelligence and Law*, 29(2), 239–269.  
<https://doi.org/10.1007/s10506-020-09274-0>
- Kosinski, M., Stillwell, D., & Graepel, T. (2013). Private traits and attributes are predictable from digital records of human behavior. *Proceedings of the national academy of sciences*, 110(15), 5802–5805.

## W13 Sentiment Analysis & Stance-based sentiment analysis

- What is sentiment? And what is stance?

-  TM09 Sentiment Analysis - Stance-based

### Reference

- Cao, R., Luo, X., Xi, Y., & Qiao, Y. (2022). **Stance detection for online public opinion awareness: An overview.** *International Journal of Intelligent Systems*, 37(12), 11944–11965. <https://doi.org/10.1002/int.23071>
- Bestvater, S. E., & Monroe, B. L. (2022). **Sentiment is Not Stance: Target-Aware Opinion Classification for Political Text Analysis.** *Political Analysis: An Annual Publication of the Methodology Section of the American Political Science Association*, 1–22. <https://doi.org/10.1017/pan.2022.10>
- Liu, B. (2007). Web Data Mining. Exploring Hyperlinks, Contents, and Usage Data. Heidelberg: Springer-Verlag Berlin (不用看)

### CODING GUILDLINE

- (Good) [Beginner's Guide to Sentiment Analysis for Simplified Chinese using SnowNLP | by Ng Wai Foong | Towards Data Science](#)

## W14 Document Classification

- Submit to here:  
<https://www.notion.so/cclab-tm/SMOA1112-a763f1d9fa854927bbd4961bcd1a5eec?pvs=4>
- <https://gamma.app/docs/Using-OpenAI-ChatGPT-for-sentiment-analysis-and-stance-detection-rx8d2z6uzf6qqj5>
- <https://deepnote.com/workspace/cclab-4442-7cab8c2d-d96f-4b95-b47b-bbd0a02160a2/project/SMOA-6e31c6ed-05e6-4a51-9714-660efc847b14/notebook/Using%20ChatGPT-4b7a2514c2814af8bbf3eb0cd21ec241>

- Project link:  
<https://deepnote.com/workspace/cclab-4442-7cab8c2d-d96f-4b95-b47b-bbd0a02160a2/project/SMOA-6e31c6ed-05e6-4a51-9714-660efc847b14>
- [Sentiment Analysis with ChatGPT, OpenAI and Python — Use ChatGPT to build a sentiment analysis AI system for your business | by Courtlin Holt-Nguyen | Data And Beyond | Apr, 2023 | Medium](#)
- 
- 

### Short Assignment: Embedding+Classification

這個練習是想測試，如果缺少標記過後的情緒資料，有沒有可能單靠embedding就完成分類任務？

請你找約100個有「答案」的情緒分類資料(可以上Kaggle找、Google找、或你也可以自己編)，通常情緒分類資料會有Positive、Negative、Neutral三種標籤。希望你找的這100筆資料是上述資料集中，至少40%是Positive、40%是Negative的資料集。

然後用sbert或simpletransformer，直接建立正向與負向的embedding後，測試每一個句子(最好是一個Tweet)和正負向文字的距離來決定該文章或該句子的正負向。

- [Text Representation Examples - Simple Transformers](#)
- [Computing Sentence Embeddings — Sentence-Transformers documentation \(sbert.net\)](#)

### Reference

- Grimmer, J., & Stewart, B. M. (2013). **Text as data: The promise and pitfalls of automatic content analysis methods for political texts.** Political analysis, 21(3), 267-297.
- Barberá, P., Boydston, A. E., Linn, S., McMahon, R., & Nagler, J. (2021). **Automated Text Classification of News Articles: A Practical Guide.** Political Analysis: An Annual Publication of the Methodology Section of the American Political Science Association, 29(1), 19–42. <https://doi.org/10.1017/pan.2020.8>
- Di Cocco, J., & Monechi, B. (2021). **How Populist are Parties? Measuring Degrees of Populism in Party Manifestos Using Supervised Machine Learning.** Political Analysis: An Annual Publication of the Methodology Section of the American Political Science Association, 1–17. <https://doi.org/10.1017/pan.2021.29>

## W15 Concluding

- <https://cclab-tm.notion.site/ChatGPT-Testing-8257ceba06a44c91acbffd6b1d667e64>
- Talk20230315-Computational Political Science?
- TM07 Document classification
- TM08 Classification
- TM08 Transformer for doc classification and clustering
  - Simpletransformer on stance detection.  
[https://colab.research.google.com/drive/1zGg\\_EeExlKuTTCIIK5T-qAhE9--4KleQ?usp=sharing](https://colab.research.google.com/drive/1zGg_EeExlKuTTCIIK5T-qAhE9--4KleQ?usp=sharing)
  - Simpletransformer on aspect-based opinion analysis.  
<https://colab.research.google.com/drive/1YD8mjrpEH6-0rYMY0L1xFFKpJ2jGY6nT?usp=sharing>
- 戴淨妍(2023)新冠記者會中的疫情政治學

- <https://cclabtw.notion.site/Automated-Framing-Analysis-A-Rule-Based-System-for-New-s-Media-Text-128692ee863242deb7f122403e24fb2b>
- 

## Reference

- Di Cocco, J., & Monechi, B. (2021). **How Populist are Parties? Measuring Degrees of Populism in Party Manifestos Using Supervised Machine Learning.** *Political Analysis: An Annual Publication of the Methodology Section of the American Political Science Association*, 1–17. <https://doi.org/10.1017/pan.2021.29>
- Baden, C., Pipal, C., Schoonvelde, M., & van der Velden, M. A. C. G. (2021). **Three Gaps in Computational Text Analysis Methods for Social Sciences: A Research Agenda.** *Communication Methods and Measures*, 1–18.
  - → Lucas, C., Nielsen, R. A., Roberts, M. E., Stewart, B. M., Storer, A., & Tingley, D. (2015). Computer-assisted text analysis for comparative politics. *Political Analysis*, 23(2), 254-277. ...STM's team, just for promoting STM, closer to comparative politics, far away from communication. I don't like it.

## Week 17

-  TM09 Topic Modeling in Python
- <https://colab.research.google.com/drive/14GDXgdGUOos2Jr82ly0Y4OZOcE0vfSQZ?usp=sharing>

## Assignments

### AS: Topic modelings

針對短文本做Topic Modeling, 取16個topics和128個topics。針對128個topics抽取關鍵字做Clustering, 比較Clustering和16 topics的結果。用PMI做Collocation, 並透過網絡分析方法觀察關鍵字群集, 觀察PMI做Collocation和Clustering和16 Topics的結果有何差異。

針對以下Topic-modeling進行方法的比較與測試。

1. R: 測試BTM [bnosac/BTM: Biterm Topic Modelling for Short Text with R \(github.com\)](#)
2. R: 測試STM並使用LDAVis來觀察結果
3. Python: 測試biterm
4. Py: 測試並比較shorttext和LDA的效果
5. Py: 測試Gensim、Scikit-learn的效果
6. Py: 測試Scikit-learn和BERTopic的效果

### 相關套件

- R: topicModels: LDA
  - [NLP: Use Topic Modeling Results in Predictive Modeling - The Analytics Lab](#) with very good visualization
- R: [bnosac/BTM: Biterm Topic Modelling for Short Text with R \(github.com\)](#)
- [biterm · PyPI](#)
- [shorttext · PyPI](#)
- [rwalk/gsdmm: GSDMM: Short text clustering \(github.com\). Example01](#)

- [Short-Text Topic Modelling: LDA vs GSDMM | by Richard Pelgrim | Towards Data Science](#)
- [Short Text Topic Modeling. Intuition and \(some\) maths to... | by Matyas Amrouche | Towards Data Science](#) (Implement by hand)
- [Gensim Topic Modeling - A Guide to Building Best LDA models \(machinelearningplus.com\)](#)
  - [Frontiers | Using Topic Modeling Methods for Short-Text Data: A Comparative Analysis | Artificial Intelligence \(frontiersin.org\)](#)
- [BERTopic Topic Modeling Using The BERTopic Library | Towards Data Science](#)
- [NLP Tutorial: Topic Modeling in Python with BerTopic | HackerNoon](#)

## AS: Collocation

Collocation and POS, word embeddings, and topic modeling (2 of 3)

## AS: Embeddings

雖然我們常說「偏見」和「刻板印象」廣存一般的新聞文本和社群輿論中，但這些存有偏見或刻板印象的文章、段落、詞句，可能會有什麼樣的特徵？例如，當帶有刻板印象／偏見的句子，和其他文章相較下，在用詞上可能會有什麼差異？在你的文本中，哪些文章或句子會有刻板印象／偏見呢？你選擇哪一種單位來分析呢？(文章、段落、句子)

- Word embeddings and text classification (demo)
- AS2: (1) Building reserved vocabulary by wikipedia DB with crawler and HMM; (2) Verifying tokenization efficacy and accuracy; (3) Building personal word2vec CBOW and Skip-gram.

## Essay

### Pitching

每個人僅需報告5~10分鐘說明自己的研究主題。

- 研究題目：簡報上，第一頁為你的主題想三個題目，然後Highlight你覺得想得最好的題目或方向。可以用這種方法思考你的研究標題[How to Write a Research Paper Title with Examples - Wordvice](#)。
- 研究主題：承上，加一頁，用點列的方式，清楚地闡述研究的主題及其背景。
- 研究動機：頁三，解釋為何這個研究題目值得研究？為什麼是你來做、可能潛在的貢獻是什麼？用點列回答這三個問題看看？
- 研究問題：頁四，目前有點難，但應該嘗試明確地寫三則研究問題，有助於聚焦和修整研究主題。因為研究問題很容易看得出你會做出什麼東西來，包含哪幾個面向，缺什麼。
- 樣本及資料蒐集：頁四，用一句話寫出你想要搜集什麼樣的資料，時間大概多長，搜尋以查找資料的關鍵字為何？解釋研究所選擇的樣本及其選擇的理由，以及資料蒐集的方式、方法和工具等。
- 文獻：頁五，簡列你所查找的相關文獻，盡可能是近三、五年內重要期刊的論文，可以有一兩篇碩論，以知道國內類似的議題做到哪裡，和國內有什麼相關學者。

### 暫時不用附

- 研究目的：說明研究的目的及其意義，明確提出研究所要達成的目標。

- 文獻回顧：(暫時不用附)梳理現有的相關文獻，解釋研究對象的相關理論、研究方法等。
- 研究方法：說明研究所採用的方法，例如問卷調查、實驗、案例研究等。
- 其中一頁，簡列你所查找的相關文獻，盡可能是近三、五年內重要期刊的論文，可以有一兩篇碩論，以知道國內類似的議題做到哪裡，和國內有什麼相關學者。
- 建議用Google Slides可以直接網路上存取的簡報。書面通常比較沒有結構性，不利於報告。

## Literature Reviews & Corpus

### **Literature Review**

#### **Corpus**

我們曾提及了如何定義適合作為研究對象的Corpus，界定良好範圍的Corpus有助於衍生好的研究問題，請用一個段落的書寫來回答下列問題：

1. 你的研究題目為何？
2. 你獲取資料的方式是用關鍵字搜尋或依賴於既定的資料範圍？
3. 你目前所獲得(或預計獲得)的資料範圍為何？媒體？平台？原始資料總筆數？時間範圍為何？
4. 你的資料裡面可能會有哪些是不相關的資料(Garbage data)？

找一篇近五年國外期刊的論文作為參考(要有APA)，複製貼上其對資料集／語料庫的相關描述，作為你的參考範本，並在其下撰寫你的研究題目的資料集／語料庫描述。

例文：本研究的研究議題為....，為了了解....，我(預期)透過...的方法來....；或我查找XX報標題..內文為....；或我透過...方法撈取XX版的所有XX，時間範圍設定在XX至XX，設定在這個範圍的原因因為....。搜尋的的文本資料如...(檢列一兩則標題...)。但資料中可能包含XX或XX或XX等情形，所以我設定一個篩除字組{XX、XX、XX、...}，並篩除所有字組的相關文本，(預期)共獲得XX筆資料。

## Exploratory data analysis (EDA) on personal project

1. Data volume and variables + irrelevant data report and tunrcation
2. Cleaned data coverage and justification (why the collected data is good enough for your research questions)
3. Submitting rmarkdown + rendered html or ipynb + rendered html

## Methods

## Final Essay and Presentation

## Take-home Quiz

這些問題是社群輿論探勘經常會被問到的問題，希望修課學生能夠有一套自己的想法，並且在遇到這樣的問題的時候從善如流地回答。請用論文論述的方式(找引證、提出觀點)來回答以下問題。

1. Echo-Chamber、Filter Bubble、Polarization的概念有何差異？操作型定義為何？試舉例說明之。
2. 何謂Aspect (Stance)-based sentiment analysis？其實作上的困難點在哪裏？
3. Word embeddings是一套什麼樣的技術？近年常被用來做什麼類型研究？
4. 請簡述Topic modeling之Documents-Topics-Words間的關係。Topic modeling比較適用於哪一種文本？嘗試閱讀論文歸納之。Topic modeling有不同的模型，過去針對不同的文本會如何選用模型？Topic modeling的結果要如何評估？

## 線上自學資源

### Python

只要先看到Py06即可以順利運用課堂上的教學內容。

Tutorials	Videos
<ul style="list-style-type: none"><li>❑ PyBook-Basic : Ch1</li><li>❑ Py01 Counting</li><li>❑ <a href="#">P01_counting.ipynb</a></li></ul>	<ul style="list-style-type: none"><li>● <a href="#">Py01_1 Counting</a></li><li>● <a href="#">Py01_2 Counting more</a></li><li>● <a href="#">Py01_3 Counting words</a></li></ul>
<ul style="list-style-type: none"><li>❑ PyBook-Basic : Ch2</li><li>❑ Py02_03 variables, list, dict</li><li>❑ <a href="#">P02-Python Basic</a></li><li>❑ <a href="#">P03-List and Dictionary</a></li></ul>	<ul style="list-style-type: none"><li>● <a href="#">Py02 Basic</a></li><li>● <a href="#">Py03_1 List</a></li><li>● <a href="#">Py03_2 Dictionary</a></li><li>● <a href="#">Py03_3 Accessing ubike data</a></li></ul>
<ul style="list-style-type: none"><li>❑ Py04 flow control - for, if</li><li>❑ <a href="#">P04-For-if-else</a></li></ul>	<ul style="list-style-type: none"><li>● <a href="#">Py04 for01</a></li><li>● <a href="#">Py04 for02 Ibike sum up</a></li><li>● <a href="#">Py04 for03 histogram</a></li><li>● <a href="#">Py04 for04 fib pi 9x9</a></li></ul>
<ul style="list-style-type: none"><li>❑ Py04 flow control - for, if</li><li>❑ <a href="#">P04 For-if applications</a></li></ul>	<ul style="list-style-type: none"><li>● <a href="#">P04_2_1 Traversing AQI data</a></li><li>● <a href="#">P04_2_2 Iterating data: 3 methods</a></li><li>● <a href="#">P04_2_3 Searching and Hashing</a></li><li>● <a href="#">P04_2_4 Rescale01 by nested if-elif-else</a></li><li>● <a href="#">P04_2_4 Rescale02 by list-dict</a></li><li>● <a href="#">P04_2_5 Finding Prominant</a></li></ul>
<ul style="list-style-type: none"><li>❑ Py04 flow control - for, if</li><li>❑ <a href="#">P04 for-if applications</a></li></ul>	<ul style="list-style-type: none"><li>● <a href="#">Py04_1 if for Concept and AQI case introduction</a></li><li>● <a href="#">Py04_2 if for Accessing AQI data</a></li><li>● <a href="#">Py04_3 if for Categorizing and rescaling</a></li><li>● <a href="#">Py04_5 Sorting manually by for, if, swap, list and dict</a></li></ul>
<ul style="list-style-type: none"><li>❑ Py06 Pandas and data su...</li></ul>	<ul style="list-style-type: none"><li>● <a href="#">Py06 pandas01</a></li><li>● <a href="#">Py06 pandas02</a></li></ul>

	<ul style="list-style-type: none"> <li><a href="#">Py06 Pandas03 create new variables</a></li> <li><a href="#">Py06 Pandas04 groupby sort</a></li> <li><a href="#">Py06 Pandas05 series to df</a></li> <li><a href="#">Py06 AS Pandas how to</a></li> </ul>
Visualization	<ul style="list-style-type: none"> <li><a href="#">P07_2 Bokeh Seaborn (loom.com)</a></li> </ul>
<input checked="" type="checkbox"/> Py07 Crawler	<ul style="list-style-type: none"> <li><a href="#">Py08_1 Web request and response</a></li> <li><a href="#">Py08_2 ChromeDevTool</a></li> <li><a href="#">P08_3 tweepy (loom.com)</a></li> <li><a href="#">Py08_3 Crawler 104 - YouTube</a></li> </ul>
<a href="#">TM01 Tokenization</a> <a href="#">Chinese Processing</a>	<ul style="list-style-type: none"> <li><a href="#">TM00 Overview</a></li> <li><a href="#">TM01_1 Term frequency tokenize</a></li> <li><a href="#">TM01_2 Term frequency</a></li> <li><a href="#">TM01_3 Stemming lemmatization</a></li> <li><a href="#">TM04 Chinese segmentation</a></li> </ul>

## R

如果要複習R的話，建議先看到**III Typd Data**結束後，跳到**V. Text Mining**看完Trump的案例即可。

<b>I. R Basic</b>			
1	<a href="#">R01_1 loading_data.Rmd</a> <a href="#">R01_2 vector.Rmd</a> <a href="#">R01_3 dataframe.Rmd</a>	<input checked="" type="checkbox"/> R01_2 basic	<ul style="list-style-type: none"> <li><a href="#">R01_2 Storing data</a></li> <li><a href="#">R01_3 vector</a></li> <li><a href="#">R01_4 dataframe</a></li> </ul>
<b>II. Reading Data</b>			
2	<a href="#">R02_1p readxl paid maternal leave.Rmd</a> <small>讀取Excel: 產假支薪(gitbook)</small>	<input checked="" type="checkbox"/> R02_1 Paid maternal L...	<ul style="list-style-type: none"> <li><a href="#">R02_1_0 Intro_paid_maternal_leave</a></li> <li><a href="#">R02_1_1 Read Excel</a></li> <li><a href="#">R02_1_2: Selecting, Filtering and Cleaning data</a></li> <li><a href="#">R02_1_3 : Plotting</a></li> </ul>
3	<a href="#">R02_2p read_csv pivot_on_tptheft.Rmd</a> (final ver.)	<input checked="" type="checkbox"/> R02_2 Read CSV and ...	<ul style="list-style-type: none"> <li><a href="#">R02_2_1 Introduction to CSV and the TP theft case</a></li> <li><a href="#">R02_2_2 Reading CSV</a></li> <li><a href="#">R02_2_3 Pivot Analysis on TP theft</a></li> <li><a href="#">R02_2_4 mosaicplot() on crosstable</a></li> </ul>
<b>III. Tidy Data</b>			
4	<a href="#">R03_1 base to dplyr maternal leave.Rmd</a> <a href="#">R03_2p base to dplyr_tptheft.Rmd</a>	<a href="#">R04_1 Data Manipulation dplyr and ggplot</a>	<ul style="list-style-type: none"> <li><a href="#">R02_3 based to dplyr01 Intro. to dplyr</a></li> <li><a href="#">R02_3 base to dplyr02: using pipeline</a></li> <li><a href="#">R02_3 base to dplyr03: on TP house theft</a></li> <li><a href="#">R02_3 base to dplyr04: Paid maternity leave</a></li> </ul>
5	Open new .Rmd by yourself	<a href="#">R04_3 Join data MOI -</a>	<ul style="list-style-type: none"> <li><a href="#">R04_1_1 Loading MOI demographic data</a></li> <li><a href="#">R04_1_2 Summarizing demographic data</a></li> <li><a href="#">R04_1_3 Aggregating village to town level stat</a></li> <li><a href="#">R04_1_4 Joining demographic and referendum</a></li> </ul>
7	<a href="#">R04_3p viz ggplot.Rmd</a>		<ul style="list-style-type: none"> <li><a href="#">R04_ggplot01_basic</a></li> <li><a href="#">R04_ggplot02_plot paras</a></li> <li><a href="#">R04_ggplot03 chinese coorflip</a></li> </ul>

			<ul style="list-style-type: none"> <li><a href="#">R04 ggplot04 highlight</a></li> </ul>
<b>V. Text Mining</b>			
9	<a href="#">R05_2p_trump_tweet_dplyr.Rmd</a>	<input type="checkbox"/> R05_2 Text data manip...	<ul style="list-style-type: none"> <li><a href="#">4.1 Tweet analysis on trump tweets (Case Introduction)</a></li> <li><a href="#">4.3 Tweet analysis doc level (for RE &amp; ggplot)</a></li> </ul>
13	<a href="#">R05_2p_trump_tweet_dplyr.Rmd</a> <a href="#">R05_3p_tm_typhoon.Rmd</a> <a href="#">R05_4p_President_Speech.Rmd</a>	<input type="checkbox"/> R05_2 Text data manip... <input type="checkbox"/> R05_3 Chinese proces...	<ul style="list-style-type: none"> <li><a href="#">4.4.1 Tweet analysis word level I</a></li> <li><a href="#">4.4.2 Tweet analysis word level II</a></li> <li><a href="#">4.5 Text analysis for Chinese text</a></li> <li><a href="#">4.6 Chinese Information Retrieval</a></li> </ul>
15	<a href="#">R05_1_regular_expression.Rmd</a> <a href="#">R05_1p_regular_expression.Rmd</a>		<ul style="list-style-type: none"> <li><a href="#">R04_2_1 Regular Expression</a></li> <li><a href="#">R04_2_2 RE Extract</a></li> <li><a href="#">R04_2_3 RE applications</a></li> </ul>

## Deprecated 1102

### Student Report

- 高若如:  0530\_\_社群媒介輿論分析報告\_\_高若如
- 李密:  Social Media\_Gender Differences in Using Language
- 陳柔安:  News Diff\_related work
- Otherness & Pronoun\_Papar review\_0516 by EvelynYang & CCSun
- Keyness report: [Work Summary Plan \(canva.com\)](#) by BZHsieh

### Practice

**Readings** - 從以下PaperPile的連結中，在Thesis Folder中找到以下三篇論文並下載 <https://paperpile.com/shared/HzGKFN>。這個Prep任務是為了快速瀏覽文本探勘、內容分析相關研究的研究設計。這三本都是相當高分的碩論，但除卻議題本身，他們在理論與方法的運用與結合上，大相逕庭，和傳統方法的綜合運用和比較上、就方法而言、和理論對話上，你覺得他們有什麼樣的差異？請針對各篇文獻簡要用「100」字做個摘要。

- 江攷. (2021). Populism and Political Communication on Facebook\_ 2020 Taiwan Presidential Election. [https://doi.org/10.1016/S1541-4612\(21\)00035-5](https://doi.org/10.1016/S1541-4612(21)00035-5)
- 黃子暉. (2019). 從媒體微博之限制評論設置看中國網路輿論限縮 [國立臺灣大學]. <https://doi.org/10.6342/NTU201901458>
- 齊若堯. (n.d.). Media representation and bias against Southeast Asia: Analysis of Taiwan's media coverage of Filipinos.

繳交方式：

- 直接打在[Draft]的文件裡，以標題一將「Week 4 Prep - Comparing research designs」作為標題，然後在標題一下方開始書寫。
- 無需另外繳交(但下週上課前老師會先看看每個人寫的東西，務必準時完成)