

Show
me the
data!

Week01: Introduction

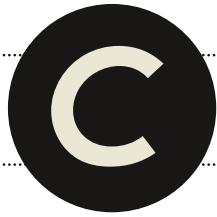
Big Data & Social Analysis R

Instructors: Chung-pei Pien

ZU1942001/266868001/Z23937001/ZM1941001



International College of
INNOVATION
National Chengchi University
國立政治大學創新國際學院



CONTENTS

- 1 Introduction
- 2 Evaluation
- 3 Data Project
- 4 R and R-Studio

Introduction

01

Introduction

The Level of This Course

Domain Knowledge

Sociology

Economy

Politics

Data Science

Statistics

This Course

Basic R Review

Coding
Application

Project-driven
Process

Other GTIM

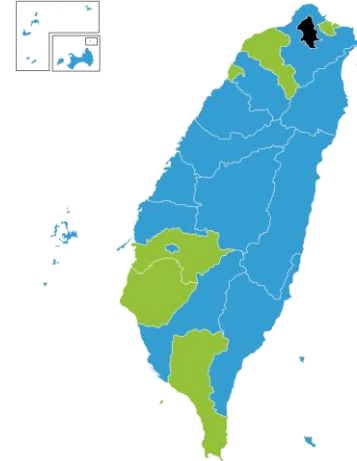
Capstone

Data Science

Business Data
Analysis

I will negotiate with two media for the capstone course in Fall 2022.

2022
MIDTERM



There are two crucial elections in this November: US midterm election and TW major election.

The participants of the capstone course can produce these two elections' data journalism and publish them in these media.

English Media: Taiwan Major Election

Chinese Media: US Midterm Election

To prepare the capstone course, I will use my previous election data projects to teach you coding in this course.

01

Introduction

The Structure of This Course

Election Data and R

Data
Exploration

Regular
Expressions

Loop

Data
Cleaning

Web Crawler

Data
Visualization



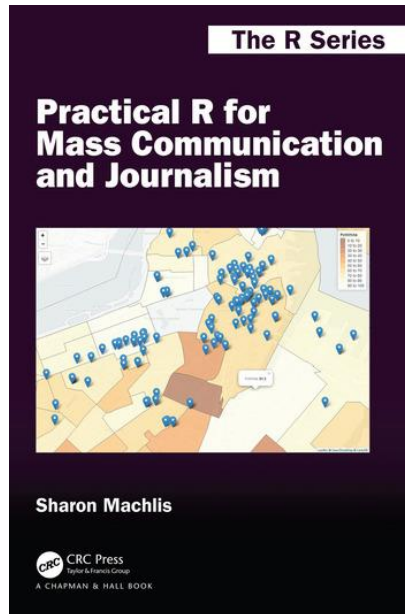
Advanced Data Projects

Taiwan Presidential
Election Data Project

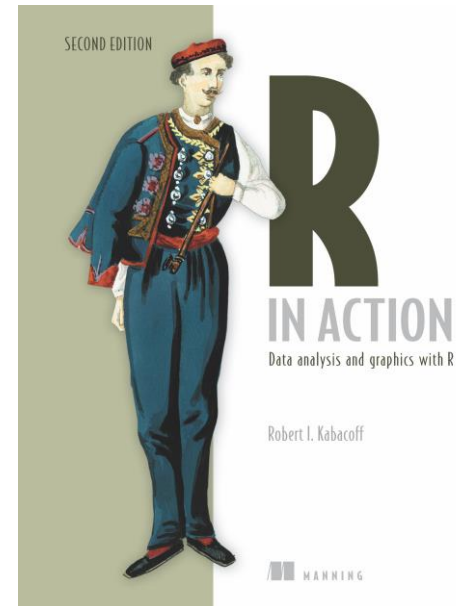
US Midterm Election
Data Project

Taiwan Mayor
Election Data Project

Twitter and US
Presidential Election



Machlis, Sharon.
2019. Practical R for
Mass
Communication
and Journalism.
CRC Pres



Kabacoff, Robert.
2015. R in Action,
Data Analysis and
Graphics with R.
Manning
Publications; 2
Edition

Reading

9

Evaluation

02

Evaluation

| Items | Points | % |
|--------------------|---------------|--------|
| Assignments | 10 x 30 = 300 | 23.08% |
| To-do List | 100 | 7.69% |
| Midterm Exam | 200 | 15.38% |
| Project Proposal | 150 | 11.54% |
| Final Presentation | 200 | 15.38% |
| Final Paper | 350 | 26.92% |
| Total | 1300 | 100% |

| Items | Points | % |
|-------------|----------------------|--------|
| Assignments | $10 \times 30 = 300$ | 23.08% |

- 10 weeks have assignments. I will upload them (R script file) to Moodle.
- Print and submit them in the beginning of the next weeks' classes

```
← →  [Icons] Source on Save  [Icons]  → F
1  #Class: Week 01
2  #Course: Big Data and Social Analysis
3  #Lesson: Introduction
4  #Instructor: Chung-pei Pien
5  #Organization: ICI, NCCU
6
7  ▾ ### Student Information -----
8
9  #Name: Bill Chen
10 #ID: 1234567
11 #E-mail: bill_chen@nccu.edu.tw
12
13 ▾ ### Questions -----
14
15 #Question 1: (10 points)
16 #Please calculate 1 + 1 in R
17
18 1 + 1
19
20 ##Calculate one plus one for this homework
21
22 #Question 2: (15 points)
23 #Create an object X that involves "home" and "ICI"
24 |
25 X <- c("home", "ICI")
26
27 ##This is a list that includes two characters: home and ICI
24:1  # Questions ▾
```

| Items | Points | % |
|-------------|----------------------|--------|
| Assignments | $10 \times 30 = 300$ | 23.08% |

- **No comments no points**
- To meet the ICI dean's attendance request, assignments will link to your attendance:
Students who do not ask for leave to skip the class can't submit assignments.

| Items | Points | % |
|------------|--------|-------|
| To-do List | 100 | 7.69% |

- **In week 3, you should decide your teammates and email me your team members' lists.**
- I will provide you a to-do list guideline.
- Your team can follow my suggestions to design your data projects.
- **In Week 6, your team should submit your to-do lists.**
I will return them in Week 7.

02

Evaluation

| Items | Points | % |
|------------|--------|-------|
| To-do List | 100 | 7.69% |

| | Team Points | PM bonus | Writers bonus |
|----|-------------|----------|---------------|
| A+ | 95 | 3 | 2 |
| A | 92 | 3 | 2 |
| A- | 90 | 2 | 2 |
| B+ | 85 | 2 | 1 |
| B | 82 | 2 | 1 |
| B- | 80 | 1 | 1 |
| C+ | 75 | 0 | 0 |
| C | 70 | 0 | 0 |

| Items | Points | % |
|--------------|--------|--------|
| Midterm Exam | 200 | 15.38% |

- Midterm exam will be held on April 21 (week 10)
- Open books and reference cards
- Only can open R-studio or R-studio cloud on your laptop

| Items | Points | % |
|------------------|--------|--------|
| Project Proposal | 150 | 11.54% |

- Every team will be required to present your project proposal in 5-10 slides in the Week 13 for the final presentation and report.
- I will hand out a project proposal guideline in Week 7. You should follow the guideline to discuss your project with your teammates and me.

| Items | Points | % |
|------------------|--------|--------|
| Project Proposal | 150 | 11.54% |

| | Team Points | PM bonus | Writers bonus |
|----|-------------|----------|---------------|
| A+ | 143 | 5 | 3 |
| A | 138 | 4 | 3 |
| A- | 135 | 3 | 3 |
| B+ | 127 | 3 | 1.5 |
| B | 123 | 3 | 1.5 |
| B- | 120 | 1.5 | 1.5 |
| C+ | 112 | 0 | 0 |
| C | 105 | 0 | 0 |

| Items | Points | % |
|--------------------|--------|--------|
| Final Presentation | 200 | 15.38% |

- Week 18, this course and other GTIM courses have a joint presentation party on June 15 (Wed.) 12:00pm.
- Every team has a booth (a table and a poster frames) and use a poster or a laptop's ppt to present your final project.
- The location may be at the 3F lobby of International building.

| Items | Points | % |
|--------------------|--------|--------|
| Final Presentation | 200 | 15.38% |

| | Team Points | PM bonus | Writers bonus |
|----|-------------|----------|---------------|
| A+ | 190 | 6 | 4 |
| A | 184 | 6 | 4 |
| A- | 180 | 4 | 4 |
| B+ | 170 | 4 | 2 |
| B | 164 | 4 | 2 |
| B- | 160 | 2 | 2 |
| C+ | 150 | 0 | 0 |
| C | 140 | 0 | 0 |

| Items | Points | % |
|-------------|--------|--------|
| Final Paper | 350 | 26.92% |

- The final reports should be submitted to my mailbox or emailed before the deadline.
 1. A4 10 pages (not include reference)
 2. Times New Roman 12pt
 3. Double spaces
 4. Margin 1 inch

| Items | Points | % |
|-------------|--------|--------|
| Final Paper | 350 | 26.92% |

| | Team Points | PM bonus | Writers bonus |
|----|-------------|----------|---------------|
| A+ | 332.5 | 9 | 6 |
| A | 322 | 9 | 6 |
| A- | 315 | 6 | 6 |
| B+ | 297.5 | 6 | 3 |
| B | 287 | 6 | 3 |
| B- | 280 | 3 | 3 |
| C+ | 262.5 | 0 | 0 |
| C | 245 | 0 | 0 |

Data Project



Deconstructing: The Sexiest Job of the 21st Century

Deconstructing and Patil 2012

Data analysis process: Six independent but related steps

1. Identify Goals and
Plans



2.Data Collection



3.Data Cleaning

4.Data Analysis

5.Data Presentation

6.Writing Report



3

Data Project

1. Data cleaning spends 80% project time: tired and frustrated
2. There is no data or model to achieve your goals
3. Examining your budgets and ability
4. Team work



Vicki Boykis
@vboykis

跟隨

Have been extremely curious about this for a while now, so I decided to create a poll.
"As someone titled 'data scientist' in 2019, I spend most of (60%+) my time:"
("Other") also welcome, add it in the replies.

6% Picking features/models

67% Cleaning data/Moving data

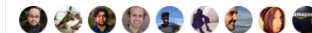
4% Deploying models in prod

23% Analyzing/presenting data

2,116 票 • 最終結果

上午8:17 - 2019年1月28日

118 次轉推 212 個喜歡



38

118

212

Data Collection

Data cleaning

Data analysis

Data Presentation

Writing
report



3

Data Project

1. Domain Knowledge: Social issues
2. Coding: R (this courses), Python, and etc
3. Models: Data science and statistics courses
4. Graphics: Art designers/adobe illustrators
5. Presentation: Writing

Data Collection



Data cleaning



Data analysis



Data Presentation



Writing
report

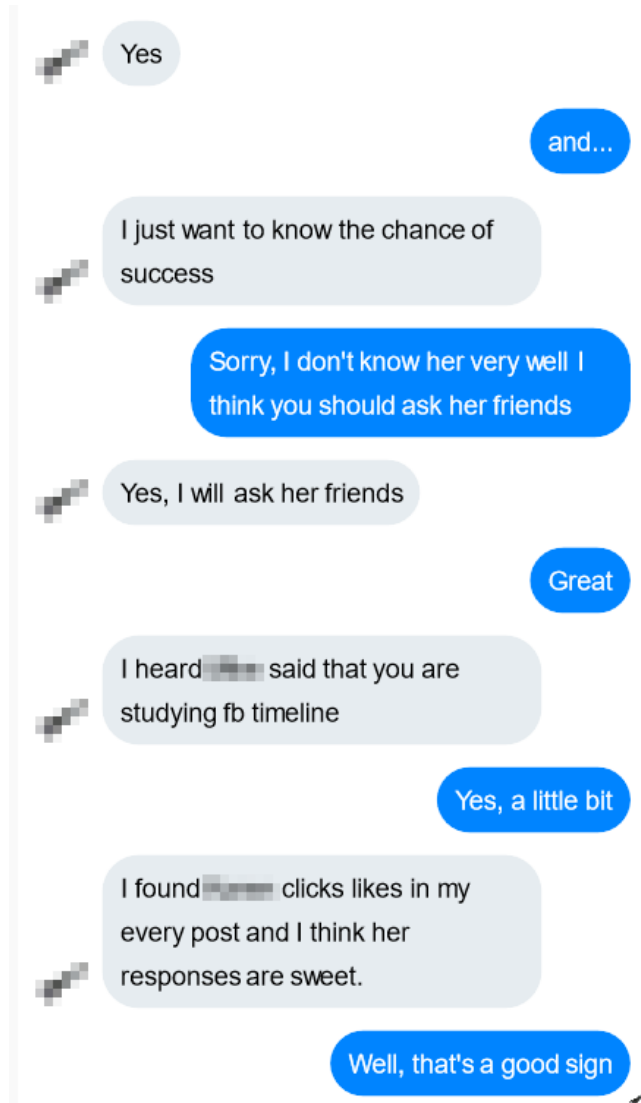
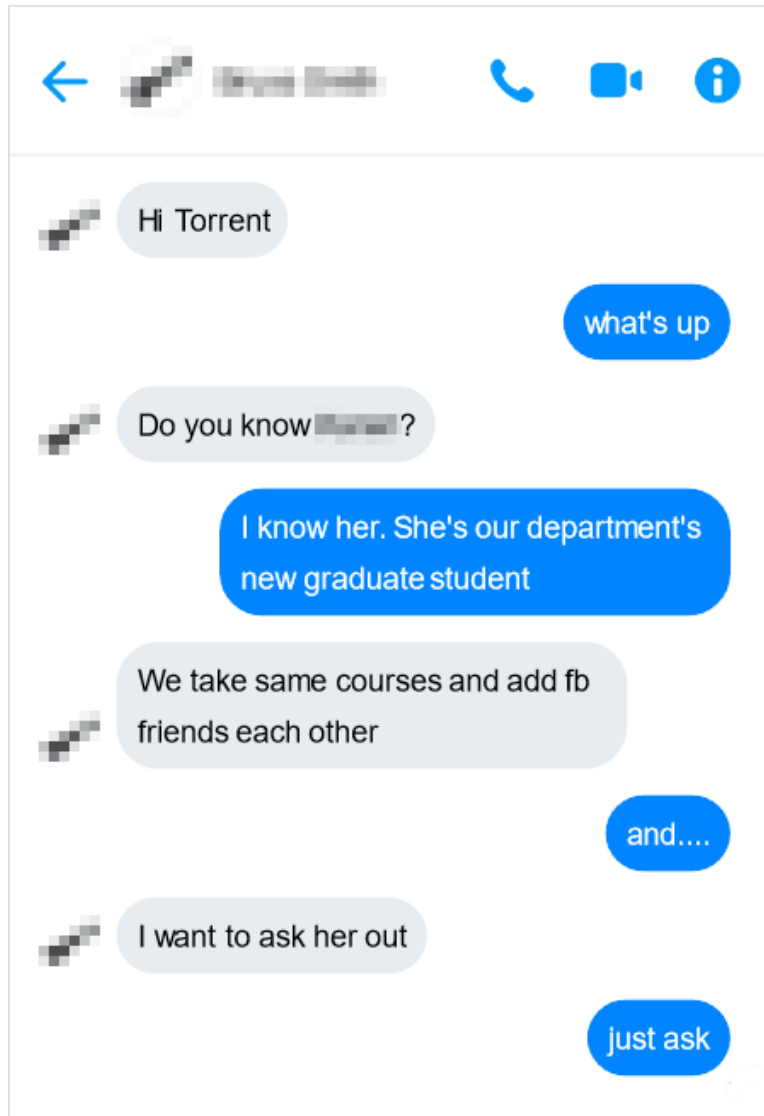


One day, a friend of mine raised a request in Facebook messenger.....

3

Data Project

Facebook Analysis of User Activities



3

Data Project

Facebook Analysis of User Activities



Do you think his demand – ensuring a girl likes him – is technically possible?

3

Data Project

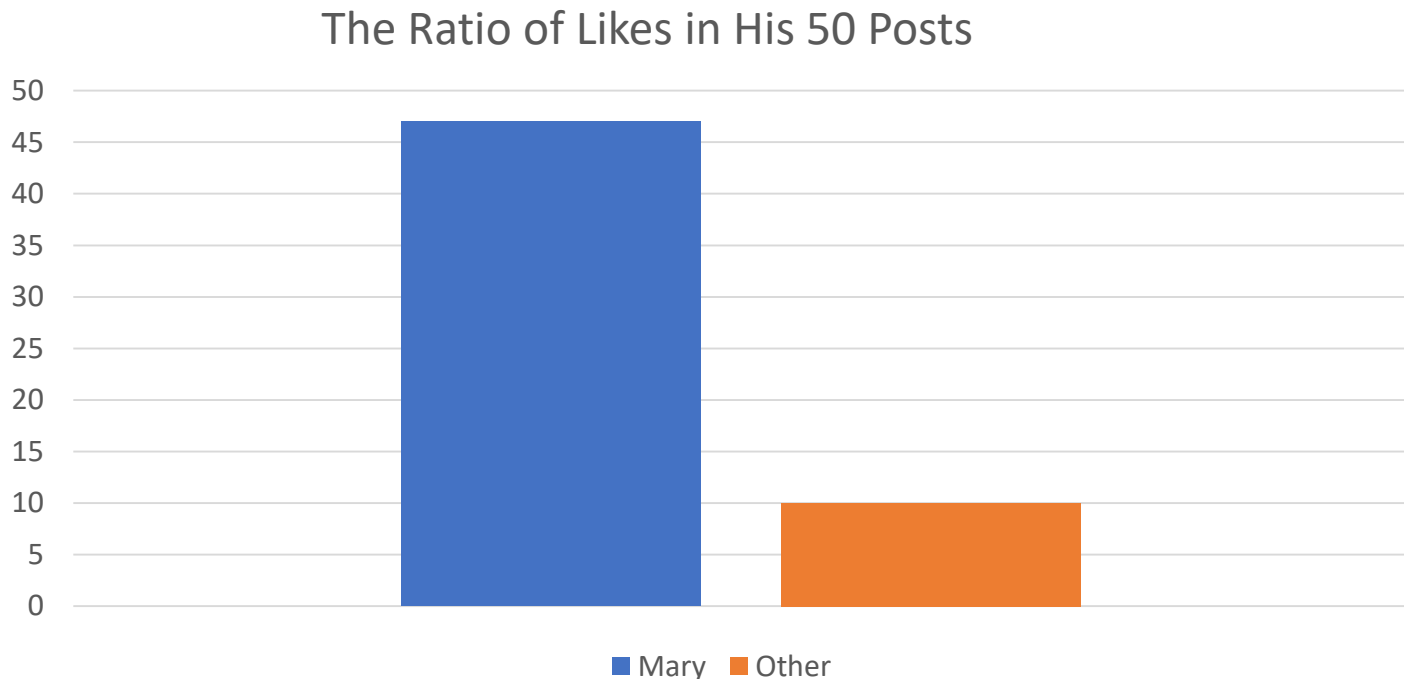
Facebook Analysis of User Activities

- Using FB data to ensure that the girl likes him.



Mary's activities are different than his other friends.

- The ratio of her likes is larger than other friends, significantly.



Mary's activities are different than his other friends.

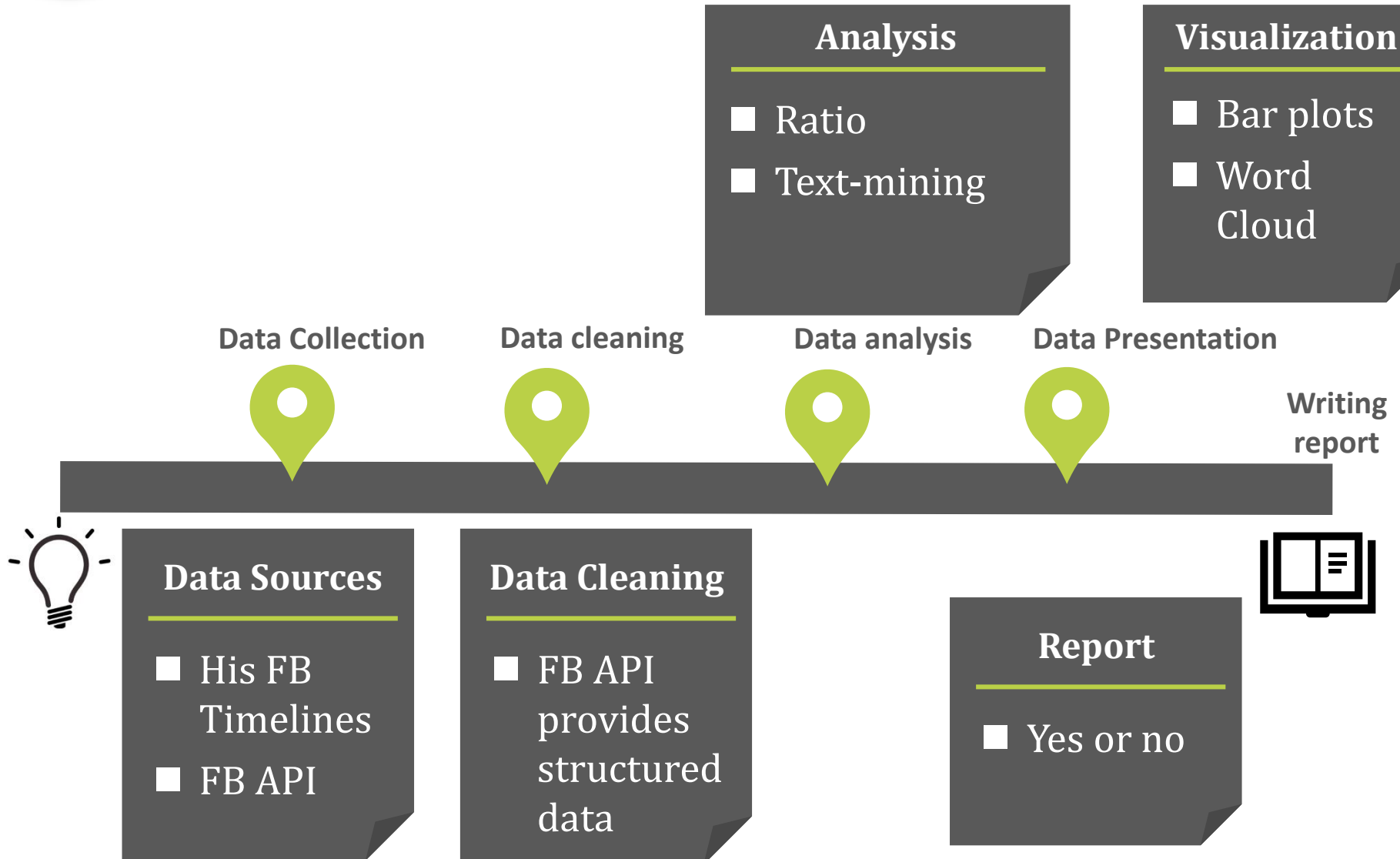
- Comparing other Facebook friends, her responses' key words are sweet.



3

Data Project

Facebook Analysis of User Activities

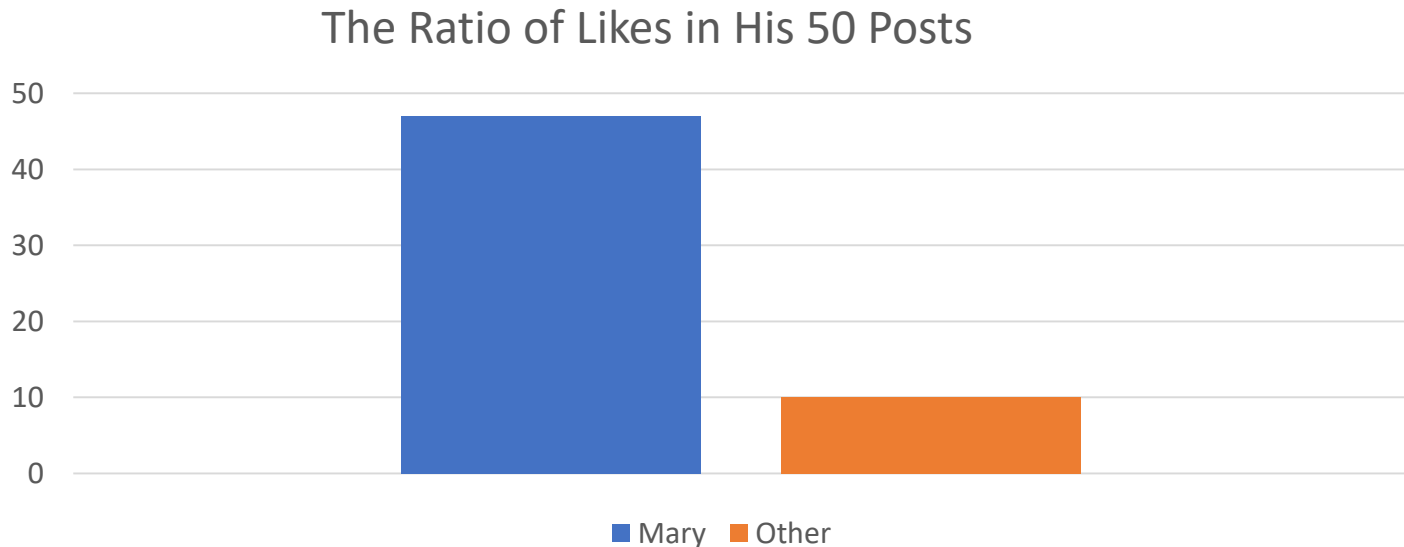


3

Data Project

Facebook Analysis of User Activities

The ratio of Mary's likes is larger than other friends, significantly.



- Maybe she also has a high clicking like rate in her friends' timelines.
- We can't (or should not) get her friends' timelines.

Comparing other Facebook friends, her responses' key words are sweet.



- There are online dictionaries to check a word's positive or negative.
- Positive words do not mean "LIKE."
- You can do your dictionary, but you need a lot of data to do machine learning.

In 2018, after Trump's 2 years' presidency, Democrats might flip the control of congress in the US midterm election.

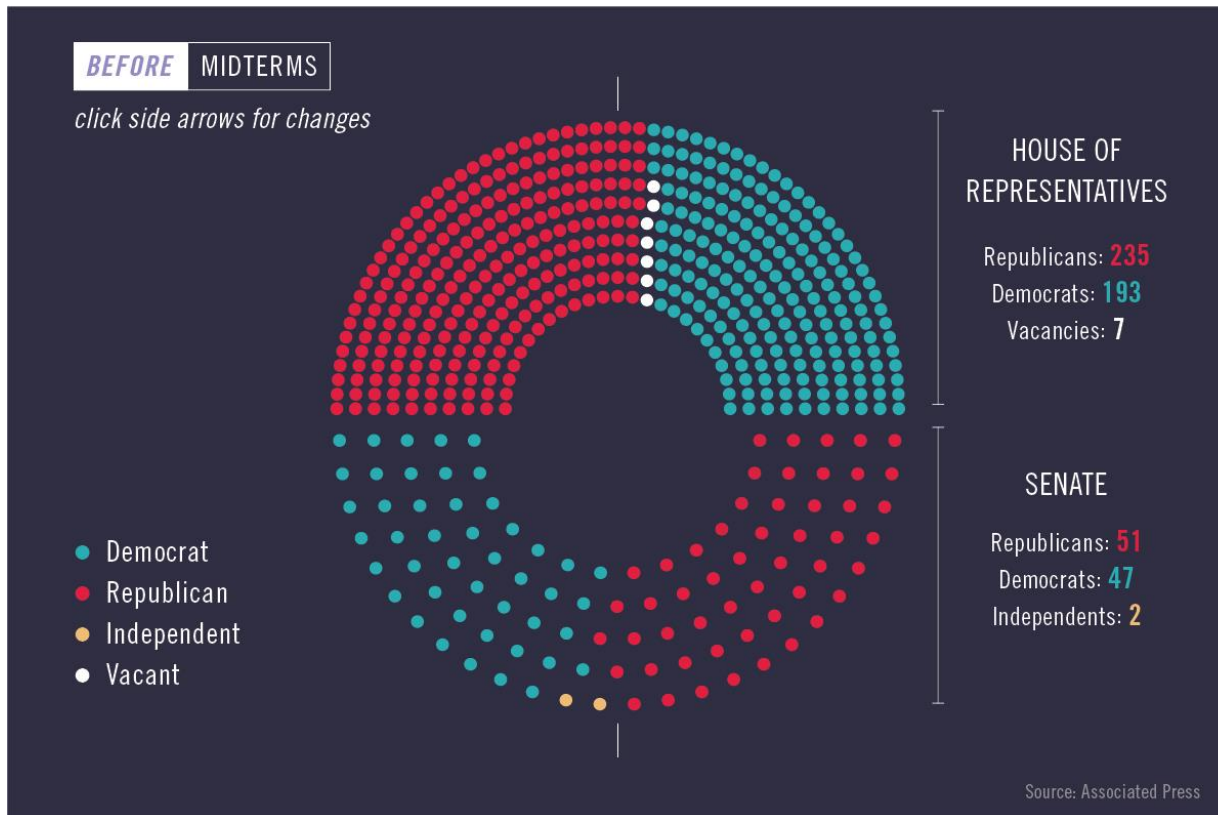
**Pro-Taiwan lawmakers would be replaced?
New congress would support Taiwan?**

3

Data Project

US Congressmen and Taiwan

POLITICAL PARTY

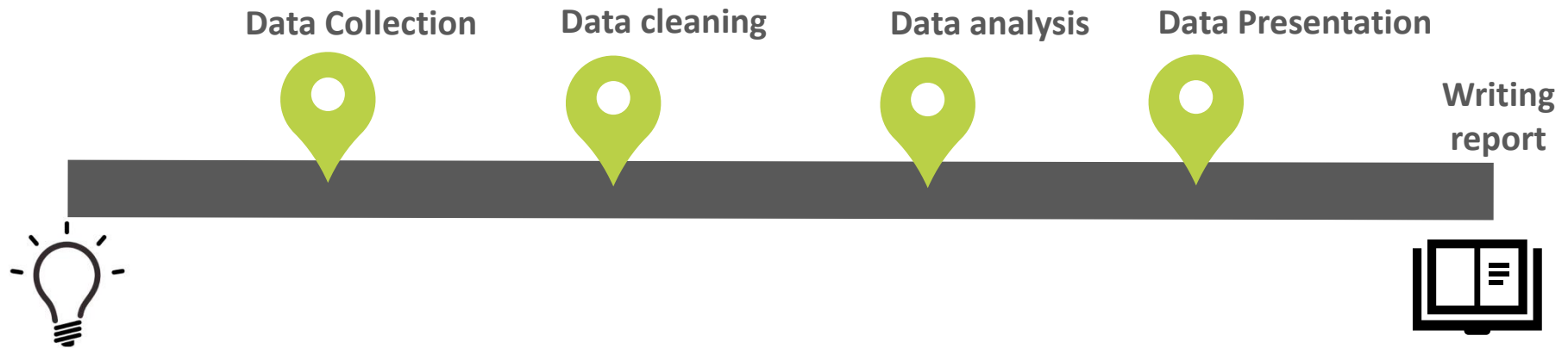


2018

House: 435

Senate: 100

- All US lawmakers' attitude toward Taiwan before and after the midterm election



The primary difficulty of this project is data:

What kinds of data can represent a lawmaker to support Taiwan or not?????

News stories?

Twitters?

Congressional records?

3

Data Project

US Congressmen and Taiwan

| | News Stories | Twitters | Congressional Records |
|-------------------------------------|---|--|--|
| Easy to get? | Is it possible to get all lawmakers' Taiwan opinions in media? | Does all lawmakers have Twitter account. Twitter's | US congress website has csv file or API? |
| Easy to clean? | News is unstructured data. How to evaluate pro-Taiwan or anti-Taiwan? | How to evaluate pro-Taiwan or anti-Taiwan | How to evaluate pro-Taiwan or anti-Taiwan? |
| The data can be counted or modeled? | It's very difficult to evaluate the level like or dislike TW | It's very difficult to evaluate the level like or dislike TW | Pro-Taiwan acts |
| The result can be explained? | Difficult | Difficult | Much better |

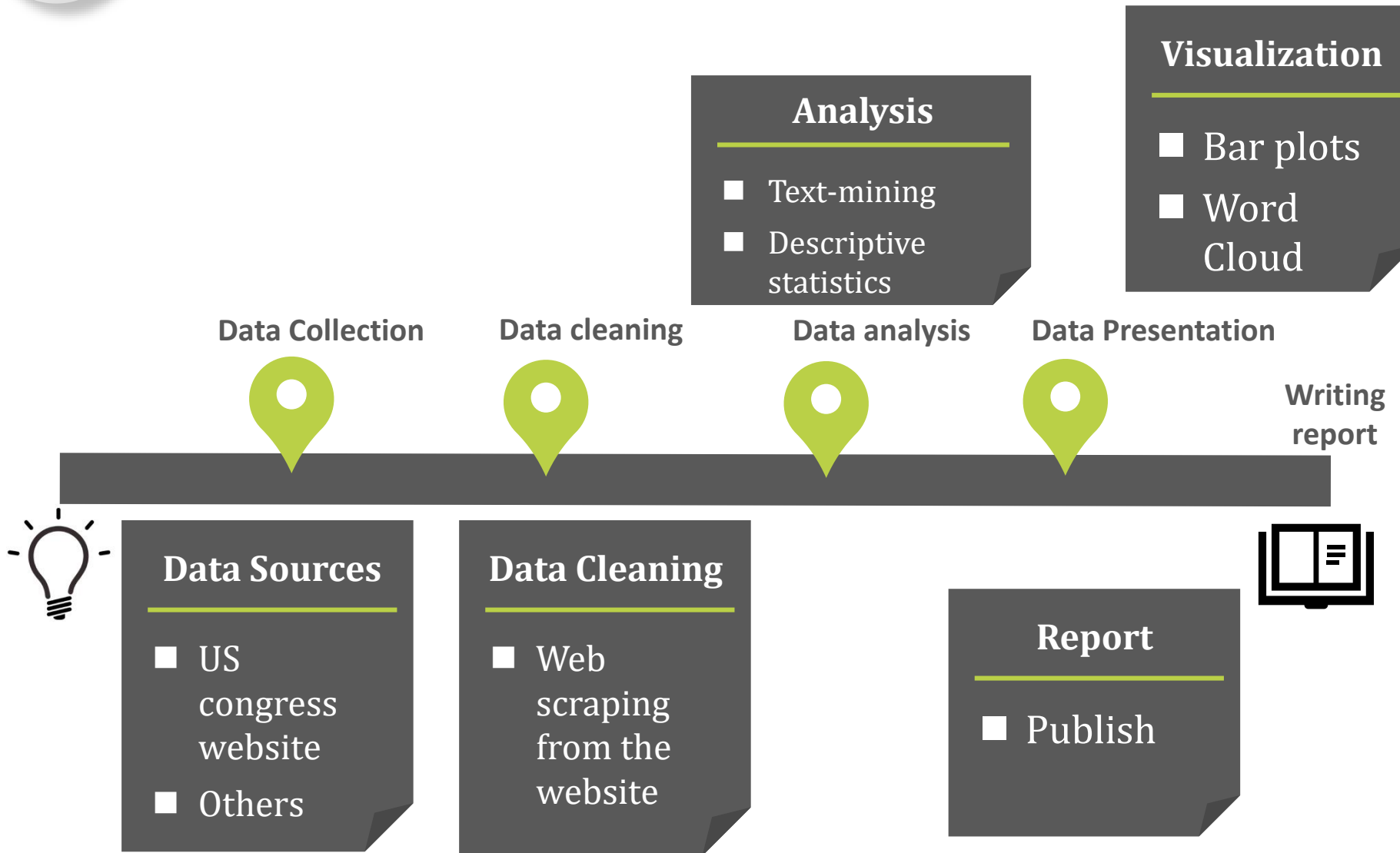
I decided to use congressional records:

1. Pro-Taiwan acts: sponsors and voters
2. Statements
3. Lobby

3

Data Project

US Congressmen and Taiwan



首頁 追蹤

端傳媒 邀請好友 搜索 通知 用戶

最新 國際 大陸 香港 台灣 評論 科技 風物 圓桌 廣場 | 2022端陪你過年 2021年終專展

台灣 深度 美國中期選舉

數據帶你看：美國期中選舉，親台議員誰主浮沉？

擷開美國國會第113屆、第114屆及115屆的法案提案、議員公開新聞稿、國會記錄三組重要數據，我們試圖比較川普任內的第115屆國會與過去第113屆、第114兩屆歐巴馬任內的國會，在對台相關議題的立場和實際行動，呈現何種消長趨勢。



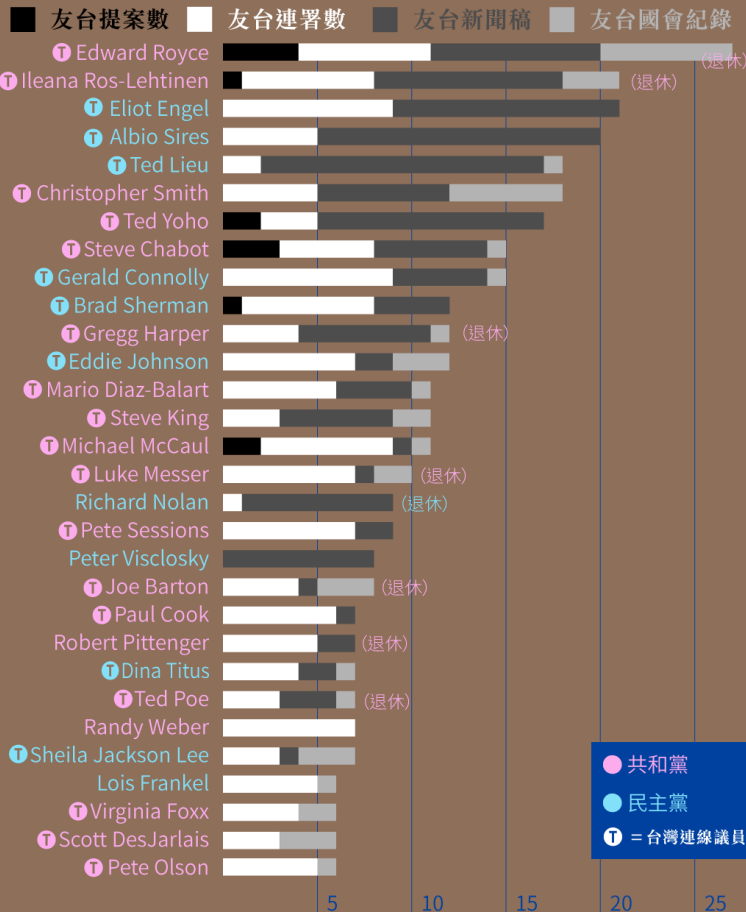


3

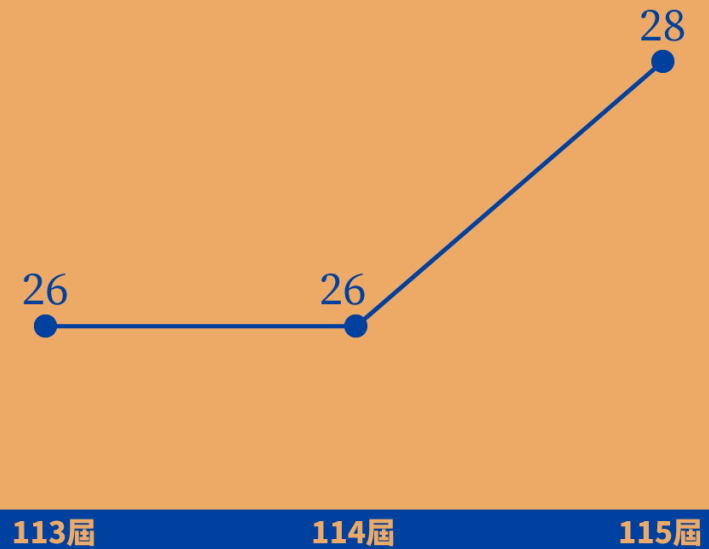
Data Project

US Congressmen and Taiwan

美國親台議員相關統計（眾議院）



近五年來，美國國會的台灣
相關提案數增加還是減少了？

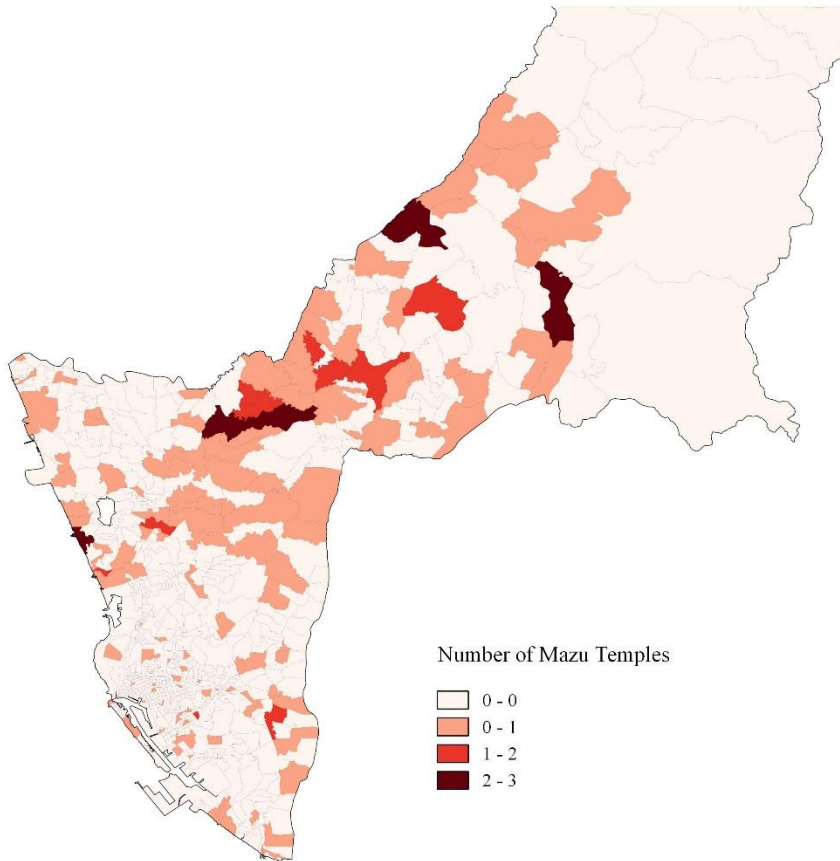


In 2018, KMT candidate Han Kuo-yu (韓國瑜) won the Kaohsiung mayor election.

There are many arguments to explain his win.

One important argument is:

Many Mazu temples organized pilgrimages to the birthplace of Mazu. The Chinese government could influence their political attitudes and they tended to elect KMT.



Can we use data analysis to prove this hypothesis?

The Kaohsiung citizens who are influenced by Mazu temples are more likely to elect Han.

Chien-yuan Sher, Chung-pei Pien, and Yu-hsi Liu have attempted to examine this hypothesis and write a research paper form 2018.

- KH li which have more Mazu temples voted for Han?

Dependent variable: Voting rate for KMT

Independent variable: The number of Mazu temple

Analysis

- Aggregate multi-level nested logit model

Visualization

- GIS
- Regression Tables

Data Collection

Data cleaning

Data analysis

Data Presentation

Writing report



Data Sources

- Ministry of the Interior
- Central Election Commission

Data Cleaning

- 6 Months to merge data

Report

- Research Paper
- Data Journalism




[首頁](#) [追蹤](#)

端傳媒

[邀請好友](#)

[最新](#) [國際](#) [大陸](#) [香港](#) [台灣](#) [評論](#) [科技](#) [風物](#) [圖桌](#) [廣場](#) | [2022端陪你過年](#) [2022北京冬奧會](#)

台灣 深度 2020台灣大選

2020台灣大選

神明不投票：「宮廟影響選舉」的可能與不可能

隨著韓流興起，宮廟力量也在台灣政治討論中被神秘化。被境外勢力滲透？變信徒為選票？直接左右選舉？現實比傳說複雜許多。《端傳媒》深入走訪高雄多家宮廟，結合三位學者2014年與2018年兩次九合一選舉的數據調研，挖掘宮廟與選舉之間的真實關係。



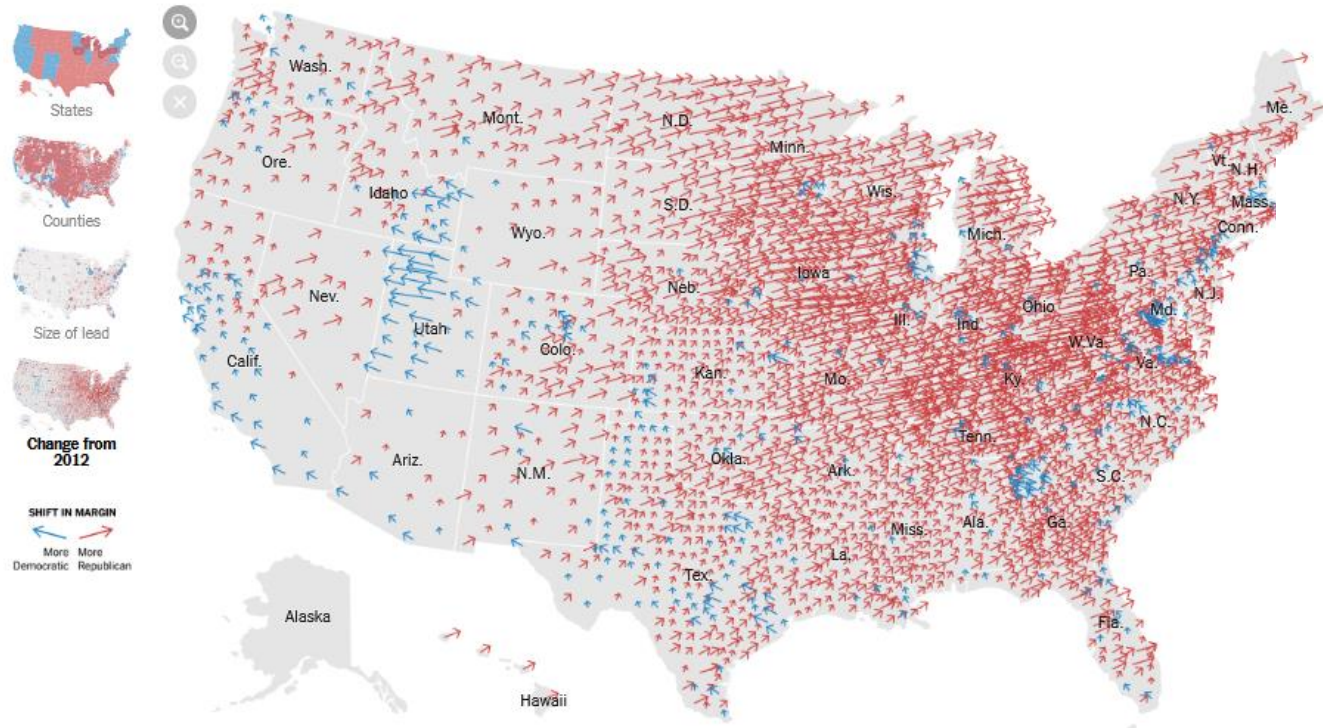
1. You can do election data projects to prepare the capstone course
2. Or you can choose any topics for your final project.

Exercise: Learn from classic projects

2016 Presidential Election Results

AUG. 9, 2017, 9:00 AM ET

In 2016, Donald J. Trump won [the Electoral College](#) with 304 votes compared to 227 votes for Hillary Clinton. Seven electors voted for someone other than their party's candidate. Visit our [2020 election results pages](#) for the latest updates.



3

Data Project

Classic Projects

- Think about how to imitate the same plot:



Analysis

Visualization

Data Collection

Data cleaning

Data analysis

Data Presentation

Writing report

Data Sources

Data Cleaning

Report



Do some ~~google~~ research and answer the following questions in Moodle:

1. Project aim:
2. Data sources and levels:
3. The skills or codes to clean the data in R:
4. The way to calculate to change from 2012-2016:
5. The way to create a similar plot in R:

R

Any problems to install R and R-Studio?

4

R and R-Studio

I may create R-Studio Cloud for this course in the next week.

