# Week 01
# *R* You Ready?

POLI3148. Data Science in Politics and Public Administration

Dr. Haohan Chen          HKU-PPA

# Today

- About Me
- About the course (syllabus)
- *R* Software Setup

# About Me

# About Me



Dr. CHEN Haohan  陳昊瀚

Interested in …

- Political communication
- Computational methods

# My Research

Chinese context

- Understanding the ideological spectrum of Chinese intellectuals using cultural product reviews
- Using audio of news program to infer power dynamic and policy agenda

US context

- Perception on Covid-19 among the US public using Twitter data
- Monitoring political polarization in the US using Twitter data

# How We Work Together Going Forward

- I would love to get to know each one of you in person.
- Language: English, Mandarin, Cantonese
- How we communicate
  - **Office hours**: Calendly appointment system
  - **CampusWire** Course Forum
- Appointments outside office hours possible. Email me.
- If I do not reply to your email within two days, kindly send me a nudge.

# About **DaSPPA**

# To discuss

- Topics
- Readings
- Output and Assessment

**Topics**

**_R_ + Data**
**Model**
**Text Mining**

| Week | Lecture (1st half) | Lecture (2nd half) | Due |
|---|---|---|---|
| 2 | Welcome | _R_ you ready? | |
| 3 | _R_ Basics (1) | _R_ Basics (2) | |
| 4 | _R_ Basics (3) | Data Wrangling (1) | |
| 5 | Data Wrangling (2) | Machine Learning Overview | |
| 6 | Data Visualization (1) | Linear Regression | |
| 7 | Data Visualization (2) | Classification | |
| 8 | _Reading week. No class._ | | |
| 9 | Data Visualization (3) | Resampling Methods | A1 |
| 10 | Data Visualization (4) | Model Selection and Regularization | B1 |
| 11 | Text Mining (1) | Tree-Based Methods | |
| 12 | Text Mining (2) | Unsupervised Learning | |
| 13 | Text Mining (3) | Text Mining (4) | A2 |
| 14 | Putting Everything Together | Debriefing and Q&A | |
| R | _DaSPPA Festival!_ | | G1 |
| A | _Group Final Project Replication Dossier Due_ | | G2 |
| A | _Personal DaSPPA Portfolio Submission Due_ | | A3 |

# R + Data: R Basics

- Basic knowledge of R and Rstudio
- Assume no prior experience with R
- But if you have prior experience, it will be a good review
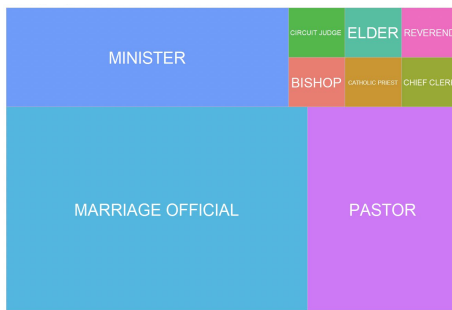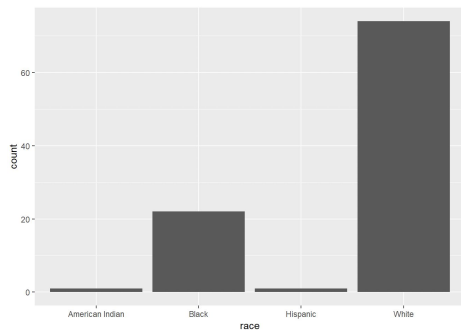
# R + Data: Data Wrangling

Reshape



Subset rows and columns

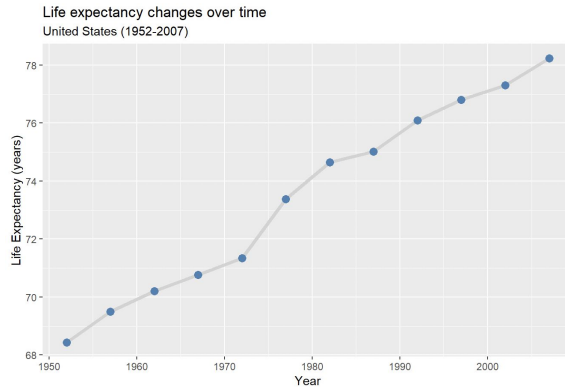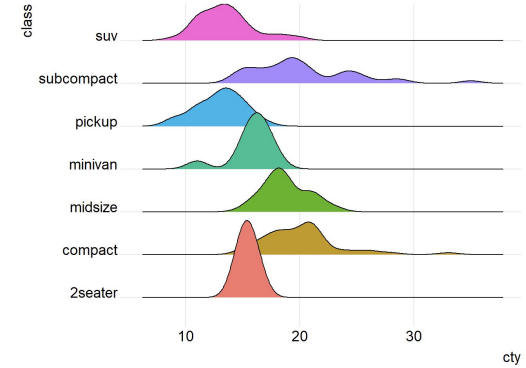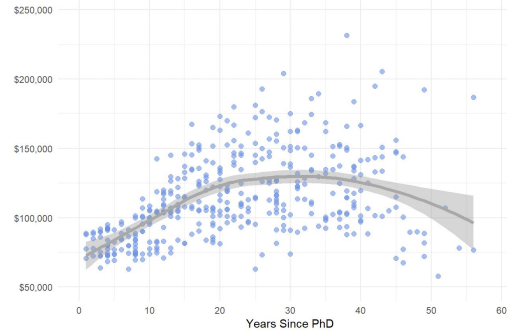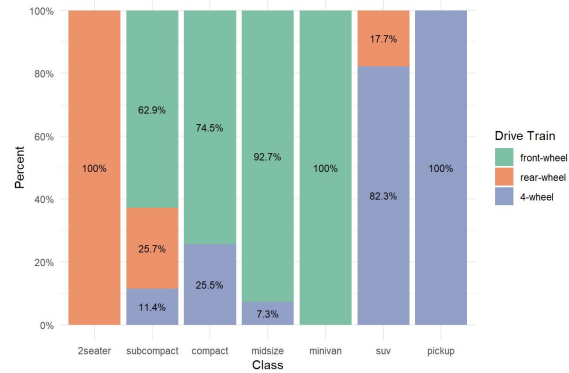

Summarize and mutate variables



Combine datasets

# R + Data: Visualize one variable

# R + Data: Visualize two variables
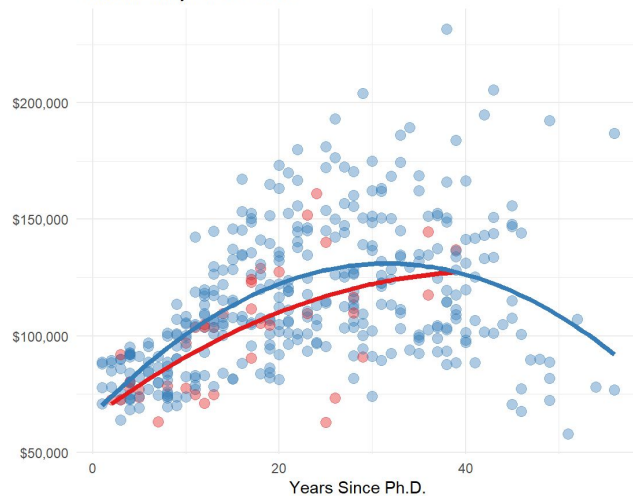
# R + Data: Visualize multiple variables

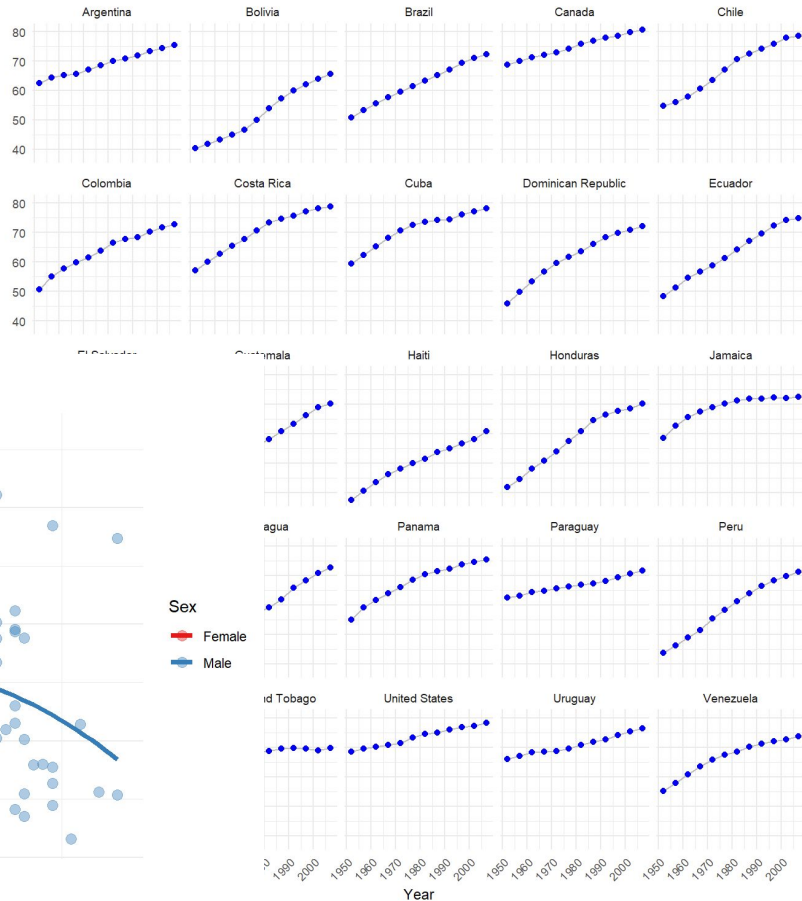# R + Data: Visualize data wrt space and time



Life expectancy by country
Gapminder 2007 data

Years
[39.6 to 49.3)
[49.3 to 55.3)
[55.3 to 62.7)

US Population by age
1900 to 2002

source: https://www.gap

Age Group
>64
55-64
45-54
35-44
25-34
15-24
5-14
<5

source: U.S. Census Bureau, 2003, HS-3

# R + Data: Other cool visualization
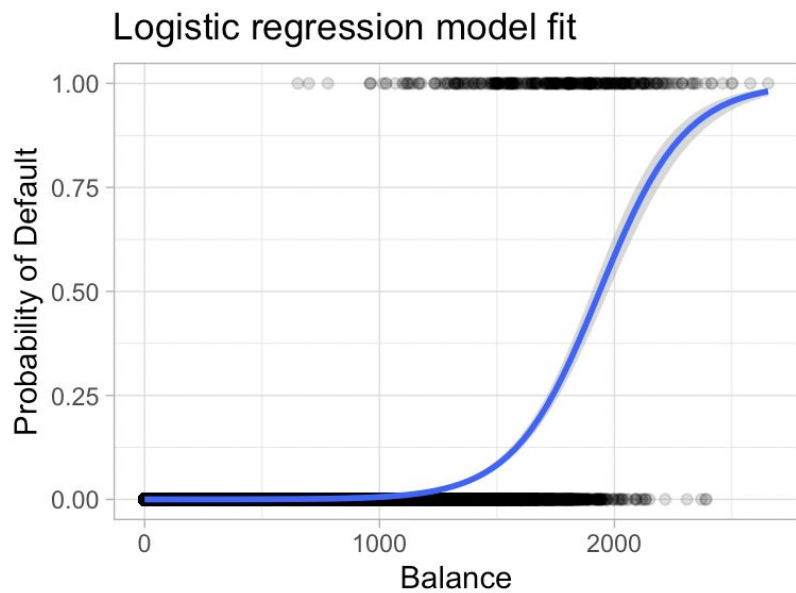


Life Expectancy in Asia

Mammal size and sleep characteristics

# Model: Linear Regression



Fitted regression line (with residuals)

# Model: Classification
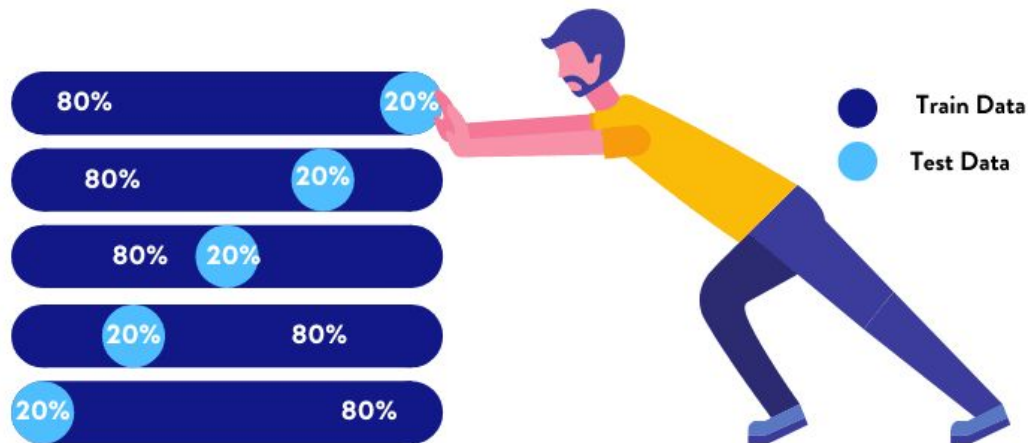


Logistic regression model fit

# Model: Model Selection and Regularization

- Criteria to evaluate how "good" a machine learning model is
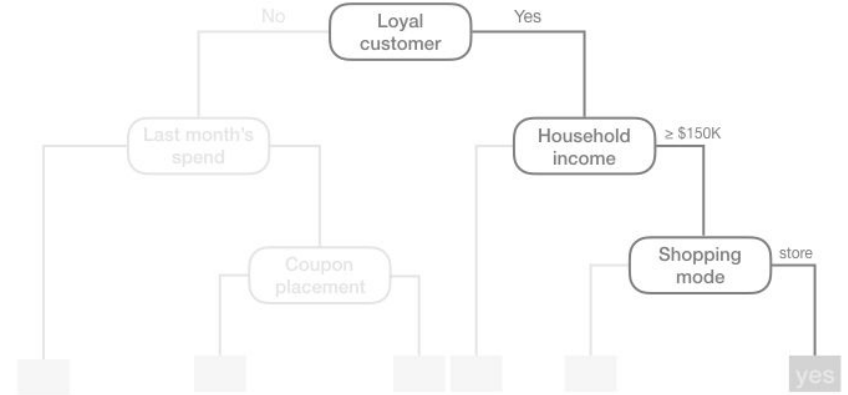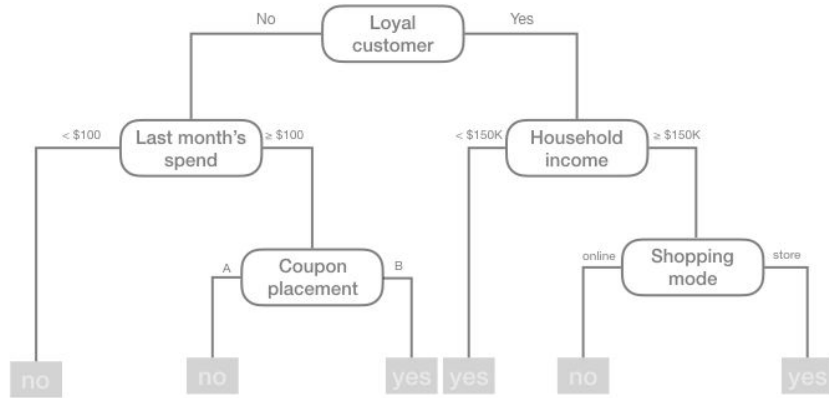- Forms of Linear Regression and Logistic Regression when you have too many predictors/ independent variables
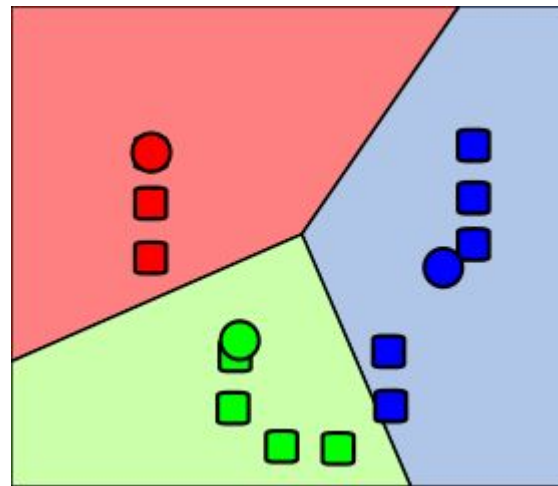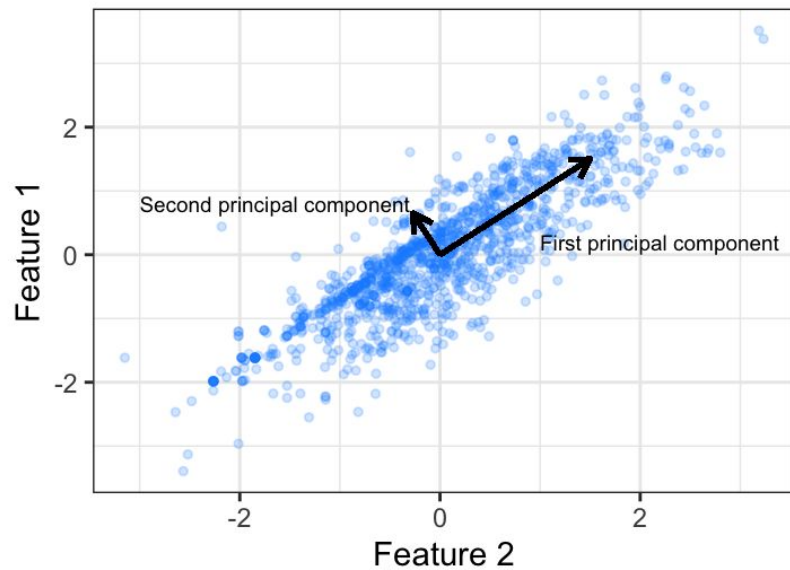
# Model: Resampling Methods

# Model: Tree-Based Methods
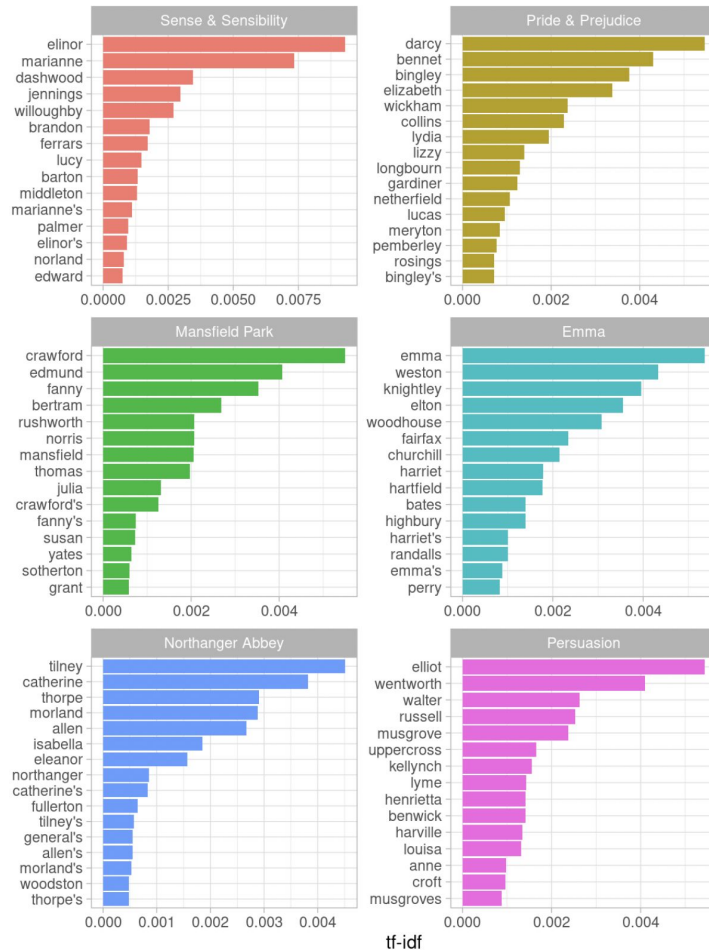
# Model: Unsupervised Learning

# Text Mining: Basics



Figure 3.4: Highest tf-idf words in each Jane Austen novel
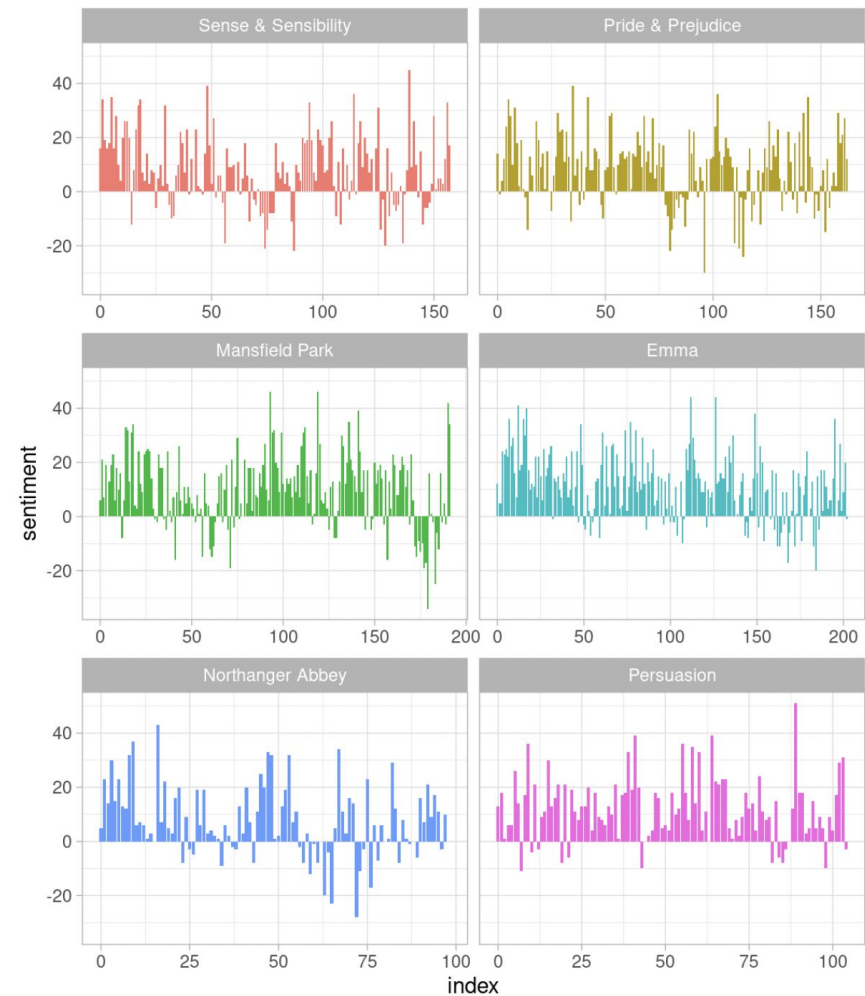
# Text Mining: Sentiment Analysis



Figure 2.2: Sentiment through the narratives of Jane Austen's novels

# Text Mining:
# Text Summarization and Information Extraction

# Readings

- Most are hands-on materials
  - Read
  - Try the code yourself
  - Tweak the code and see what happens
- Expect familiarity with the reading materials before class
- Clarify and extend in class
- Strongly encourage review and taking notes after class

# Output and Assessment

See the syllabus.

# R Setup

# In-class Exercise 1: Setup R

- Install R
- Install RStudio
- Open RStudio
- Run the following code in R Console

  ```
  install.packages("tidyverse")
  ```

- Post a screenshot of your Rstudio interface on CampusWire