

# WB21-Machine Learning with Big Data Day 1

Akitaka Matsuo

Essex IADS

# About me

- Aki Matsuo
  - Research fellow in Institute for Analytics and Data Science (IADS) and Department of Government, University of Essex
  - Ph.D in Political Science, Rice University
  - Research Interest
    - Political methodology (text analysis, scaling)
    - Legislative politics
    - Social media analysis
      - UK Politics / Japan v Korea
  - Member of quanteda project

# Your turn

- Name
- School/Company
- Discipline
- Position
- Your motivation
  - Big Data or ML

# Course schedule

- Day 1: (Big) Data Management
- Day 2: Database + Machine Learning Basics
- Day 3: Regression
- Day 4: Classification
- Day 5: Tree-based Method + Sparklyr

# Schedule Day 1

- Course info
- Discussion
- RStudio setting up
- tidyverse exercise
  - How it works
  - I will give you some time to work on coding
    - Group work
  - Then I will show the answer

# Course information

- Each class:
  - A bit of review and coding exercises
- If you are taking this course for a credit, please attend all lectures (or let me and Anna know when you will miss a class)
- For additional credit, do the take home assignment
  - Find the data, then carry out the analysis
- We will have office hours
  - What's your timezone?
  - 20 min slots
  - Google spreadsheet
    - <https://docs.google.com/spreadsheets/d/1nh7b3TjYmHIUmDRGL0MjoR-h0ykXweUmh3IEICZS65Q>
  - One zoom room for office hour

# Discussion

- Three-V of big data
  - Volume, Velocity and Variability
  - Which-V do the data of your interest fit in?
- Grimmer (2014) suggests two direction of ML in political science
  - Better causal inference
  - Improved measurement
  - How do your research interests relate to this argument?
- Lazar and Radford (2017)
  - Nowcasting. Should we be interested?
    - Fake news (U.S. election, vaccination misinformation)
  - Big data + Field experiments. What should be allowed?
    - You don't know about the boundary of the influence