

**ODI Workshop**

November 22, 2022



# The Opportunities of Big Data in Social Science Research

**Milena Tsvetkova**

London School of Economics and Political Science



# Outline

1. Overview
  1. Big data
  2. Data science methods
2. Examples
  1. Contagion and epidemics
  2. Political misinformation and polarization
  3. Poverty and inequality
3. Demo: Measuring socioeconomic differences in consumption patterns using the Facebook Marketing API
4. Problems, debates, and next steps



# Big Data

- Big data
  - High volume
  - High velocity
  - High variability/complexity
- New formats
  - Networks
  - Text
  - Images
  - Audio
  - Video





# Big Data

- New sources
  - Digital trace data
  - Satellite imagery
  - Digitized historical archives
  - Census/register data
  - A/B testing
  - Crowdsourced data donation
  - Browser tracking data
  - Surveillance cameras

The image displays three examples of how big data is used:

- Population Density Map:** A map of a city area where each colored dot represents one person. The legend indicates:
  - White (Blue)
  - Black (Green)
  - Asian (Red)
  - Hispanic (Orange)
- Facebook Election Reminder:** A screenshot of a Facebook message box. The message says "Today is Election Day" and "Find your polling place on the U.S. Politics Page and click the "I Voted" button to tell your friends you voted." It shows a "VOTE" button and a count of "01155376 People on Facebook Voted".
- Moral Machine:** A screenshot of the Moral Machine website. The top navigation bar includes Home, Judge, Classic, Design, Browse, About, and Feedback. The main content asks "What should the self-driving car do?" and shows two scenarios: one where the car must choose between hitting five people or one person, and another where it must choose between hitting a person or a group of people. Below each scenario is a "Show Description" button.

<https://www.moralmachine.net/>



# Data Science Methods

- Collecting data
  - Web scraping
  - APIs (Application Programming Interface)
- Storing and managing data
  - Relational databases – e.g., SQL
  - Non-relational databases – document/file collections
- Processing data
  - Cleaning – e.g., lemmatization
  - Labeling – e.g., crowdsourced image tagging and tweet classification
  - Matching – e.g., named entity recognition with regular expressions



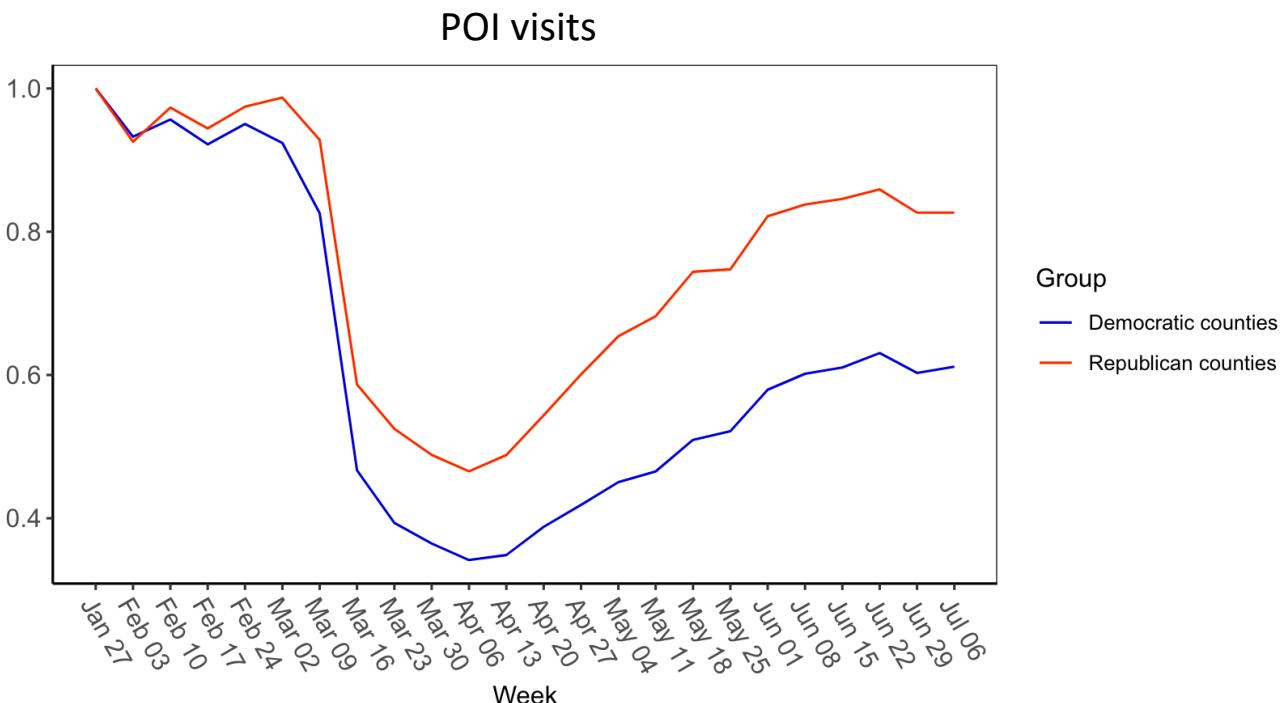
# Data Science Methods

- Analyzing and visualizing data
    - Data mining – discovering patterns
    - Machine learning – making predictions
    - Causal inference – identifying causal relationships
    - Natural language processing – understanding text
    - Network analysis – understanding relational patterns

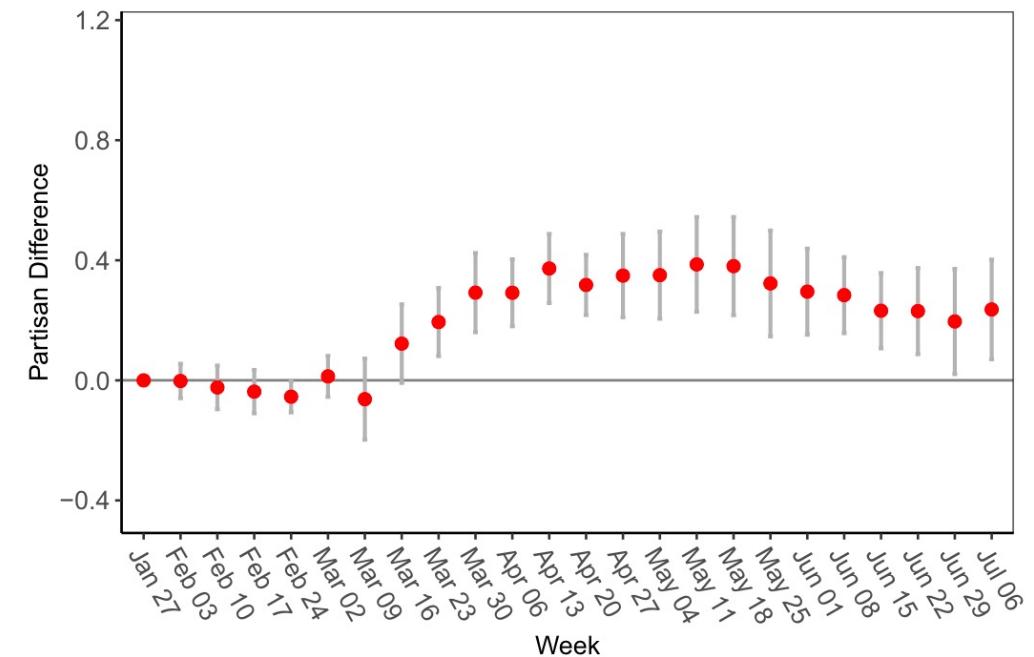


# Example: Contagion and Epidemics

- Were there partisan differences in social distancing during COVID?
- Mobile location data from  **SAFE GRAPH**



Estimated coefficient for effect of Republican vote share in county for  $\log(\text{POI visits})$

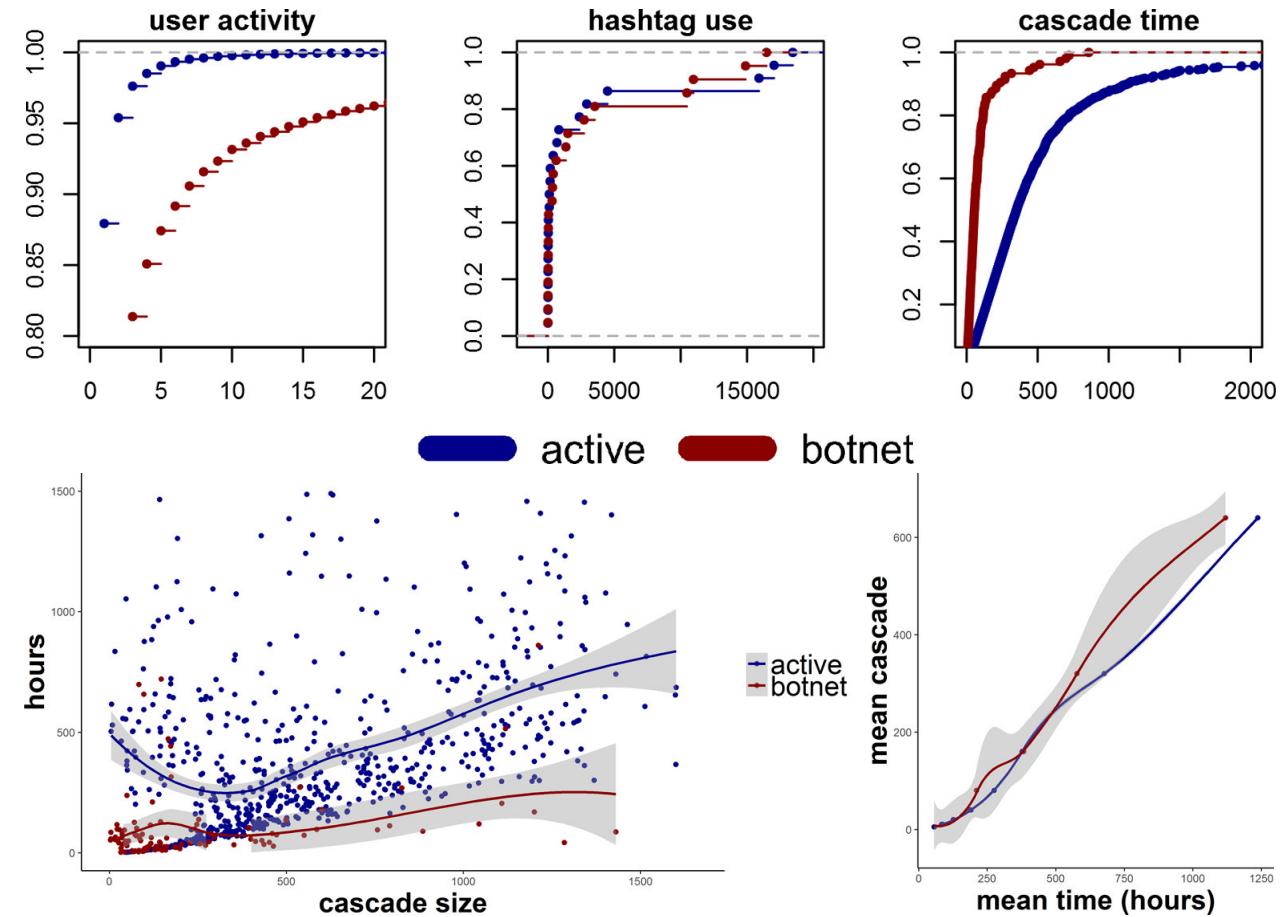
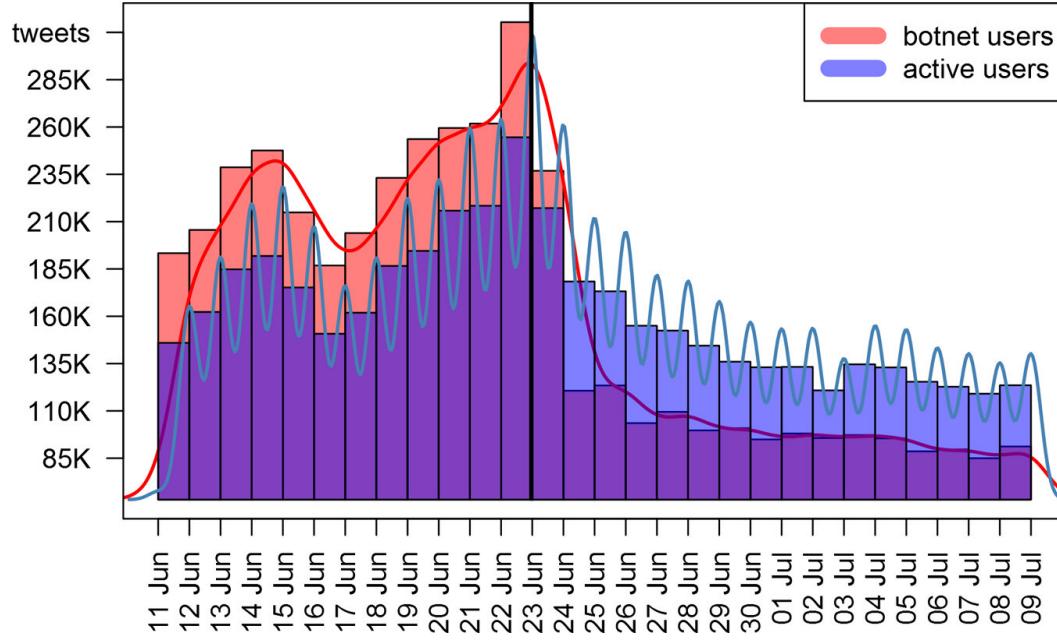


Allcott, H., Boxell, L., Conway, J., Gentzkow, M., Thaler, M., & Yang, D. (2020). Polarization and public health: Partisan differences in social distancing during the coronavirus pandemic. *Journal of Public Economics*, 191, 104254.



# Example: Political Misinformation

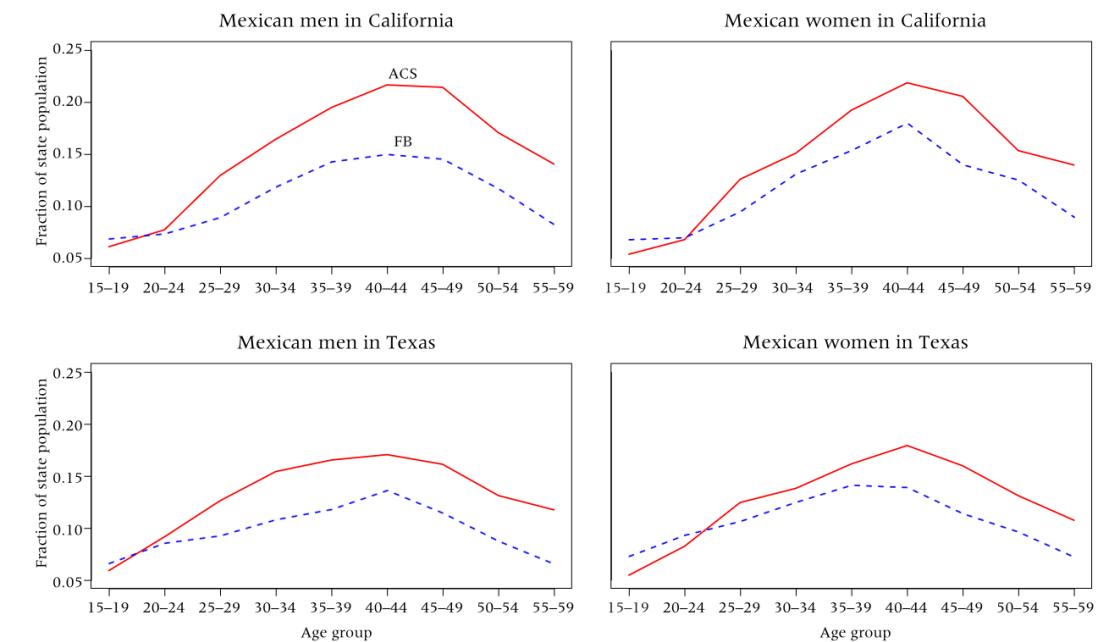
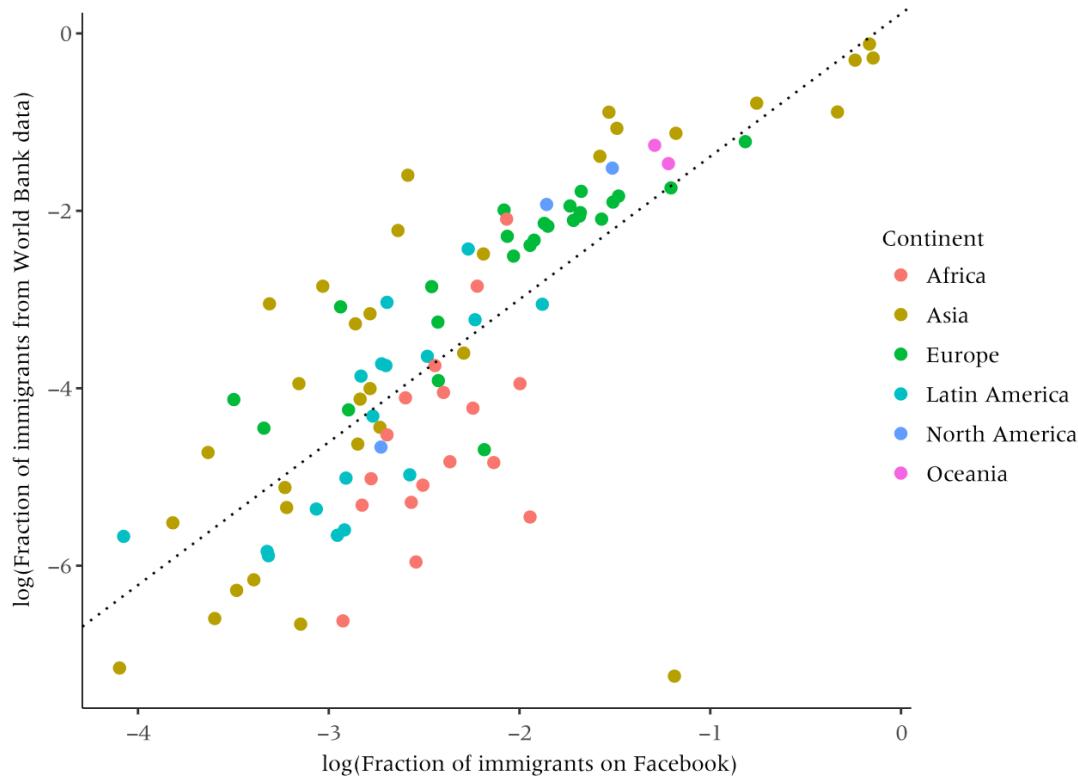
- Can bots on Twitter influence political opinion?





# Example: Poverty and Inequality

- Can we use Facebook ad data to estimate stocks of migrants within and across countries?



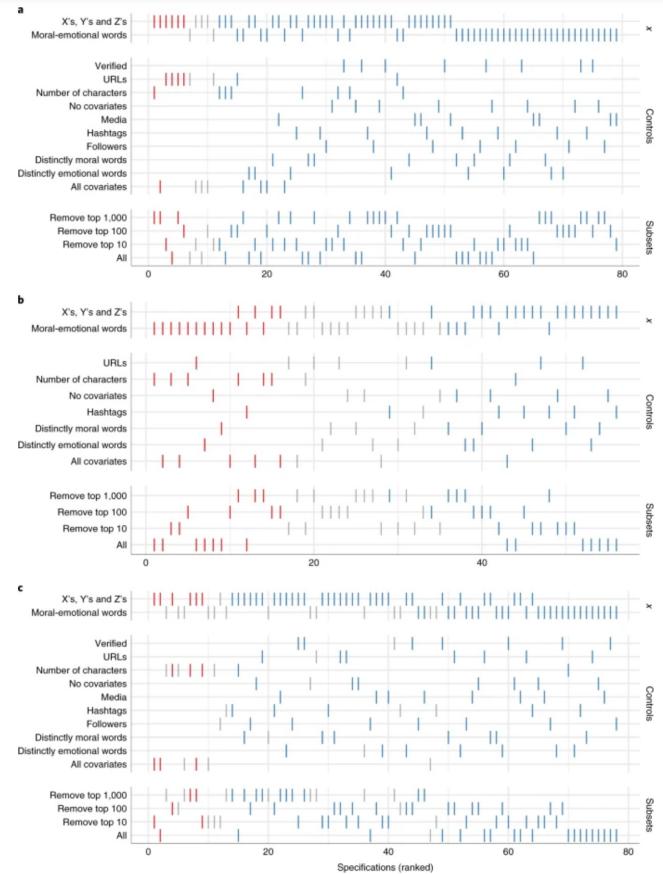


# Demonstration

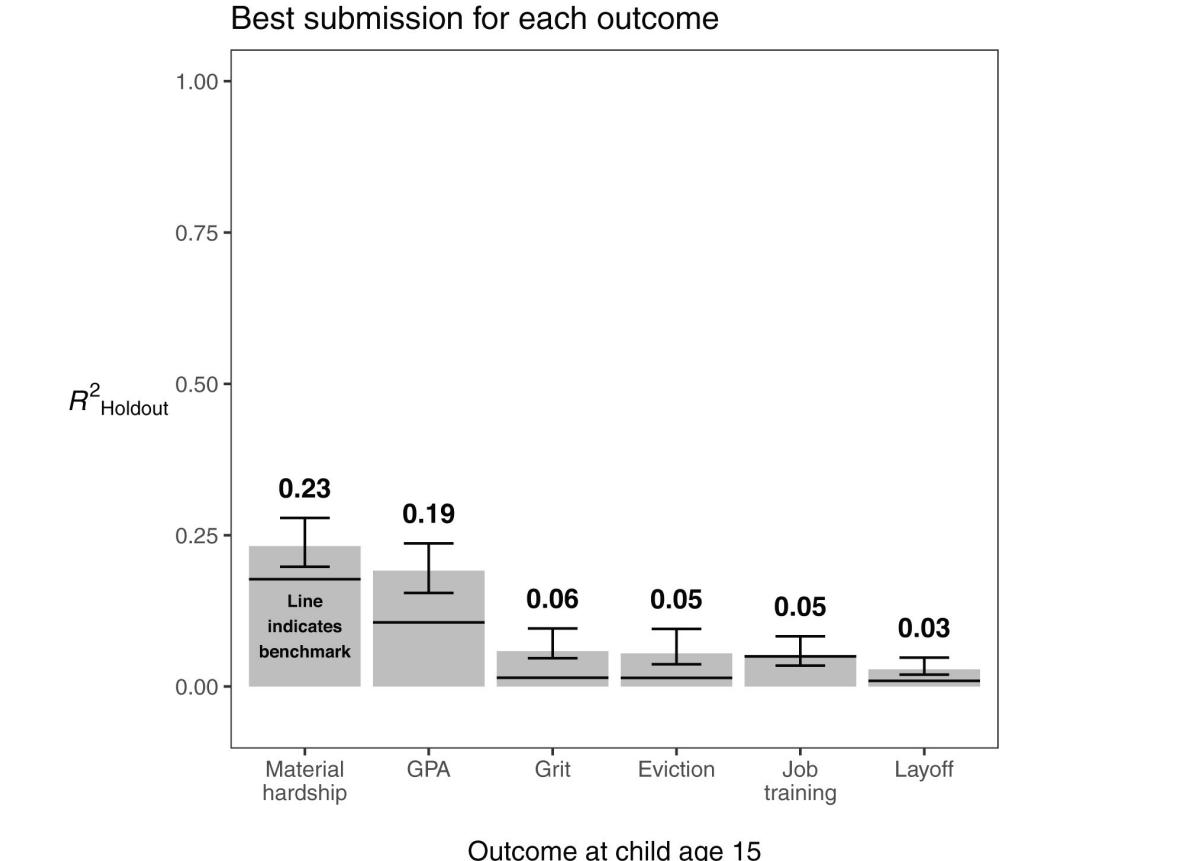
- View/work locally
  - Install python and VS Code on your machine
  - Go to <https://github.com/social-research/big-data-workshop> and download the repository
  - Open the (unzipped) folder in VS Code
- Alternatively, view/work remotely
  - Go to <https://colab.research.google.com/github/social-research/big-data-workshop>



# Problems and Debates: Reproducibility



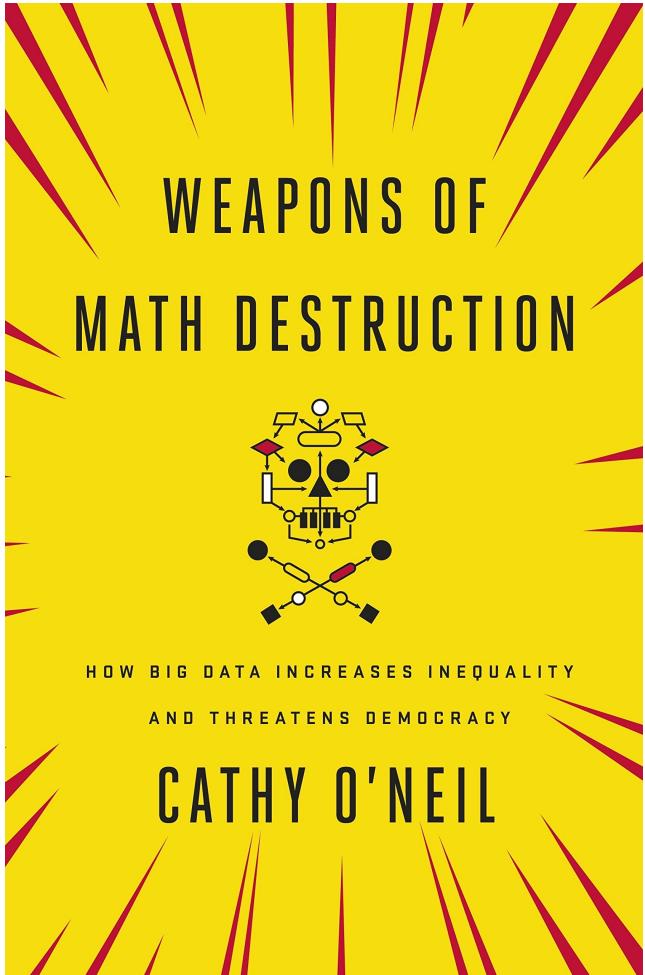
Burton, J. W., Cruz, N., & Hahn, U. (2021). Reconsidering evidence of moral contagion in online social networks. *Nature Human Behaviour*, 5(12), 1629-1635.



Salganik, M. J., et al. (2020). Measuring the predictability of life outcomes with a scientific mass collaboration. *Proceedings of the National Academy of Sciences*, 117(15), 8398-8403.



# Problems and Debates: Accountability



O'Neil, C. (2016). *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. Broadway Books.

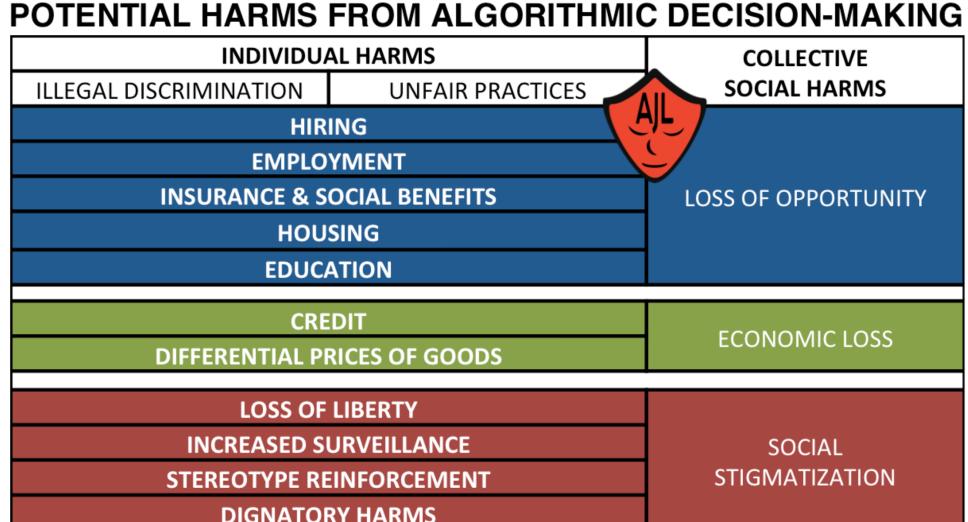


Chart Contents Courtesy of Megan Smith, Former CTO of the United States



Buolamwini, J., & Gebru, T. (2018, January). Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on Fairness, Accountability and Transparency* (pp. 77-91). PMLR.

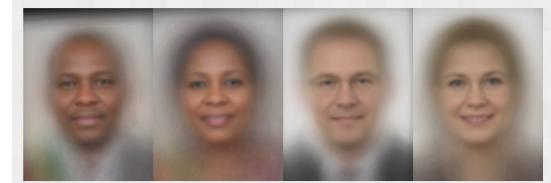
Rosalia @BonKamona

I saw a tweet saying "Google unprofessional hairstyles for work". I did. Then I checked the 'professional' ones 😊😊😊

2:04 PM - 5 Apr 2016

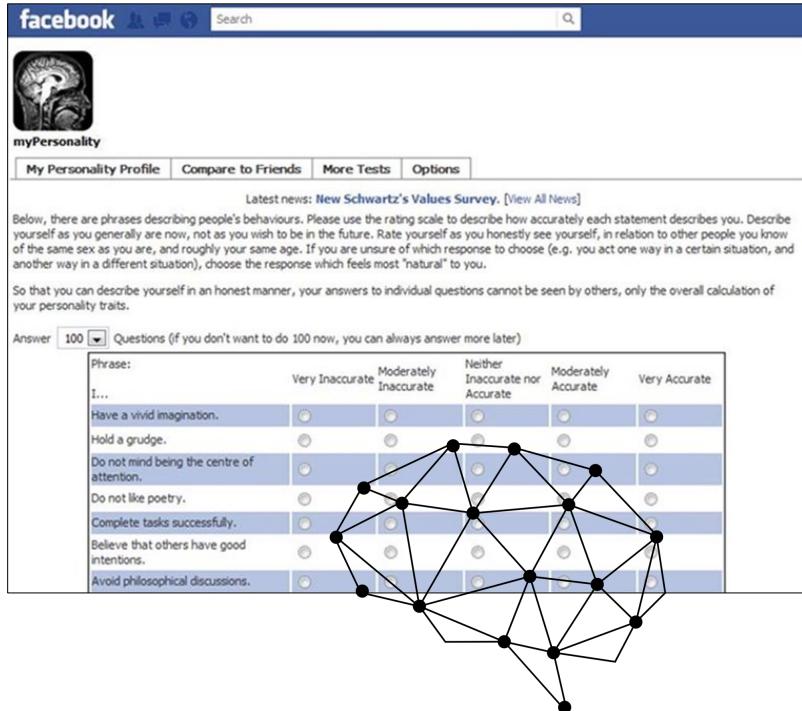


<http://gendershades.org>



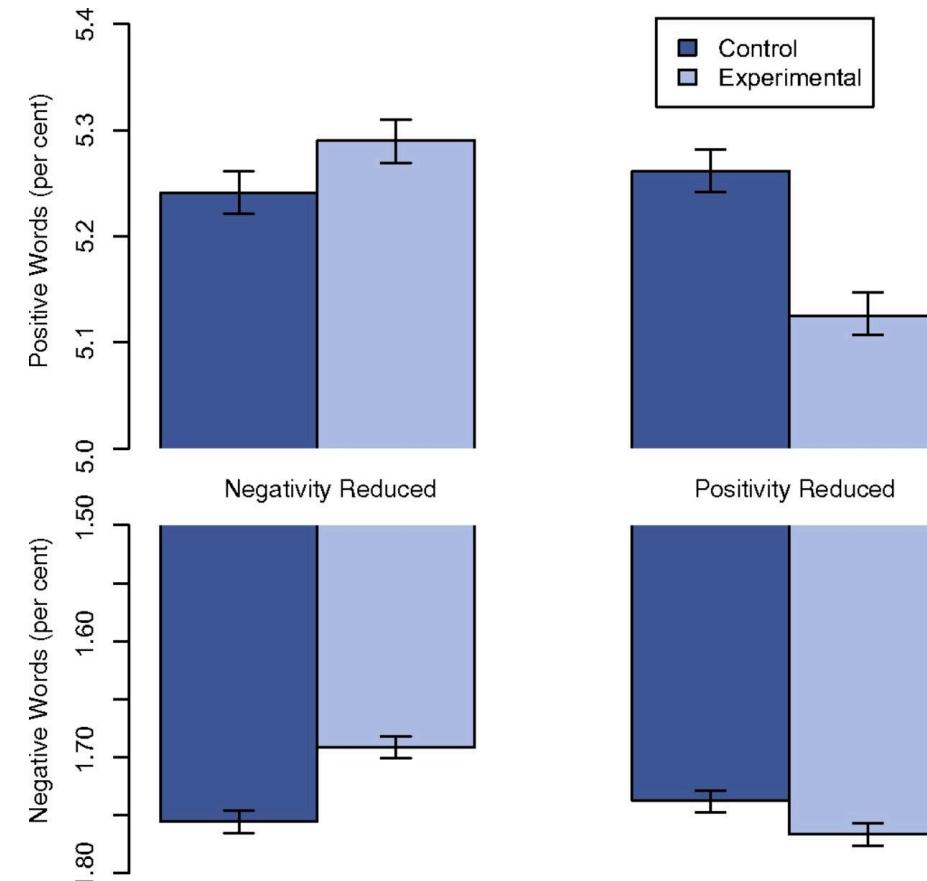


# Problems and Debates: Research Ethics



## Cambridge Analytica

Cadwalladr, C., & Graham-Harrison, E. (2018). Revealed: 50 million Facebook profiles harvested for Cambridge Analytica in major data breach. *The Guardian*, 17, 22.



Kramer, A. D., Guillory, J. E., & Hancock, J. T. (2014). Experimental evidence of massive-scale emotional contagion through social networks. *Proceedings of the National Academy of Sciences*, 111(24), 8788-8790.



# Next Steps

- **Summer course:** [LSE ME314: Introduction to Data Science and Machine Learning \(taught in R\)](#)
- **Online program:** [LSE Data Analytics Career Accelerator](#)
- **MSc program:** [LSE MSc in Applied Social Data Science](#)
- **Books**
  - Salganik, M. J. (2019). *Bit by Bit: Social Research in the Digital Age*. Princeton University Press.
  - Hogan, B. (2022). *From Social Science to Data Science: Key Data Collection and Analysis Skills in Python*. Sage.
- **Conference:** [International Conference on Computational Social Science](#)