Bit by Bit: Social Research in the Digital Age

Matthew J. Salganik Department of Sociology Princeton University

> Behave Lab University of Milan March 21, 2019



Isn't computational social science a fad?

Isn't computational social science a fad? No



Social Scientists \longleftrightarrow Data Scientists











Readymades



Custommades



Readymades



Custommades

https://commons.wikimedia.org/wiki/File:Duchamp_Fountaine.jpg https://commons.wikimedia.org/wiki/File:%27David%27_by_Michelangelo_JBU0001.JPG

Predicting poverty and wealth from mobile phone metadata

Joshua Blumenstock,¹* Gabriel Cadamuro,² Robert On³

https://doi.org/10.1126/science.aac4420









survey data











Readymade + Custommade



 ${\sf Readymade} + {\sf Custommade}$

Custommade



Readymade + Custommade

Custommade

10 times faster50 times cheaper

Blumenstock et al. (2015), Figure 3



Readymades



Custommades

https://commons.wikimedia.org/wiki/File:Duchamp_Fountaine.jpg https://commons.wikimedia.org/wiki/File:%27David%27_by_Michelangelo_JBU0001.JPG



$\hat{\beta}$ & \hat{y}

Mullainathan and Spiess (2017): http://dx.doi.org/10.1257/jep.31.2.87



Fragile Families Challenge Matthew Salganik, Ian Lundberg, Alex Kindel, Sara McLanahan, and the participants in the Fragile Families Challenge

Funding for FFCWS provided by NICHD (R01HD36916, R01HD39135, R01HD40421) and a consortium of private foundations, including the Robert Wood Johnson Foundation. Funding for FFC provided by the Russell Sage Foundation and the Overdeck Fund. FFC Board of Advisors: Jeanne Brooks-Gunn, Kathryn Edin, Barbara Engelhardt, Irwin Garfinkel, Moritz Hardt, Dean Knox, Nicholas Lemann, Karen Levy, Sara McLanahan, Arvind Narayanan, Timothy Nelson, Matthew Salganik, and Duncan Watts.

$$Y = \mathsf{E}\left(Y \mid \vec{X}\right) + \epsilon$$

$$Y = \mathsf{E}\left(Y \mid \vec{X}\right) + \epsilon$$
Attainment

$$\mathbf{Y} = \mathsf{E}\left(\mathbf{Y} \mid \vec{X}\right) + \epsilon$$

Attainment

- Academic
 - achievement
- Occupation
- Income

$$Y = \underbrace{\mathsf{E}\left(Y \mid \vec{X}\right)}_{\mathsf{Attainment}} + \epsilon$$
Attainment
- Academic
achievement
Predictable
component

- Occupation
- Income

$$Y = \underbrace{\mathsf{E}\left(Y \mid \vec{X}\right)}_{\mathsf{Attainment}} + \epsilon$$
Attainment
- Academic
achievement
Predictable
component

- Occupation
- Income



- Occupation
- Income

$$Y = \underbrace{\mathsf{E}\left(Y \mid \vec{X}\right)}_{\mathsf{Attainment}} + \epsilon$$
Attainment
- Academic
achievement
Predictable
component

- Occupation
- Income



- Occupation
- Income



- Occupation
- Income



Theories focus on the predictable component, but empirically the unpredictable component dominates
Scientific reasons

- Scientific reasons
 - Basic social fact
 - Discovery

- Scientific reasons
 - Basic social fact
 - Discovery
- Policy reasons







- Birth cohort panel study
- \blacktriangleright \approx 5,000 children born in 20 U.S. cities with an over-sample of non-marital births
- Followed from birth through age 15
- Already used in hundreds of papers and dozens of dissertations









Outcomes

- Child: GPA (continuous), Grit (continuous)
- Household: Eviction (binary), Material hardship (continuous)
- Primary care giver: Job training (binary), Job loss (binary)

459 researchers applied to participate. Many worked in interdisciplinary teams. Goal: Make a prediction that minimizes mean square error on the hold-out set

$$MSE_{holdout} = \frac{\sum_{i \in holdout} (\hat{y}_i - y_i)^2}{n_{holdout}}$$

More on privacy and ethics audit: https://arxiv.org/abs/1809.00103 Using a large, high-quality social science dataset collected since birth and modern machine learning methods, how accurately can we predict outcomes from children, parents, and families?

$$R_{holdout}^{2} = 1 - \frac{\sum_{i \in holdout} (\hat{y}_{i} - y_{i})^{2}}{\sum_{i \in holdout} (\bar{y}_{train} - y_{i})^{2}}$$

Using a large, high-quality social science dataset collected since birth and modern machine learning methods, how accurately can we predict outcomes from children, parents, and families?

$$R_{holdout}^2 = 1 - rac{\sum_{i \in holdout} (\hat{y}_i - y_i)^2}{\sum_{i \in holdout} (\bar{y}_{train} - y_i)^2}$$

Before I show the results, let's vote . . .





Is this even better than a benchmark model?





Green line: 4 variable linear regression model

Material hardship





Eviction





What can we learn looking at the all the predictions?

Squared error predicting materialHardship Squared error 0.5 0.4 0.3 0.2 0.1 0.0 I Challenge Team





Next questions:

Is it possible to get better predictive performance for this data and prediction task? Next questions:

- Is it possible to get better predictive performance for this data and prediction task?
- Why is the unpredictability so high even using modern machine learning methods and what many social scientists would consider to be large and high-quality data?

Not enough cases

- Not enough cases
- Measurement error in existing variables (particularly outcomes)

- Not enough cases
- Measurement error in existing variables (particularly outcomes)
- Important unmeasured variables

- Not enough cases
- Measurement error in existing variables (particularly outcomes)
- Important unmeasured variables

How can we learn about important measurement error and unmeasured variables?

In-depth interviews



What's next?

Next steps:

One community paper (including all code and predictions)

Next steps:

- One community paper (including all code and predictions)
- Special issue of *Socius*
 - 12 submitted manuscripts from Challenge participants (all with accompanying code and computing environment)
- One community paper (including all code and predictions)
- Special issue of *Socius*
 - 12 submitted manuscripts from Challenge participants (all with accompanying code and computing environment)
 - 3 papers from our group

- One community paper (including all code and predictions)
- Special issue of *Socius*
 - 12 submitted manuscripts from Challenge participants (all with accompanying code and computing environment)
 - 3 papers from our group
 - "Privacy, ethics, and data access: A case study of the Fragile Families Challenge" by Lundberg, Narayanan, Levy, & Salganik, https://arxiv.org/abs/1809.00103

- One community paper (including all code and predictions)
- Special issue of *Socius*
 - 12 submitted manuscripts from Challenge participants (all with accompanying code and computing environment)
 - 3 papers from our group
 - "Privacy, ethics, and data access: A case study of the Fragile Families Challenge" by Lundberg, Narayanan, Levy, & Salganik, https://arxiv.org/abs/1809.00103
 - "Improving metadata infrastructure for complex surveys: Insights from the Fragile Families Challenge" by Kindel, Catena, Hartshorne, Jaeger, Koffman, McLanahan, Phillips, Rouhani, Vinh, & Salganik, https://osf.io/93ywg/

- One community paper (including all code and predictions)
- Special issue of *Socius*
 - 12 submitted manuscripts from Challenge participants (all with accompanying code and computing environment)
 - 3 papers from our group
 - "Privacy, ethics, and data access: A case study of the Fragile Families Challenge" by Lundberg, Narayanan, Levy, & Salganik, https://arxiv.org/abs/1809.00103
 - "Improving metadata infrastructure for complex surveys: Insights from the Fragile Families Challenge" by Kindel, Catena, Hartshorne, Jaeger, Koffman, McLanahan, Phillips, Rouhani, Vinh, & Salganik, https://osf.io/93ywg/
 - "Successes and struggles with computational reproducibility in the Fragile Families Challenge" by Liu & Salganik, https://osf.io/preprints/socarxiv/g3pdb/

$\hat{\beta}$ & \hat{y}

Mullainathan and Spiess (2017): http://dx.doi.org/10.1257/jep.31.2.87



- Read: http://www.bitbybitbook.com
- Teach: http://www.bitbybitbook.com/en/teaching/ (and Italian version coming soon from il Mulino)

▶ 6 year gap between end of background data and outcome

- ▶ 6 year gap between end of background data and outcome
- large social disruption—the Great Recession—between end of background data and outcome

- 6 year gap between end of background data and outcome
- large social disruption—the Great Recession—between end of background data and outcome
- the sample design of the Fragile Families study

- 6 year gap between end of background data and outcome
- large social disruption—the Great Recession—between end of background data and outcome
- the sample design of the Fragile Families study
- outcomes measured when child was 15
- outcomes are at a relatively narrow point in time rather than average over a longer time period (e.g., grades last semester vs grades in high school)

Advances in hurricane prediction

Data from the NOAA National Hurricane Center (VHC) (J3) show that forecast errors for tropical storms and hurricanes in the Atlantic basin have failen rapidly in recent decades. The graph shows the forecast error in nautical miles (1 n mi = 1.852 km) for a range of time intervals.



Advances in hurricane prediction

Data from the NOAA National Hurricane Center (VHC) (J3) show that forecast errors for tropical storms and hurricanes in the Atlantic basin have failen rapidly in recent decades. The graph shows the forecast error in nautical miles (1 n mi = 1.852 km) for a range of time intervals.



Can we do this?

Advances in hurricane prediction

Data from the NOAA National Hurricane Center (NHC) (13) show that forecast errors for tropical storms and hurricanes in the Atlantic basin have fallen rapidly in recent decades. The graph shows the forecast error in nautical miles (1 n mi = 1.852 km) for a range of time intervals.



- Can we do this?
- Should we do this?

 Dozens already happening all over the world with interesting similarities and differences

- Dozens already happening all over the world with interesting similarities and differences
- Already strong community around each survey

- Dozens already happening all over the world with interesting similarities and differences
- Already strong community around each survey
- Code from a single Challenge can be repurposed to create many simulated Challenges

- Dozens already happening all over the world with interesting similarities and differences
- Already strong community around each survey
- Code from a single Challenge can be repurposed to create many simulated Challenges
- Data collected with informed consent under well-developed ethical frameworks

- Dozens already happening all over the world with interesting similarities and differences
- Already strong community around each survey
- Code from a single Challenge can be repurposed to create many simulated Challenges
- Data collected with informed consent under well-developed ethical frameworks
- Likely to spur useful scientific developments













