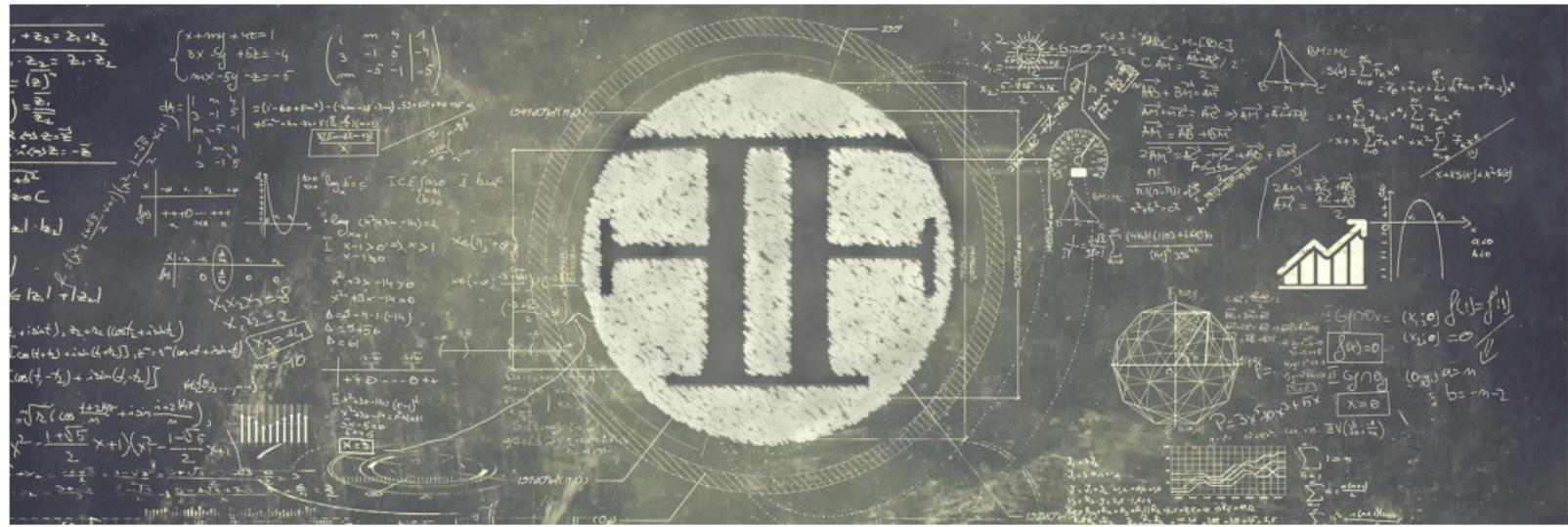


Class slides for Tuesday, November 10: Predictability of life trajectories

Matthew J. Salganik

COS 597E/SOC 555 Limits to prediction
Fall 2020, Princeton University



Fragile Families Challenge

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

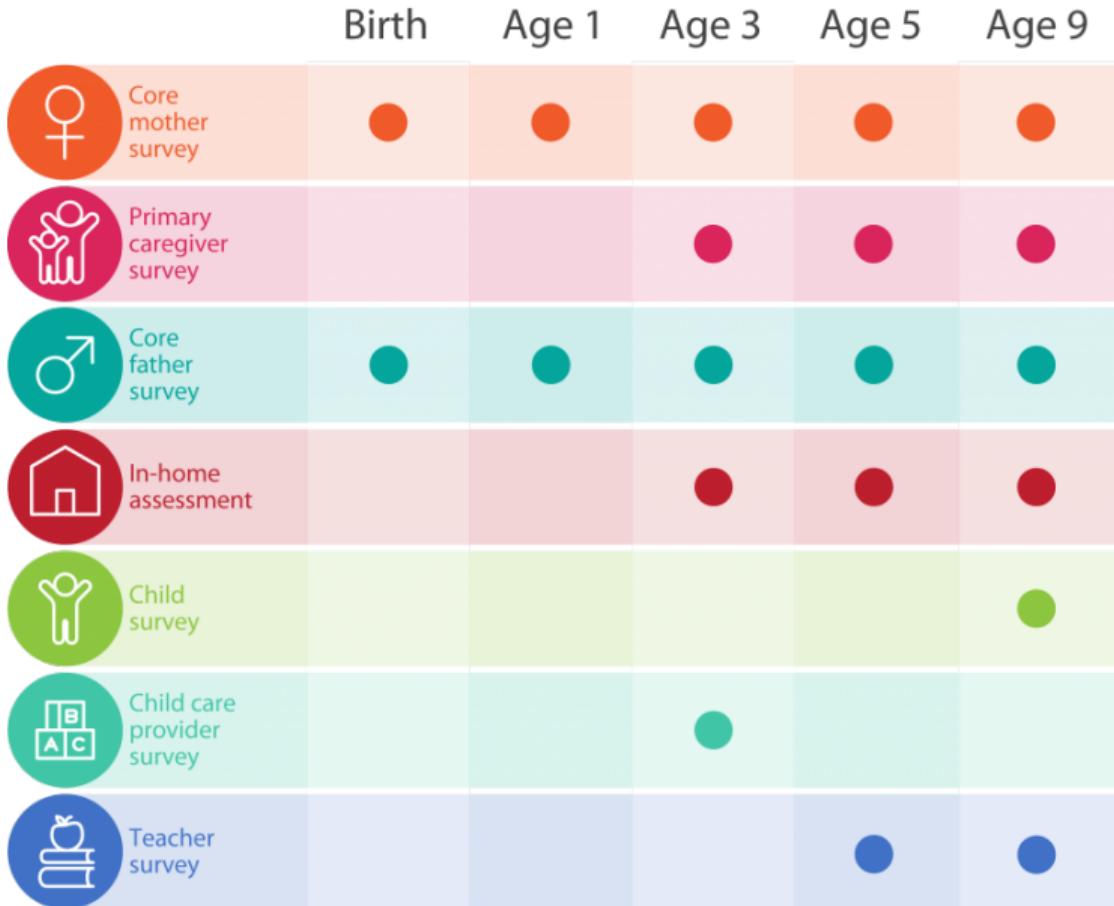
Mullainathan and Spiess (2017)

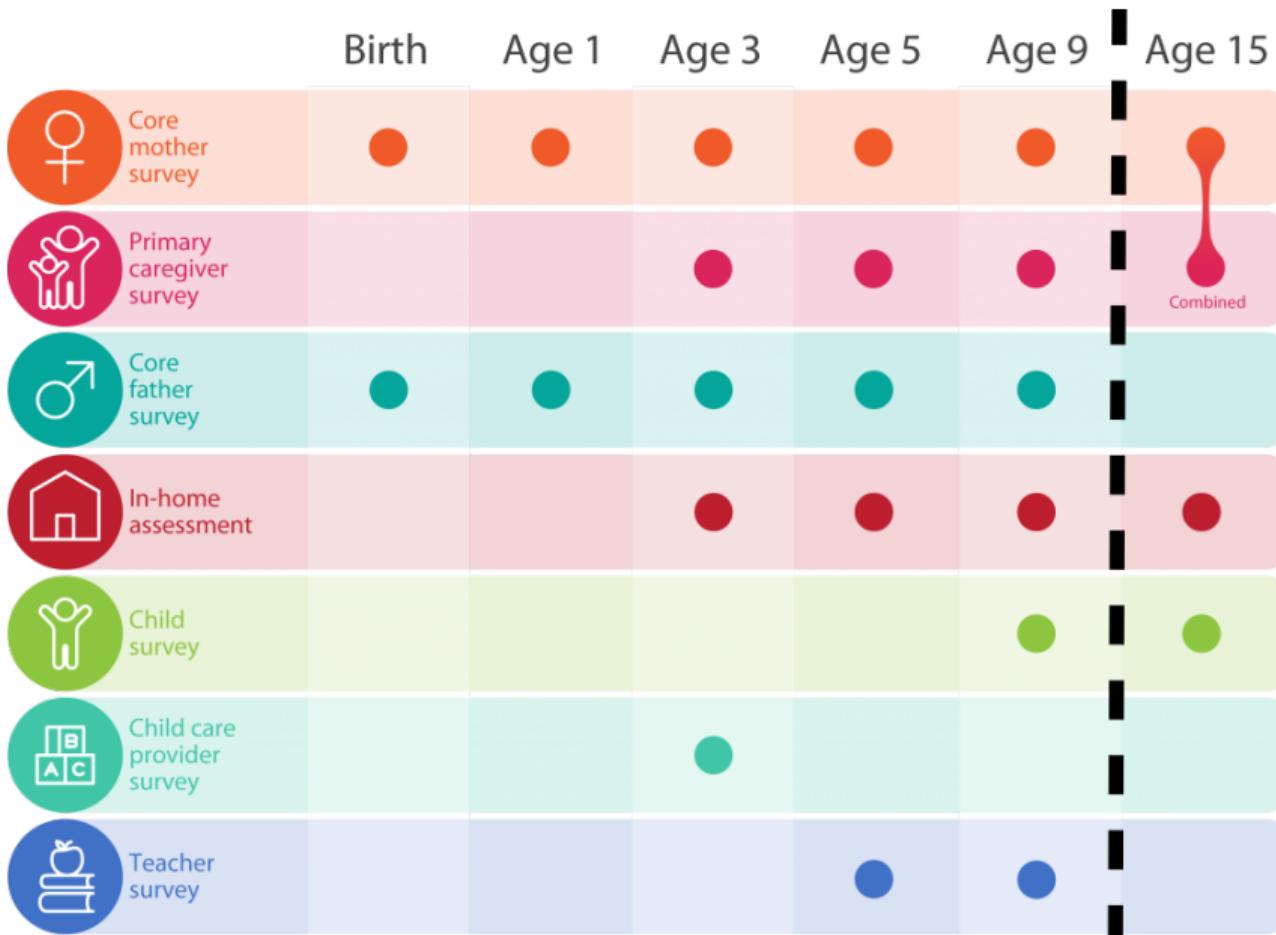
FF Fragile Families

& Child Wellbeing Study
PRINCETON | COLUMBIA



- ▶ Birth cohort panel study
- ▶ ≈ 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





4,242 families

Birth to age 9
12,942 features

Age 15
1,500 features

4,242 families

Birth to age 9
12,942 variables

Information about child and family

Background data

Age 15
6 variables

Training

Leaderboard

Holdout

Outcome
data

Outcomes

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

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

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

More on privacy and ethics audit: Lundberg et al. (2019)

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?

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

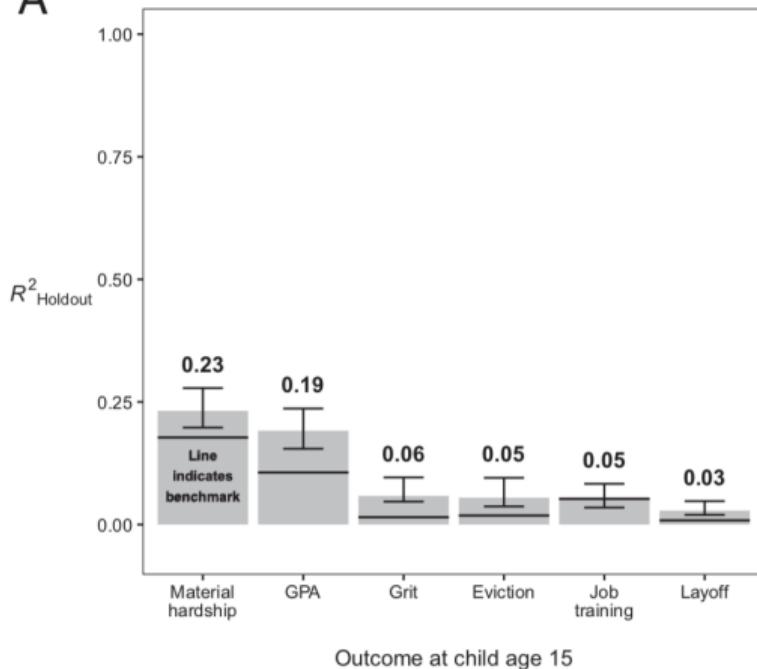
Using a large, high-quality social science dataset collected since birth and modern machine learning methods, how accurately can we predict outcomes from children, parents, and families?

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

Before I show the results, let's vote . . . (Also notice the switch from absolute to relative metric)

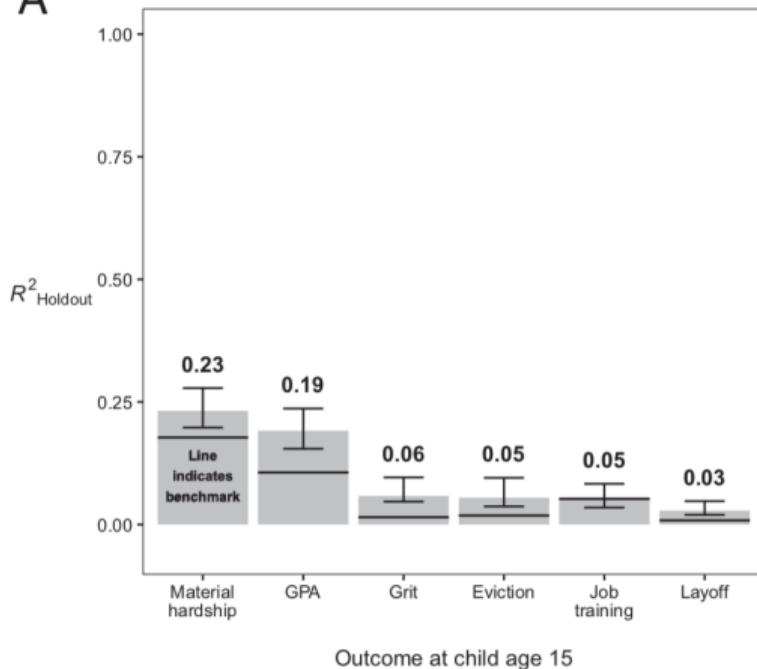
A

Best submission for each outcome

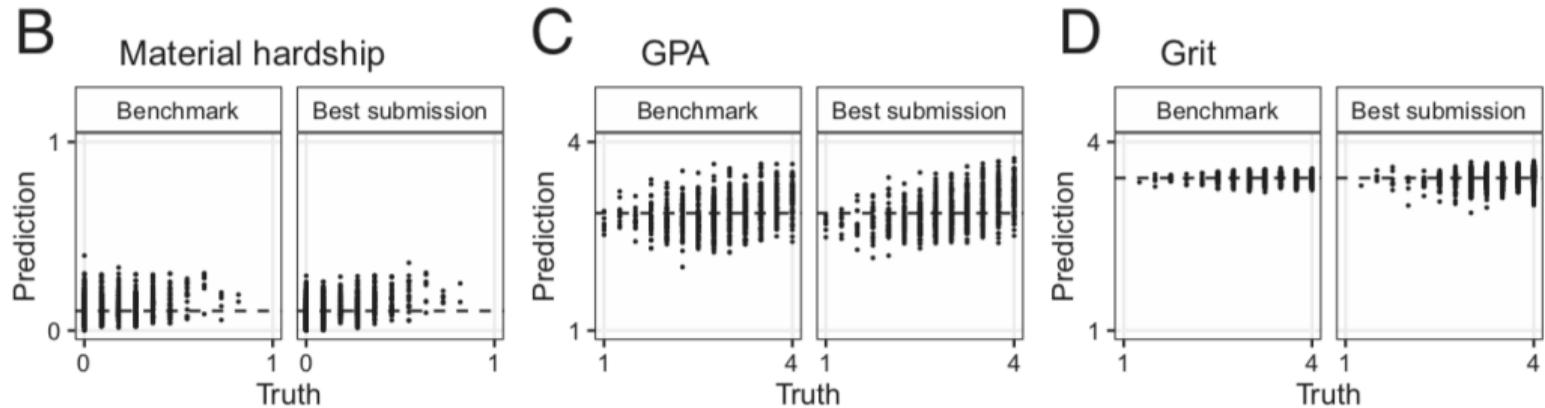


A

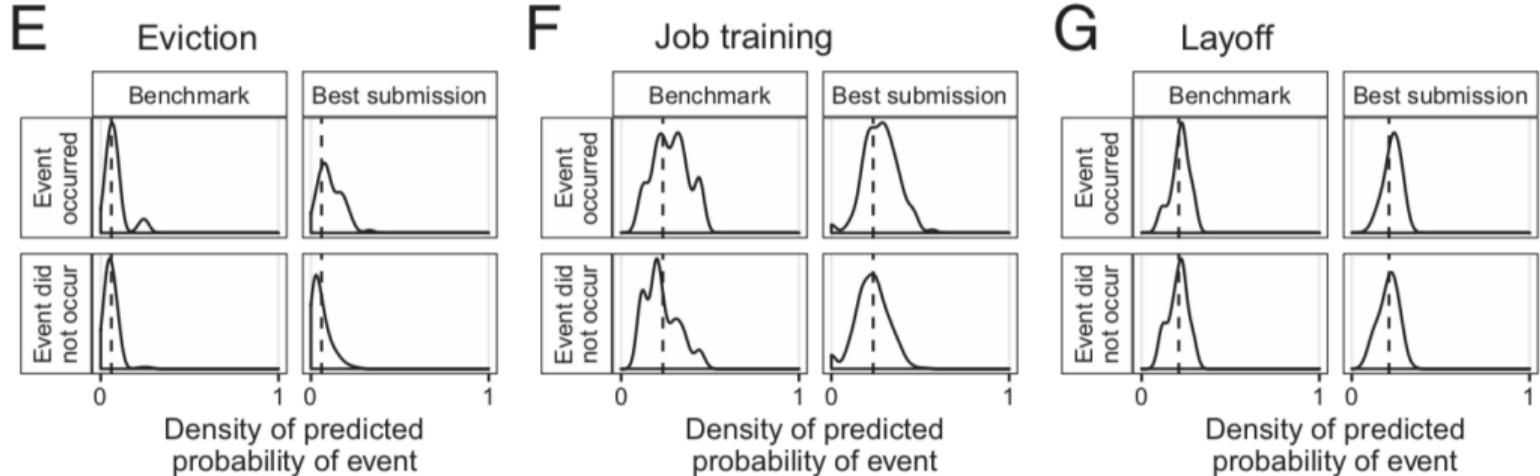
Best submission for each outcome



What is the right y-scale here?



- ▶ Best submission is like the mean of the training data with a bit of signal, not like the truth with a bit of noise
- ▶ Best submission is qualitatively similar to the benchmark model



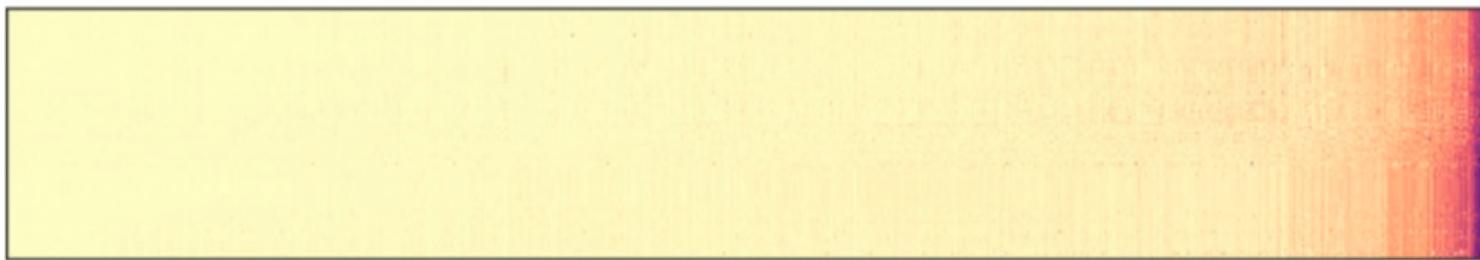
- ▶ Best submission is like the mean of the training data with a bit of signal, not like the truth with a bit of noise
- ▶ Best submission is qualitatively similar to the benchmark model

We can learn a lot by looking at all the valid submissions, not just the best

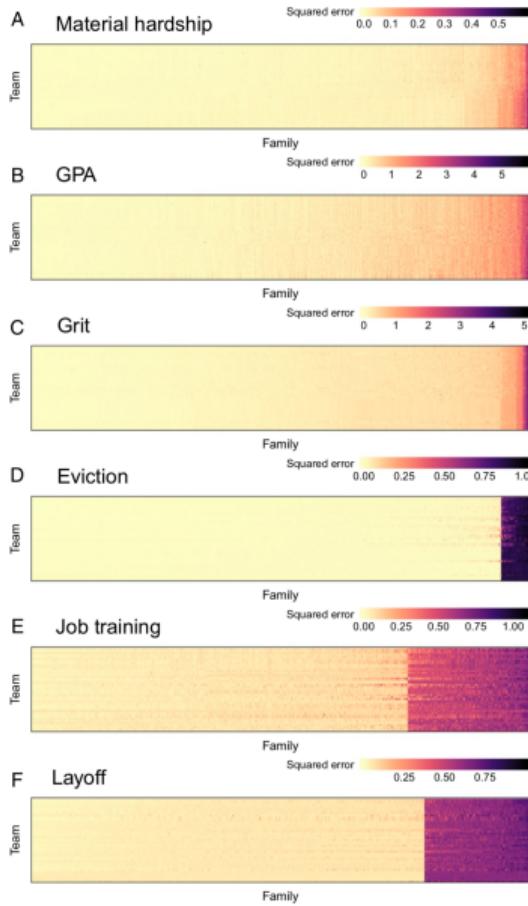
A Material hardship

Squared error
0.0 0.1 0.2 0.3 0.4 0.5

Team



Family



A Material hardship



Team



Family

Ensembling didn't improve things much. :(



Researchers must reconcile an “understanding/prediction” paradox

Researchers must reconcile an “understanding/prediction” paradox

- ▶ We don't understand much

Researchers must reconcile an “understanding/prediction” paradox

- ▶ We don't understand much
- ▶ Prediction is not a good measure of understanding

Researchers must reconcile an “understanding/prediction” paradox

- ▶ We don't understand much
- ▶ Prediction is not a good measure of understanding
- ▶ Our current understanding is correct but incomplete

Was the Fragile Families Challenge a success or failure?

Behind the scenes: Designing and running a mass collaboration

Too hard, don't do it

Use a mix of approaches to get folks to participate

WHY PARTICIPATE?

HELP THE WORLD:

The Fragile Families Challenge is designed to produce scientific knowledge that can be used to improve the lives of disadvantaged children in the United States. Even more than that, we hope the Fragile Families Challenge can serve as a model for how social scientists and data scientists can collaborate on problems of societal importance.

LEARN NEW SKILLS:

The Fragile Families Challenge blends ideas from social science and data science. Maybe you're a data scientist that wants to start working with social data? Maybe you're a social scientist that wants to learn more about machine learning? Either way, the Fragile Families Challenge is for you. This blending of ideas also makes the Fragile Families Challenge ideal to assign in a class that you are teaching.

WIN PRIZES:

We awarded prizes to participants who made important contributions to the project. All prize winners were given an all-expenses paid trip to Princeton University for a scientific workshop.

HAVE FUN:

The Fragile Families Challenge could be worked on in teams, and we hoped that participants would enjoy working with data, learning new skills, and cooperating and competing with people from all over the world.

GET INVOLVED IN SCIENTIFIC RESEARCH:

The Fragile Families Challenge is real scientific research. While working on the project, participants have a chance to interact with the other participating scientists and the distinguished researchers on our Board of Advisors.

PUBLISH PAPERS:

We are publishing the results of the Fragile Families Challenge in scientific journals, both individually and collectively. Participants who made important contributions had the opportunity to be a co-author on the paper describing the results of the Fragile Families Challenge.

“Hack” existing scientific reward mechanisms: papers and prizes.

Measuring the predictability of life outcomes with a scientific mass collaboration

Matthew J. Salganik¹, Ian Lundberg¹, Alexander T. Kindred², Caitlin E. Ahern³, Khaled Al-Ghoreini⁴, Abdullah Almarzouq^{5,6}, Drew M. Altchul⁷, Jessie E. Brand⁸, Nicla Bohne Canegiela⁹, Ryan James Compton¹⁰, Debanjan Datta¹¹, Thomas Davidson¹², Anna Filippova¹³, Connor Gilroy¹⁴, Brian J. Goode¹⁵, Eman Jahanri¹⁶, Rida Kashyap¹⁷, Arif Kirchner¹⁸, Stephen McKay¹⁹, Allison C. Merseth²⁰, Alex Pechard²¹, Karen Polens²², Heidi Price²³, Dovile Ruzicka²⁴, Michael V. Salter²⁵, Daniel M. Selsky²⁶, Valeria Sestini²⁷, Adriana Utrera²⁸, Erik H. Wang²⁹, Murva Adeny³⁰, Abdulkadir Altaraj³¹, Bednar Alshehri³², Redwan Ansari³³, Ryan R. Argote³⁴, Lilia Baer-Bouit³⁵, Morton Bodzin³⁶, Bo-Byeon Chang³⁷, William Egger³⁸, Gregory Falsetto³⁹, Zhilan Fan⁴⁰, Jeremy Freese⁴¹, Sergey Gulyayev⁴², Josh Gage⁴³, Yuxi Guo⁴⁴, Andrey Habermann-Mammen⁴⁵, Sonja F. Huskamp⁴⁶, Jennifer J. Hunt⁴⁷, Kristina I. Hwang⁴⁸, Hengyi Jiang⁴⁹, Farah M. Hosseini⁵⁰, Daniel Hsu⁵¹, Naveen Jain⁵², Kuan Jen⁵³, David Jurgens⁵⁴, Patrick Kannan⁵⁵, Ang Karapetyan⁵⁶, E. H. Kim⁵⁷, Ron Leinwand⁵⁸, Nasja Liu⁵⁹, Malte Möller⁶⁰, Andrew E. Mack⁶¹, Mayank Matihajer⁶², Noa Mandell⁶³, Helge Marathiotis⁶⁴, Olafur Mousavi-Garcia⁶⁵, Vicenç Muñoz⁶⁶, Katharina Mutschler⁶⁷, Ahmad Mumtaz⁶⁸, Ondjoum Ndi⁶⁹, William Nowak⁷⁰, Jennifer O’Donnell⁷¹, Daniel Ort⁷², Karen Oren⁷³, Kaye McPherson⁷⁴, Michael Portnoy⁷⁵, Crystal Olear⁷⁶, Terekirat Rau⁷⁷, Anshu Sargyan⁷⁸, Theresia Schaffner⁷⁹, London Schubert⁸⁰, Bryan Schoenfeld⁸¹, Ben Sender⁸², Jonathan D. Tang⁸³, Emma Tsourkoy⁸⁴, Austin von Loon⁸⁵, Onur Varol⁸⁶, Xiaofei Wong⁸⁷, Zhi Wang⁸⁸, Julia Wang⁸⁹, Yuxia Wang⁹⁰, Chandrani Venkateswaran⁹¹, Kirstie Whittaker⁹², Melia K. Winters⁹³, Wei Lee Worcester⁹⁴, Daniel Wu⁹⁵, Christopher Xiang⁹⁶, Karen Yano⁹⁷, Jason Yin⁹⁸, Bioggy Zhao⁹⁹, Chanyan Zhu¹⁰⁰, Jeanne Brooks-Gunn¹⁰¹, Barbara E. Engelhardt¹⁰², Morris Hardt¹⁰³, Deon Kross¹⁰⁴, Karen Levy¹⁰⁵, Arvind Narayanan¹⁰⁶, Preston M. Stewart¹⁰⁷, Duncan J. Watts¹⁰⁸, and Sara McLanahan¹⁰⁹

“Hack” existing scientific reward mechanisms: papers and prizes.

Measuring the predictability of life outcomes with a scientific mass collaboration

Matthew J. Salganik¹, Ian Lundberg², Alexander T. Kindel¹, Caitlin E. Ahern³, Khaled Al-Ghoreini⁴, Abdullah Almarzouq^{5,6}, Drew M. Attchell⁷, Jessie E. Brand⁸, Nicla Bohne Canegiela⁹, Ryan James Compton¹⁰, Debanjan Dasgupta¹¹, Thomas Davidson¹², Anna Filippova¹³, Connor Gilroy¹⁴, Brian J. Goode¹⁵, Eman Jahanri¹⁶, Rida Kashyap¹⁷, Amrit Kirchner¹⁸, Stephen McKay¹⁹, Allison C. Mersky²⁰, Alex Pernstein²¹, Karen Polens²², Heidi Riedl²³, Dovile Sipaviene²⁴, Michael V. Smith²⁵, Daniel M. Stein²⁶, Vicki Selsky²⁷, Adam Steuerwald²⁸, Erik H. Wang²⁹, Murva Aden³⁰, Abdulkadir Altayir³¹, Bednar Alshehri³², Redwan Ansari³³, Ryan R. Argote³⁴, Lilia Baer-Bouckaert³⁵, Morris Badih³⁶, Bo-Byeon Chang³⁷, William Egger³⁸, Gregory Falsetti³⁹, Zhihsu Fan⁴⁰, Jeremy Freese⁴¹, Sergey Gulyaev⁴², Josh Gage⁴³, Yuxi Guo⁴⁴, Andrey Halmov-Mammenov⁴⁵, Sonia F. Hashemi⁴⁶, Jennifer Heintzelman⁴⁷, Kristina Hogenboom⁴⁸, Farah M. Hosseini⁴⁹, Naveen Jain⁵⁰, Kuan Jen⁵¹, David Jorgenson⁵², Patrick Kannan⁵³, Ang Karapetyan⁵⁴, E. H. Kim⁵⁵, Ben Leinwand⁵⁶, Nasja Liu⁵⁷, Malte Möller⁵⁸, Andrew E. Mack⁵⁹, Mayank Matihajer⁶⁰, Noam Mansell⁶¹, Ondřej Nálež⁶², William Nowak⁶³, Daniel Mousavi-Garciá⁶⁴, Vicent Muñoz⁶⁵, Katerina Mavrogianni⁶⁶, Ahmad Mumtaz⁶⁷, Christopher Niel⁶⁸, Duncan J. Watts⁶⁹, Daniel Ort⁷⁰, Karen Ong⁷¹, Kay M. Pfeifer⁷², Michael Portnoy⁷³, Crystal Olear⁷⁴, Terekirat Rau⁷⁵, Anshu Sarpotdar⁷⁶, Therese Schaffner⁷⁷, London Schulze⁷⁸, Bryan Schoenfeld⁷⁹, Ben Sender⁸⁰, Jonathan D. Tang⁸¹, Emma Tsarkov⁸², Austin von Loos⁸³, Onur Varol⁸⁴, Xifei Wong⁸⁵, Zhi Wang⁸⁶, Julia Wang⁸⁷, Yuxia Wang⁸⁸, Kavindra Yakkunawala⁸⁹, Kirstie Whittlesey⁹⁰, Melia K. Winters⁹¹, Wei Lee Woon⁹², Barbara E. Engelhardt⁹³, Moritz Hardt⁹⁴, Deon Kross⁹⁵, Karen Levy⁹⁶, Arvind Narayanan⁹⁷, Brandon M. Stewart⁹⁸, Duncan J. Watts⁹⁹, and Sara McLanahan¹



Special Collection: Fragile Families Challenge

Introduction to the Special Collection on the Fragile Families Challenge

Matthew J. Salganik¹, Ian Lundberg², Alexander T. Kindel¹, and Sara McLanahan¹



Social Technologies Research for
a Dynamic World
Volume 1, 1–12
© The Author(s) 2018
Article reuse guidelines:
<http://sagepub.com/author-reuse.html>
DOI: 10.1177/2329004918761080
<http://journals.sagepub.com>
ISSN: 2329-0049

“Hack” existing scientific reward mechanisms: papers and prizes.

Measuring the predictability of life outcomes with a scientific mass collaboration

Matthew J. Salganik¹, Ian Lundberg², Alexander T. Kindell¹, Caitlin E. Ahern³, Khaled Al-Ghoreini⁴, Abdullah Almarzouq^{5,6}, Drew M. Altchul⁷, Jessie E. Brand⁸, Nicla Bohne Canegiela⁹, Ryan James Compton¹⁰, Debanjan Das¹¹, Thomas Davidson¹², Anna Filippova¹³, Connor Gilroy¹⁴, Brian J. Goode¹⁵, Eman Jahanri¹⁶, Rida Kashyap¹⁷, Amrit Kirchner¹⁸, Stephen McKay¹⁹, Allison C. Mersky²⁰, Alex Perner²¹, Karen Polens²², Heidi Raskin²³, Daniel Reiter²⁴, Michael V. Rosenbaum²⁵, Daniel M. Rossen²⁶, Adam Schaeffer²⁷, Erik H. Wang²⁸, Murva Adeny²⁹, Abdulkadir Altay³⁰, Bednarz Anis³¹, Ryan B. Argote³², Lilia Baer-Bouitif³³, Morris Badih³⁴, Bo-Byeon Chang³⁵, William Egger³⁶, Gregory Falsetti³⁷, Zhihs Fan³⁸, Jeremy Freese³⁹, Sajerwany Gade⁴⁰, Josh Gage⁴¹, Yuxi Guo⁴², Andrew Habermann-Maenner⁴³, Sonja F. Huskamp⁴⁴, Daniel J. Kishinevsky⁴⁵, Kristina Kishinevsky⁴⁶, Heng Li⁴⁷, Jason M. Hommerich⁴⁸, Naveen Jain⁴⁹, Kuan Jen⁵⁰, David Jorgenson⁵¹, Patrick Kannan⁵², Ang Karapetyan⁵³, E. H. Kim⁵⁴, Ben Leinwand⁵⁵, Nasja Liu⁵⁶, Molte Miser⁵⁷, Andrew E. Mack⁵⁸, Mayank Matihajra⁵⁹, Noah Mandell⁶⁰, Helge Marquart⁶¹, Daniel Meissner-Garcia⁶², Vicela Milicevic⁶³, Kathleen Maunder⁶⁴, Abraed Muller⁶⁵, Ondrejnik Nek⁶⁶, William Nowak⁶⁷, Jennifer Odegaard⁶⁸, Daniel Ort⁶⁹, Kamala Oberoi⁷⁰, Kayla Peacock⁷¹, Pauline Perner⁷², Crystal Olari⁷³, Terekirat Ratt⁷⁴, Anshut Sarpay⁷⁵, Theresia Schaffner⁷⁶, London Schulze⁷⁷, Bryan Schoenfeld⁷⁸, Ben Sender⁷⁹, Jonathan D. Tang⁸⁰, Emma Tsarkov⁸¹, Austin von Loon⁸², Onur Varol^{83,84}, Xufei Wong⁸⁵, Zhi Wang^{86,87}, Julia Wang⁸⁸, Yuxia Wang⁸⁹, Chandrani Venkatesan⁹⁰, Kirstie Whittle⁹¹, Melia K. Winters⁹², Wei-Lee Woser⁹³, Barbara Wu⁹⁴, Christopher Xie⁹⁵, Karen Yen⁹⁶, Bioggy Zhao⁹⁷, Chanyan Zhu⁹⁸, Jeannine Brooks-Gurin^{99,100}, Barbara E. Engelhardt¹⁰¹, Moritz Hardt¹⁰², Deon Kross¹⁰³, Karen Levy¹⁰⁴, Arvind Narayanan¹⁰⁵, Brandon M. Stewart¹⁰⁶, Duncan J. Watts^{107,108}, and Sara McLanahan¹



Special Collection: Fragile Families Challenge

Introduction to the Special Collection on
the Fragile Families Challenge



Social Technologies Research for
a Dynamic World
Volumes 1–11
© The Author(s) 2018
Article reuse guidelines:
<http://sagepub.com/page/authors/reuse.html>
DOI: 10.1177/2329004918761080
<http://journals.sagepub.com>
SAGE

Matthew J. Salganik¹, Ian Lundberg², Alexander T. Kindell¹, and Sara McLanahan¹

The image shows a screenshot of a video player interface. At the top, the title "FRAGILE FAMILIES CHALLENGE SCIENTIFIC WORKSHOP" is displayed in large, bold, black capital letters. Below the title, the date "NOVEMBER 17, 2017" is shown in a smaller, black font. In the bottom left corner of the video frame, there is a small logo of a shield with a crown on top. The bottom right corner of the video frame features a circular icon with a stylized letter "T". The overall background of the video player is white.



Privacy, Ethics, and Data Access: A Case Study of the Fragile Families Challenge

Socius: Sociological Research for a Dynamic World

Volume 5: 1–25

© The Author(s) 2019

Article reuse guidelines:

sagepub.com/journals-permissions

DOI: 10.1177/2378023118813023

srd.sagepub.com



Ian Lundberg¹ , Arvind Narayanan¹, Karen Levy²,
and Matthew J. Salganik¹

Privacy, Ethics, and Data Access: A Case Study of the Fragile Families Challenge

Socius: Sociological Research for a Dynamic World
Volume 5: 1–25

© The Author(s) 2019

Article reuse guidelines:
sagepub.com/journals-permissions
DOI: 10.1177/2378023118813023
srd.sagepub.com



Ian Lundberg¹ , Arvind Narayanan¹, Karen Levy²,
and Matthew J. Salganik¹

Improving Metadata Infrastructure for Complex Surveys: Insights from the Fragile Families Challenge

Socius: Sociological Research for a Dynamic World
Volume 5: 1–24

© The Author(s) 2019

Article reuse guidelines:
sagepub.com/journals-permissions
DOI: 10.1177/2378023118817378
srd.sagepub.com



Alexander T. Kindel¹, Vineet Bansal¹, Kristin D. Catena¹,
Thomas H. Hartshorne¹, Kate Jaeger¹, Dawn Koffman¹,
Sara McLanahan¹, Maya Phillips¹, Shiva Rouhani¹, Ryan Vinh¹,
and Matthew J. Salganik¹

Privacy, Ethics, and Data Access: A Case Study of the Fragile Families Challenge

Socius: Sociological Research for a Dynamic World

Volume 5: 1–25

© The Author(s) 2019

Article reuse guidelines:

sagepub.com/journals-permissions

DOI: 10.1177/2378023118813023

srd.sagepub.com

Ian Lundberg¹ , Arvind Narayanan¹, Karen Levy²,
and Matthew J. Salganik¹

Improving Metadata Infrastructure for Complex Surveys: Insights from the Fragile Families Challenge

Socius: Sociological Research for a Dynamic World

Volume 5: 1–24

© The Author(s) 2019

Article reuse guidelines:

sagepub.com/journals-permissions

DOI: 10.1177/2378023118817378

srd.sagepub.com

Alexander T. Kindel¹, Vineet Bansal¹, Kristin D. Catena¹,
Thomas H. Hartshorne¹, Kate Jaeger¹, Dawn Koffman¹,
Sara McLanahan¹, Maya Phillips¹, Shiva Rouhani¹, Ryan Vinh¹,
and Matthew J. Salganik¹

Successes and Struggles with Computational Reproducibility: Lessons from the Fragile Families Challenge

Socius: Sociological Research for a Dynamic World

Volume 5: 1–21

© The Author(s) 2019

Article reuse guidelines:

sagepub.com/journals-permissions

DOI: 10.1177/2378023119849803

srd.sagepub.com

Competition vs collaboration

- ▶ We can learn from looking at the full set of submissions not just the best

Competition vs collaboration

- ▶ We can learn from looking at the full set of submissions not just the best
- ▶ We can open source everything to seed future research

Competition vs collaboration

- ▶ We can learn from looking at the full set of submissions not just the best
- ▶ We can open source everything to seed future research
- ▶ Just feels right to me

Break into discussion groups

Getting ready for Thursday

Thoughts about reading the interviews

- ▶ Starts simple, gets more personal. This is what happens when you talk to people.

Thoughts about reading the interviews

- ▶ Starts simple, gets more personal. This is what happens when you talk to people.
- ▶ Difficulty of calibration.

Thoughts about reading the interviews

- ▶ Starts simple, gets more personal. This is what happens when you talk to people.
- ▶ Difficulty of calibration.
- ▶ What is the goal of the reading? Filling out the sheet and helping us make sense of the results from the Challenge

Thoughts about reading the interviews

- ▶ Starts simple, gets more personal. This is what happens when you talk to people.
- ▶ Difficulty of calibration.
- ▶ What is the goal of the reading? Filling out the sheet and helping us make sense of the results from the Challenge
- ▶ If you don't have access yet, let me know

Thoughts about reading the interviews

- ▶ Starts simple, gets more personal. This is what happens when you talk to people.
- ▶ Difficulty of calibration.
- ▶ What is the goal of the reading? Filling out the sheet and helping us make sense of the results from the Challenge
- ▶ If you don't have access yet, let me know
- ▶ We will have guess on Thursday from the FFCWS data team and the dark matter interview team

Thoughts about reading the interviews

- ▶ Starts simple, gets more personal. This is what happens when you talk to people.
- ▶ Difficulty of calibration.
- ▶ What is the goal of the reading? Filling out the sheet and helping us make sense of the results from the Challenge
- ▶ If you don't have access yet, let me know
- ▶ We will have guess on Thursday from the FFCWS data team and the dark matter interview team

Class slides for Tuesday, November 10: Predictability of life trajectories

Matthew J. Salganik

COS 597E/SOC 555 Limits to prediction
Fall 2020, Princeton University