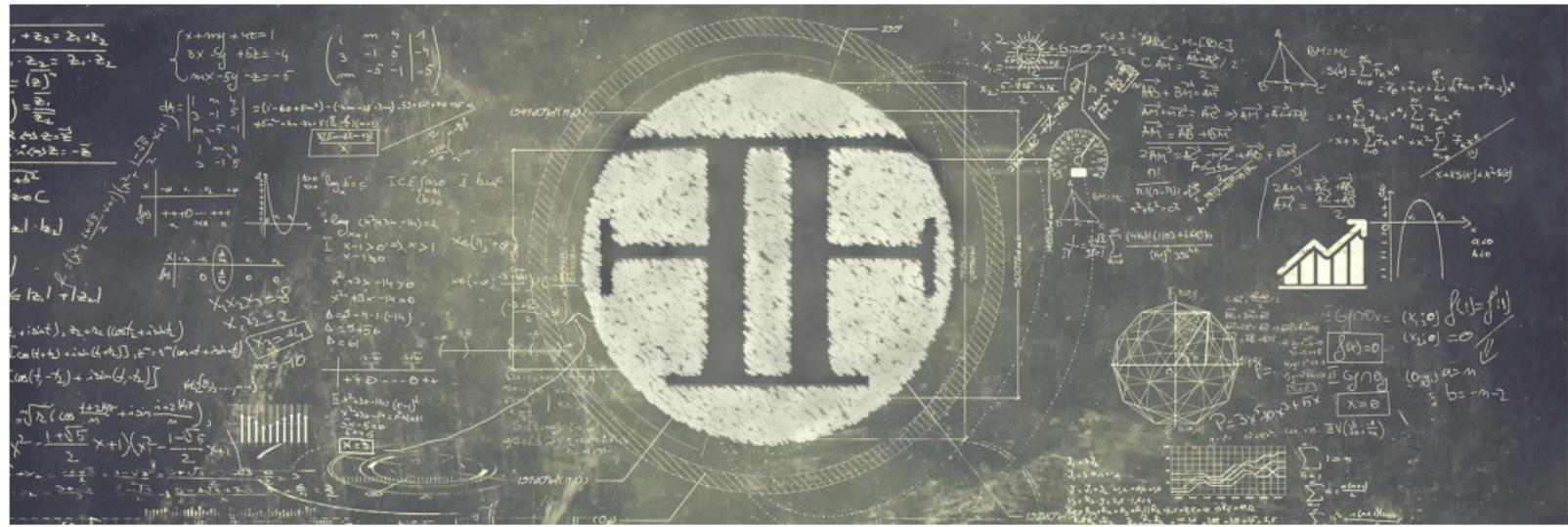


Class slides for Tuesday, November 10: Predictability of life trajectories

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)

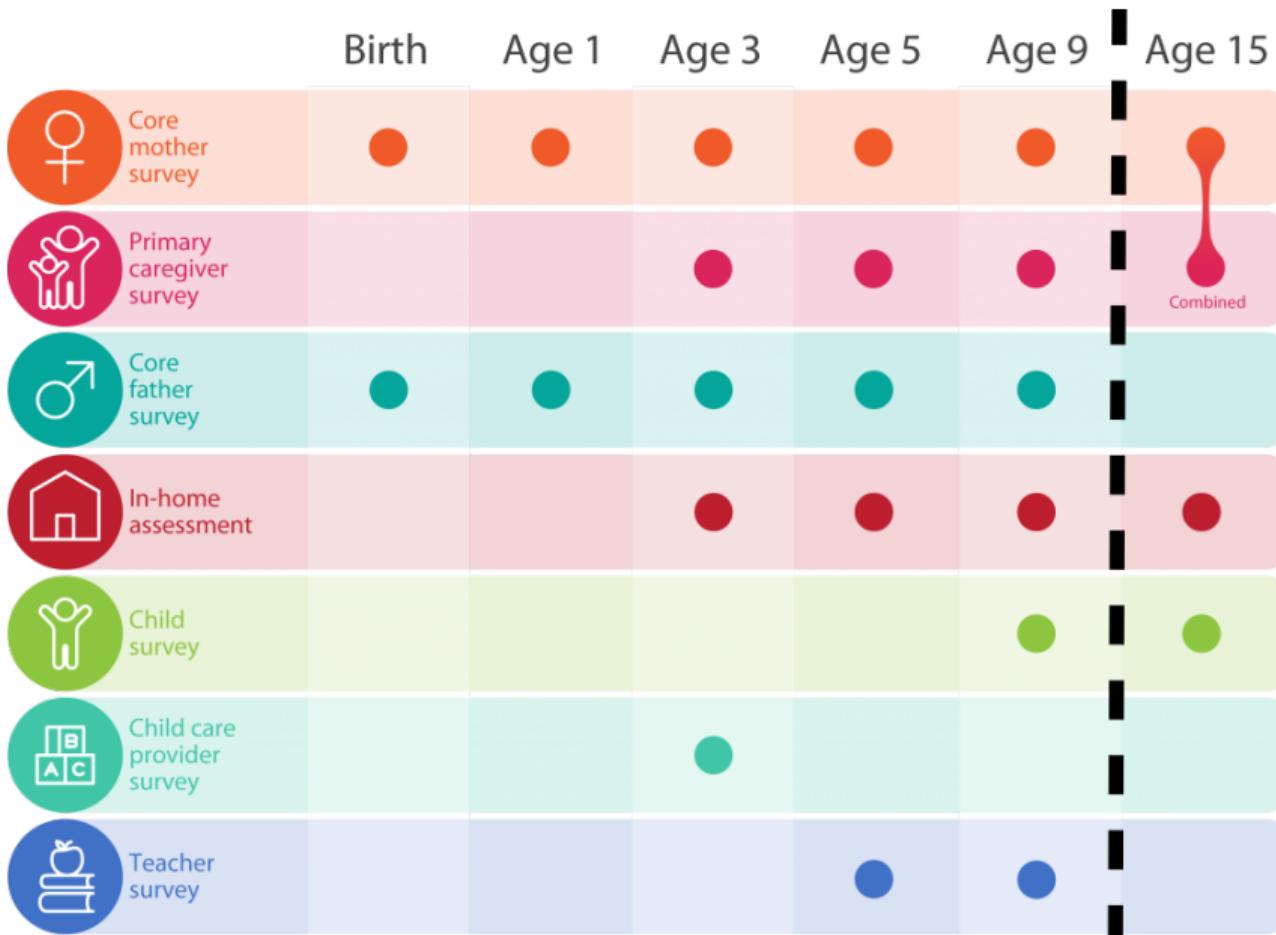
FF Fragile Families

& 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

	Birth	Age 1	Age 3	Age 5	Age 9
 Core mother survey	●	●	●	●	●
 Primary caregiver survey			●	●	●
 Core father survey	●	●	●	●	●
 In-home assessment			●	●	●
 Child survey					●
 Child care provider survey			●		
 Teacher survey				●	●



4,242 families

Birth to age 9
12,942 features

Age 15
1,500 features

Birth to age 9
12,942 variables

Age 15
6 variables

4,242 families

Information about child and family

Background data

Training

Leaderboard

Holdout

Outcome
data

This is not inferring the present or forecasting. What is it?

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?

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

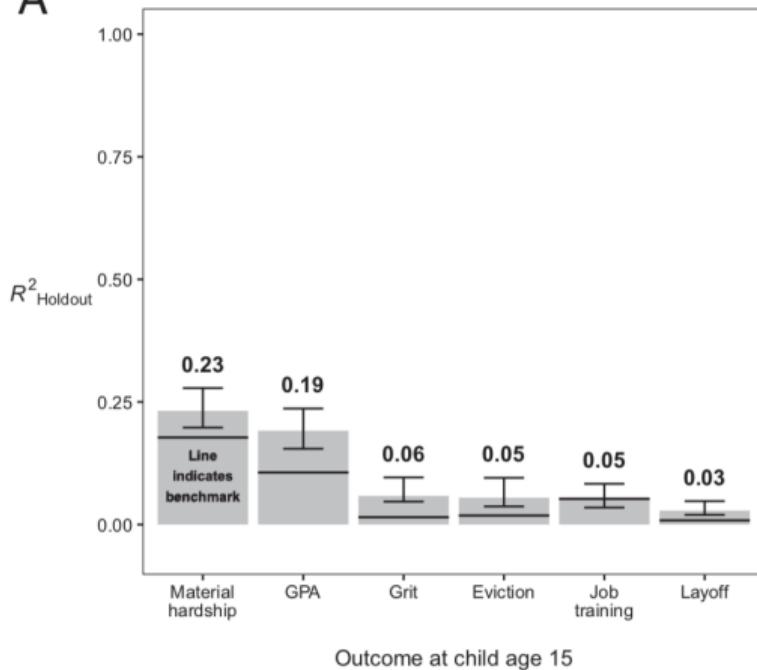
Using a large, high-quality social science dataset collected since birth and modern machine learning methods, how accurately can we predict outcomes from children, parents, and families?

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

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

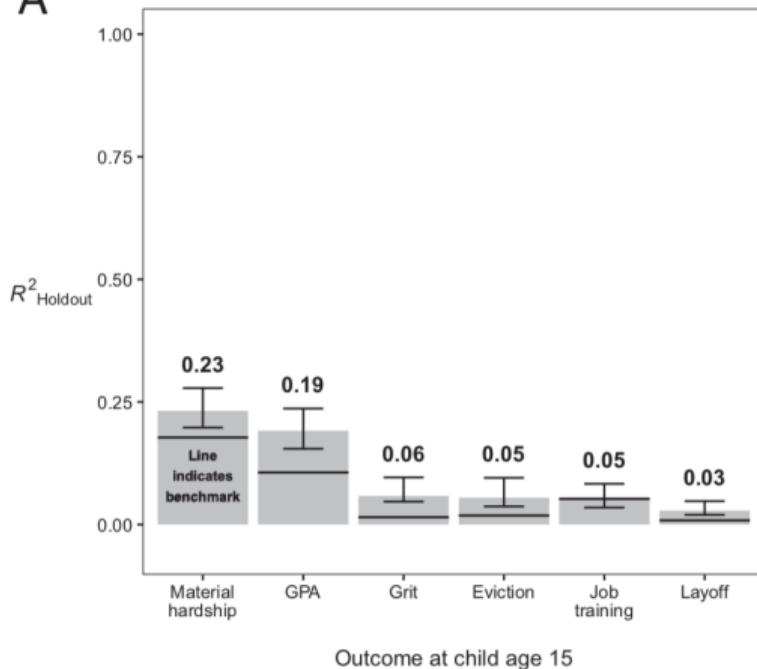
A

Best submission for each outcome

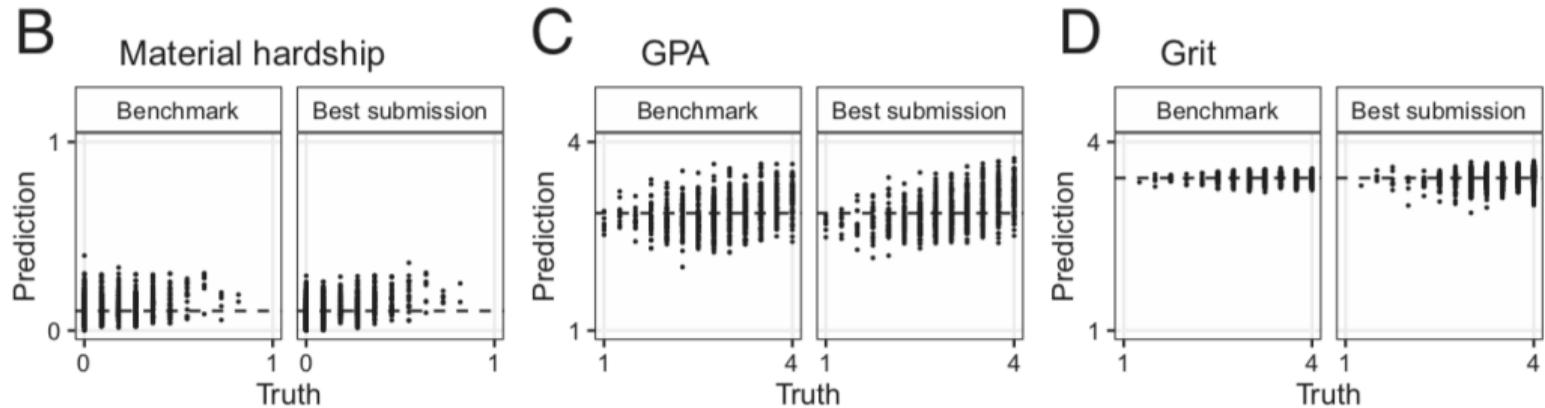


A

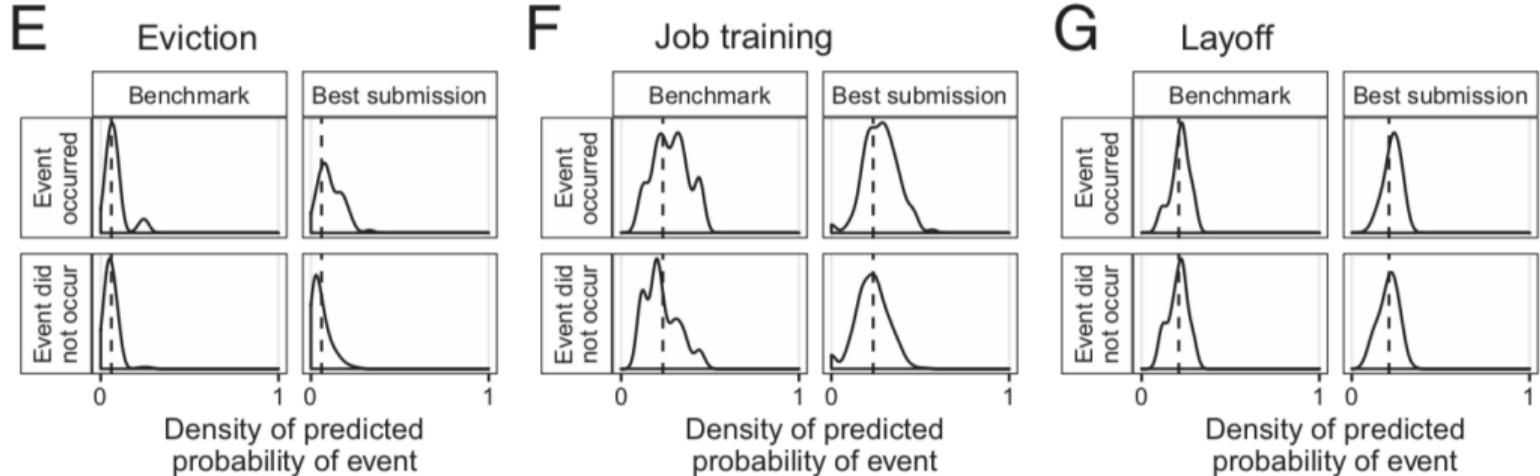
Best submission for each outcome



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



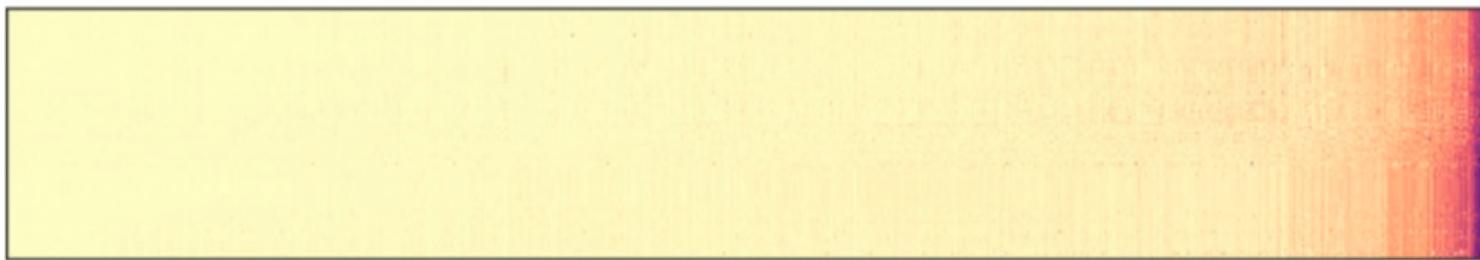
- ▶ Best submission is like the mean of the training data with a bit of signal, not like the truth with a bit of noise
- ▶ Best submission is qualitatively similar to the benchmark model

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

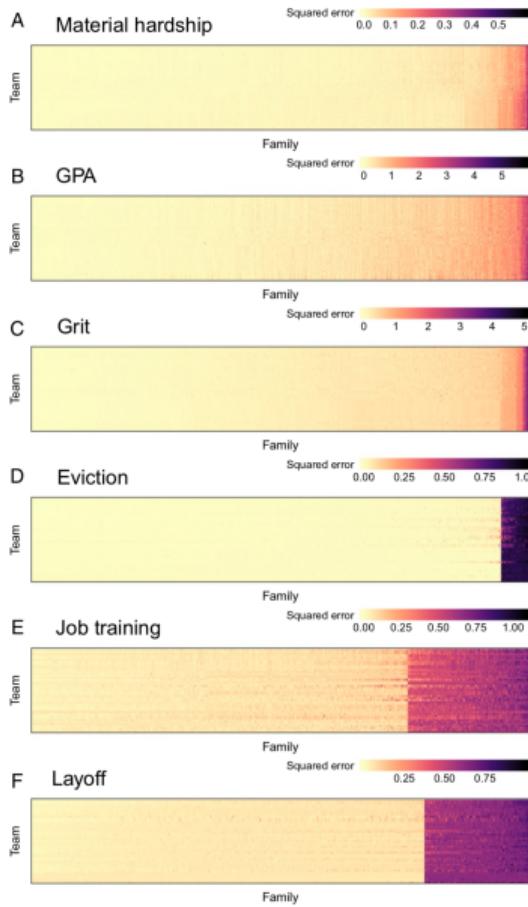
A Material hardship

Squared error
0.0 0.1 0.2 0.3 0.4 0.5

Team



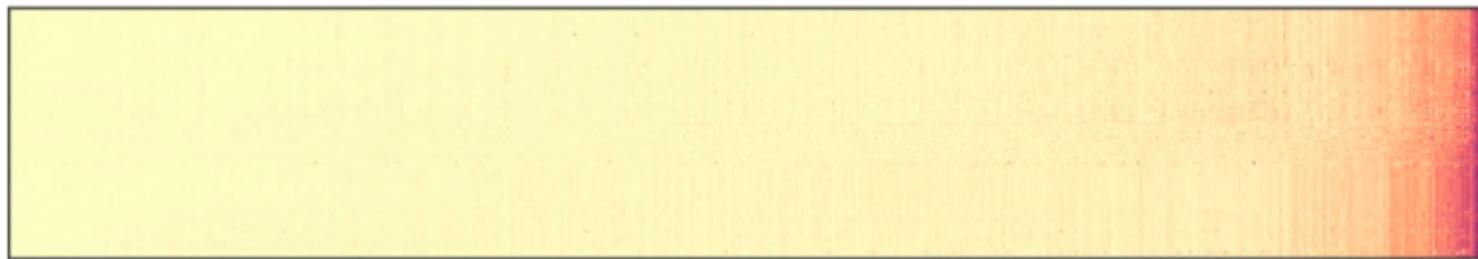
Family



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. Salganik^{1,2}, Ian Lundberg^{1,3}, Alexander T. Kindred¹, Caitlin E. Ahern², Khaled Al-Ghoreini², Abdullah Almarzouq^{2,4}, Drew M. Altchul⁵, Jessie E. Brand⁶, Nicla Bohne Canegiela⁷, Ryan James Compton¹, Debanjan Datta⁸, Thomas Davidson⁹, Anna Filippova¹⁰, Connor Gilroy¹¹, Brian J. Goode¹², Eman Jahanri¹³, Rida Kashyap¹⁴, Arif Kirchner¹⁵, Stephen McKay¹⁶, Allison C. Merseth¹⁷, Alex Pechard¹⁸, Karen Polens¹⁹, Heidi Price²⁰, Dovile Ruzicka²¹, Michael V. Salter²², Daniel M. Selsky²³, Valeria Sestini²⁴, Adriana Ureña²⁵, Erik H. Wang²⁶, Murva Adeny²⁷, Abdulkadir Altaraj²⁸, Bednar Alshehri²⁹, Redwan Ansari³⁰, Ryan R. Argote³¹, Lilia Baer-Bouit³², Morton Bodzin³³, Bo-Byeon Chang³⁴, William Egger³⁵, Gregory Falsetto³⁶, Zhilan Fan³⁷, Jeremy Freese³⁸, Sergey Gulyayev³⁹, Josh Gage⁴⁰, Yuxi Guo⁴¹, Andrey Habermann-Mammen⁴², Sonja F. Hassler⁴³, Jennifer Heidrich⁴⁴, Kristina Hogenboom⁴⁵, Farah M. Hosseini⁴⁶, Daniel Hwang⁴⁷, Naveen Jain⁴⁸, Kuan Jen⁴⁹, David Jurgens⁵⁰, Patrick Kannisto⁵¹, Ang Karapetyan⁵², E. H. Kim⁵³, Ben Leinwand⁵⁴, Nasja Liu⁵⁵, Malte Möller⁵⁶, Andrew E. Mack⁵⁷, Mayank Matihajer⁵⁸, Noa Mandell⁵⁹, Helge Marathane⁶⁰, Olafur Mousavi-Garcia⁶¹, Vicenç Muñoz⁶², Katharina Mutschler⁶³, Ahmad Mumtaz⁶⁴, Ondjoum Ndi⁶⁵, William Nowak⁶⁶, Jennifer O’Donnell⁶⁷, Daniel Ort⁶⁸, Karen Oren⁶⁹, Kaye McPherson⁷⁰, Kristin Portnoy⁷¹, Crystal Olear⁷², Terekirat Rau⁷³, Anshu Sargyan⁷⁴, Theresia Schaffner⁷⁵, London Schubert⁷⁶, Bryan Schoenfeld⁷⁷, Ben Sender⁷⁸, Jonathan D. Tang⁷⁹, Emma Tsourkoy⁸⁰, Austin von Loon⁸¹, Onur Varol⁸², Xiaofei Wong⁸³, Zhi Wang⁸⁴, Julia Wang⁸⁵, Yuxia Wang⁸⁶, Chandrika Venkateswaran⁸⁷, Kirstie Whittaker⁸⁸, Melia K. Winters⁸⁹, Wei Lee Worcester⁹⁰, Daniel Wu⁹¹, Christopher Yang⁹², Junwei Yin⁹³, Bioggy Zhao⁹⁴, Chanyan Zhu⁹⁵, Jeanne Brooks-Gunn⁹⁶, Barbara E. Engelhardt⁹⁷, Morris Hardt⁹⁸, Deon Kross⁹⁹, Karen Levy¹⁰⁰, Arvind Narayanan¹⁰¹, Preston M. Stewart¹⁰², Duncan J. Watts¹⁰³, and Sara McLanahan¹⁰⁴.

“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, Bedder Alshehri³², Redwan Ansari³³, Ryan R. Argote³⁴, Lilia Baer-Bouckaert³⁵, Morris Badih³⁶, Bo-Byeon Chang³⁷, William Egger³⁸, Gregory Falsetti³⁹, Zhilan Fan⁴⁰, Jeremy Freese⁴¹, Sergey Gulyaev⁴², Josh Gage⁴³, Yuxi Guo⁴⁴, Andrey Halmov-Mammenov⁴⁵, Sonia F. Hashemi⁴⁶, Jennifer Heintzelman⁴⁷, Kristina Hogenboom⁴⁸, Farah M. Hosseini⁴⁹, Naseem Ihsan⁵⁰, Naveen Jain⁵¹, Koen Jans⁵², David Jorgenson⁵³, Patrick Kannikka⁵⁴, Ang Karapetyan⁵⁵, E. H. Kim⁵⁶, Ben Leinwand⁵⁷, Nasja Liu⁵⁸, Malte Möller⁵⁹, Andrew E. Mack⁶⁰, Mayank Matihajer⁶¹, Noam Mansell⁶², Helge Marathane⁶³, Olafur Melsas-Garcia⁶⁴, Vicent Mir⁶⁵, Katerina Mavroudi⁶⁶, Ahmad Mumtaz⁶⁷, Ondrejuk Nek⁶⁸, William Nowak⁶⁹, Dennis O’Hearn⁷⁰, Daniel Ort⁷¹, Karen Otero⁷², Kayla Peacock⁷³, Paul Portney⁷⁴, Crystal Olear⁷⁵, Terekirat Rau⁷⁶, Anshu Sarpotdar⁷⁷, Therese Schaffner⁷⁸, London Schubert⁷⁹, Bryant Schoenfeld⁸⁰, Ben Sender⁸¹, Jonathan D. Tang⁸², Emma Tsarkov⁸³, Austin von Loese⁸⁴, Onur Varol⁸⁵, Xifei Wong⁸⁶, Zhi Wang⁸⁷, Julia Wang⁸⁸, Yuxia Wang⁸⁹, Kavindra Yakkunawala⁹⁰, Kirstie Whittlesey⁹¹, Melia K. Winters⁹², Wei Lee Woon⁹³, Barbara E. Engelhardt⁹⁴, Moritz Hardt⁹⁵, Deon Kross⁹⁶, Karen Levy⁹⁷, Arvind Narayanan⁹⁸, Brandon M. Stewart⁹⁹, Duncan J. Watts¹⁰⁰, and Sara McLanahan¹⁰¹



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: 10.1177/2329024618761080
<http://journals.sagepub.com>
ISSN: 2329-0246

“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-Gutierrez⁹⁹, Barbara E. Engelhardt¹⁰⁰, Moritz Hardt¹⁰¹, Deon Kross¹⁰², Karen Levy¹⁰³, Arvind Narayanan¹⁰⁴, Brandon M. Stewart¹⁰⁵, Duncan J. Watts^{106,107}, 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. Kindell¹, and Sara McLanahan¹



Social Technologies Research for

a Dynamic World

Volumes 1–11

© The Author(s) 2018

Article reuse guidelines:

http://sagepub.com/page/authors/reuse.html

ISSN: 1529-1016 (print)

ISSN: 1529-1024 (electronic)

DOI: 10.1177/1529101618760200

sagepub.com



The thumbnail shows a presentation slide with the title "FRAGILE FAMILIES CHALLENGE SCIENTIFIC WORKSHOP". Below the title is a date: NOVEMBER 17, 2017. The slide features a shield logo on the left and a small image of a dollar bill on the right. The overall background is white with some dark shadows at the bottom.





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.
- ▶ Year 15 interview GPA refers to 9th grade for both young adults.

Thoughts about reading the interviews

- ▶ Starts simple, gets more personal. This is what happens when you talk to people.
- ▶ Year 15 interview GPA refers to 9th grade for both young adults.
- ▶ Difficulty of calibration.

Thoughts about reading the interviews

- ▶ Starts simple, gets more personal. This is what happens when you talk to people.
- ▶ Year 15 interview GPA refers to 9th grade for both young adults.
- ▶ 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.
- ▶ Year 15 interview GPA refers to 9th grade for both young adults.
- ▶ 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.
- ▶ Year 15 interview GPA refers to 9th grade for both young adults.
- ▶ 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.
- ▶ Year 15 interview GPA refers to 9th grade for both young adults.
- ▶ 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