# Projecting the outcomes of people's lives with AI isn't so simple

March 30 2020



|  | Birth | Age 1 | Age 3 | Age 5 | Age 9 | Age 15 |
|---|---|---|---|---|---|---|
| Core mother survey | ● | ● | ● | ● | ● | ● |
| Primary caregiver survey |  |  | ● | ● | ● | ● Combined |
| Core father survey | ● | ● | ● | ● | ● |  |
| In-home assessment |  |  | ● | ● | ● | ● |
| Child survey |  |  |  |  | ● | ● |
| Child care provider survey |  |  | ● |  |  |  |
| Teacher survey |  |  |  | ● | ● |  |

The Fragile Families study captured information on children at birth and ages 1, 3, 5, 9 and 15. This information was captured through a variety of surveys, listed to the left of these ages in the above chart. The Fragile Families Challenge used data from waves one to five to predict outcomes in wave six. Credit: Matthew Salganik et al. 2020, Princeton University

The machine learning techniques scientists use to predict outcomes from large datasets may fall short when it comes to projecting the outcomes of people's lives, according to a mass study led by researchers at Princeton University in a collaboration with researchers across many institutions, including Virginia Tech.

This mass collaboration, called the Fragile Families Challenge, represents a cohort of scientists that build statistical and machine-learning models to predict and measure life outcomes for children, parents, and households across the United States.

Published by 112 co-authors in the *Proceedings of the National Academy of Sciences*, the results suggest that sociologists and data scientists should use caution when using predictive modeling, especially in the criminal justice system and social programs.

Even after using state-of-the-art modeling and a high-quality dataset containing 13,000 data points for more than 4,000 families, the best AI predictive models were not very accurate.

Brian J. Goode, a research scientist from Virginia Tech's Fralin Life Sciences Institute, was among the data and social scientists that were in the Fragile Families Challenge.
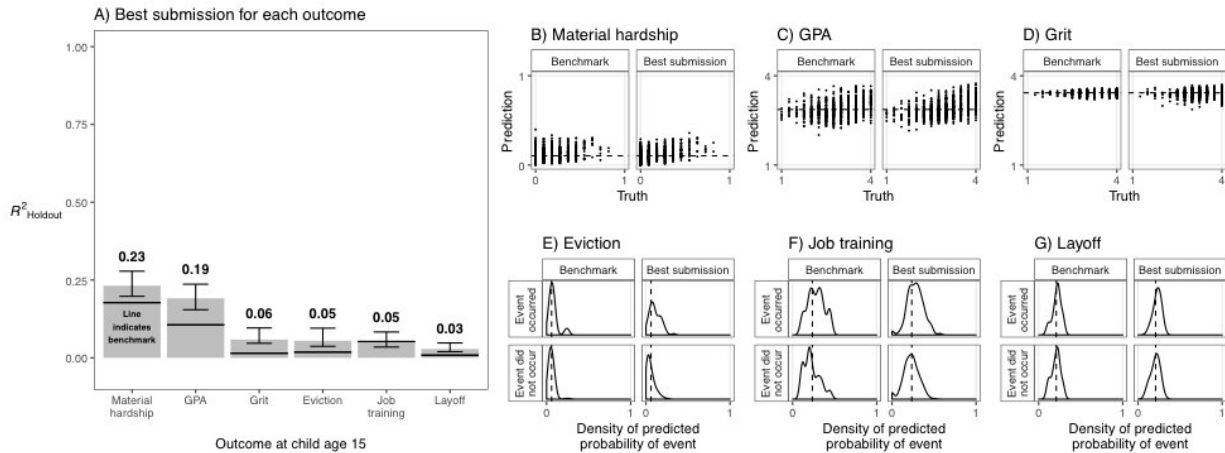
Figure A shows the difference between the best submissions for each outcome compared to the benchmark model. Figure B-G compared the predictions and the truth for each outcome. Credit: Matthew Salganik et al. 2020, Princeton University

"It's one effort to try to capture the complexities and intricacies that compose the fabric of a human life in data and models. But, it is compulsory to take the next step and contextualize models in terms of how they are going to be applied in order to better reason about expected uncertainties and limitations of a prediction. That's a very difficult problem to grapple with, and I think the Fragile Families Challenge shows that we need more research support in this area, particularly as machine learning has a greater impact on our everyday lives," said Goode.Goode's modeling was conducted through the Discovery Analytics Center at Virginia Tech. There, he teamed up with the Discovery Analytics Center's director and the Thomas L. Phillips Professor of Engineering, Naren Ramakrishnan, and Debanjan Datta, a Ph.D. student in the Department of Computer Science in the College of Engineering, who were instrumental in gathering and analyzing data.

The Virginia Tech team has also published research in a special issue of

Socius, a new open-access journal from the American Sociological Association. In order to support additional research in this area, all the submissions to the Challenge—code, predictions and narrative explanations—are publicly available.

"The study also shows us that we have so much to learn, and mass collaborations like this are hugely important to the research community," said the PNAS study co-lead author Matt Salganik, professor of sociology at Princeton and interim director of the Center for Information Technology Policy, based at Princeton's Woodrow Wilson School of Public and International Affairs.
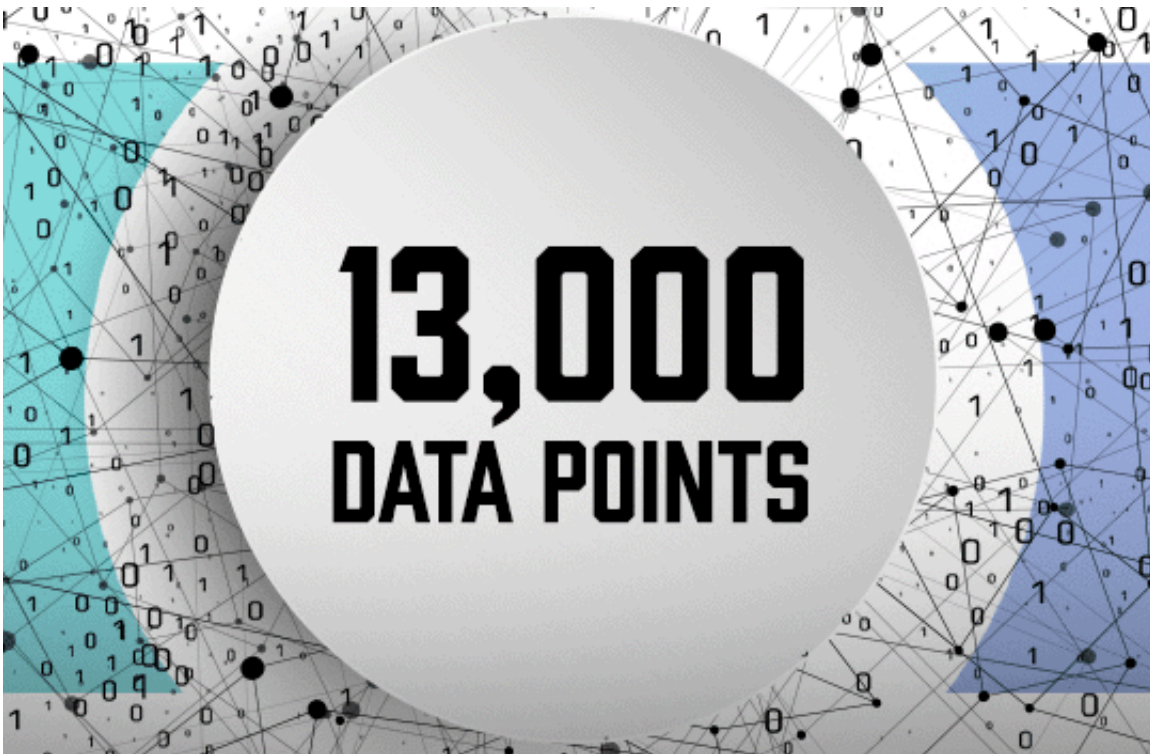
The project was inspired by Wikipedia, one of the world's first mass collaborations, which was created in 2001 as a shared encyclopedia. Salganik pondered what other scientific problems could be solved through a new form of collaboration, and that's when he joined forces with Sara McLanahan, the William S. Tod Professor of Sociology and Public Affairs at Princeton, as well as Princeton graduate students Ian Lundberg and Alex Kindel, both in the Department of Sociology.

McLanahan is principal investigator of the Fragile Families and Child Wellbeing Study based at Princeton and Columbia University, which has been studying a cohort of about 5,000 children born in large American cities between 1998 and 2000, with an oversampling of children born to unmarried parents. The longitudinal study was designed to understand the lives of children born into unmarried families.

Through surveys collected in six waves (when the child was born and then when the child reached ages 1, 3, 5, 9, and 15), the study has captured millions of data points on children and their families. Another wave will be captured at age 22.

At the time the researchers designed the challenge, data from age 15

(which the researchers call in the paper the "hold-out data) had not yet been made publicly available. This created an opportunity to ask other scientists to predict life outcomes of the people in the study through a mass collaboration.



160 research teams of data and social scientists built statistical and machine-learning models to predict measure six life outcomes for children, parents, and households. Even after using a state-of-the-art modeling and a high-quality dataset containing 13,000 data points about more than 4,000 families, the best AI predictive models were not very accurate. Credit: Egan Jimenez, Princeton University

The co-organizers received 457 applications from 68 institutions from around the world, including from several teams based at Princeton. Using the Fragile Families data, participants were asked to predict one or

more of the six life outcomes at age 15. These included child grade point average (GPA); child grit; household eviction; household material hardship; primary caregiver layoff; and primary caregiver participation in job training.

The challenge was based around the common task method, a research design used frequently in computer science but not in the social sciences. This method releases some but not all of the data, allowing people to use whatever technique they want to determine outcomes. The goal is to accurately predict the hold-out data, no matter how fancy a technique it takes to get there.

The team is currently applying for grants to continue research in this area.

The paper, "Measuring the predictability of life outcomes with a scientific mass collaboration," was published on March 30 by *PNAS*.

**More information:** Matthew J. Salganik el al., "Measuring the predictability of life outcomes with a scientific mass collaboration," *PNAS* (2020). www.pnas.org/cgi/doi/10.1073/pnas.1915006117

Provided by Virginia Tech