

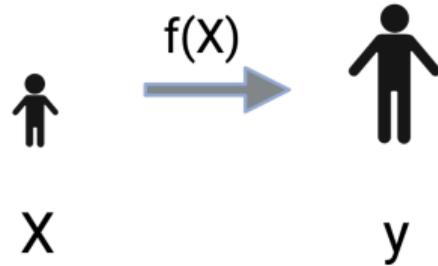
# HW2: Fragile Families Challenge

COS424/524 precept

Slides adapted from Prof. Matt Salganik's [SICSS materials](#).



Given data about a young child's circumstances, are outcomes later in life predictable?



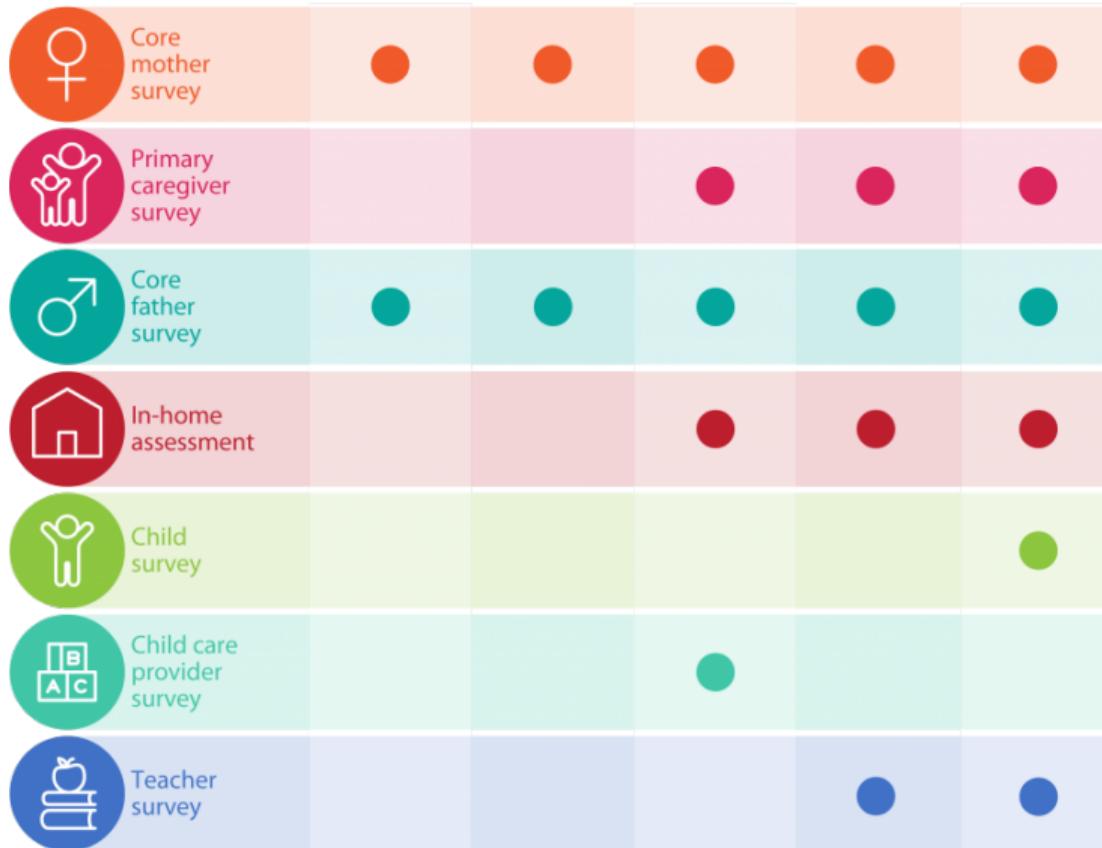
# FF Fragile Families

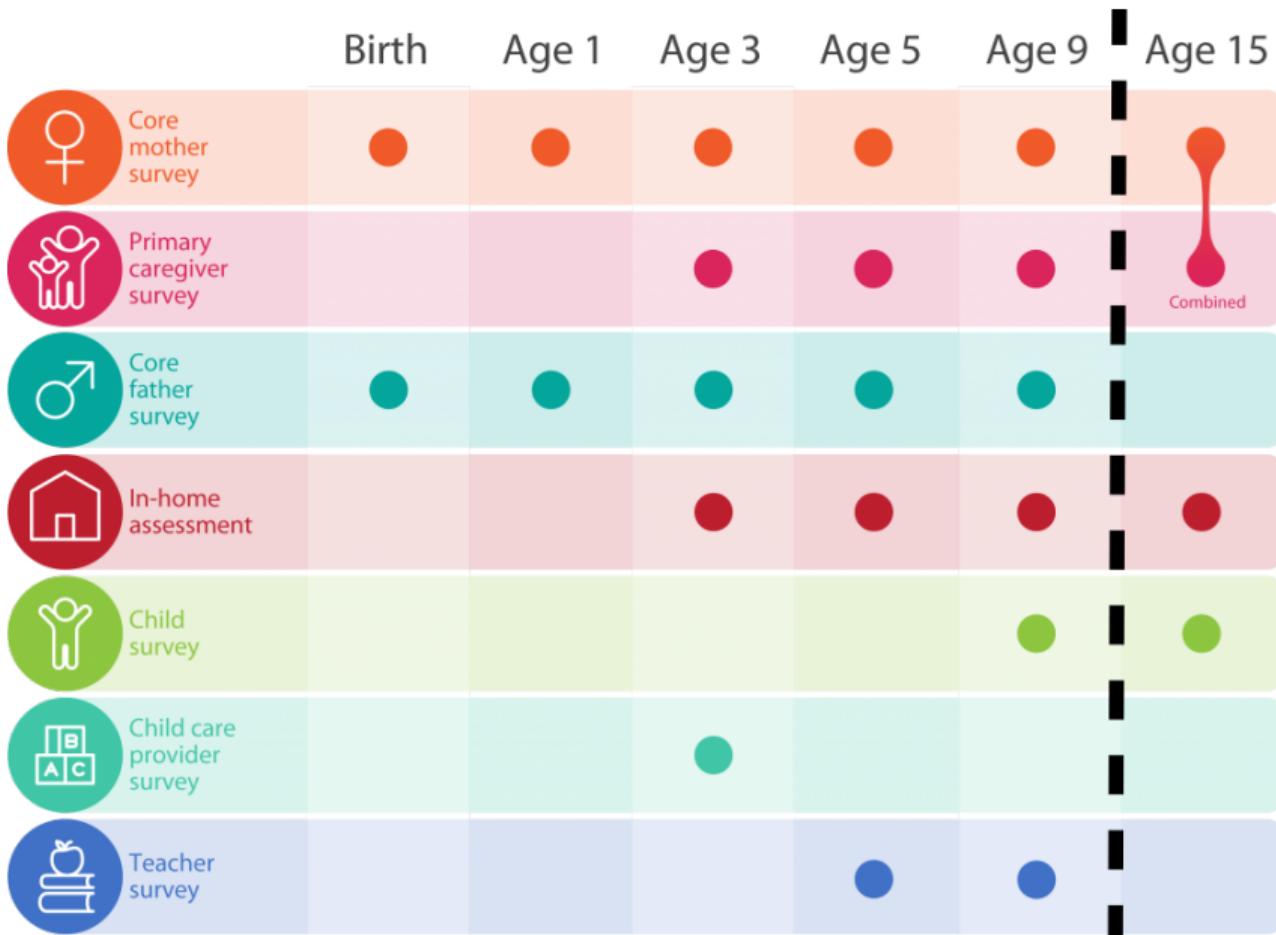
& Child Wellbeing Study  
PRINCETON | COLUMBIA



- ▶ 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
- ▶ “Fragile” because they’re at a greater risk of breaking up and living in poverty than more traditional families

Birth      Age 1      Age 3      Age 5      Age 9





## Example questions

*Mother:* How many days per week do you read stories to your child?

*Father:* How many days per week do you hug or show physical affection to your child?

*Child:* On average, how much time do you spend hanging out with friends on a weekday?

Life course researchers have:

- ▶ described social patterns
- ▶ theorized important factors
- ▶ estimated causal effects

Life course researchers have:

- ▶ described social patterns
- ▶ theorized important factors
- ▶ estimated causal effects

How well can we predict individual life outcomes?

We should care about the predictability of social outcomes

We should care about the predictability of social outcomes

- ▶ Scientific reasons

We should care about the predictability of social outcomes

- ▶ Scientific reasons
  - ▶ Basic social fact

We should care about the predictability of social outcomes

- ▶ Scientific reasons

- ▶ Basic social fact
- ▶ Discovery

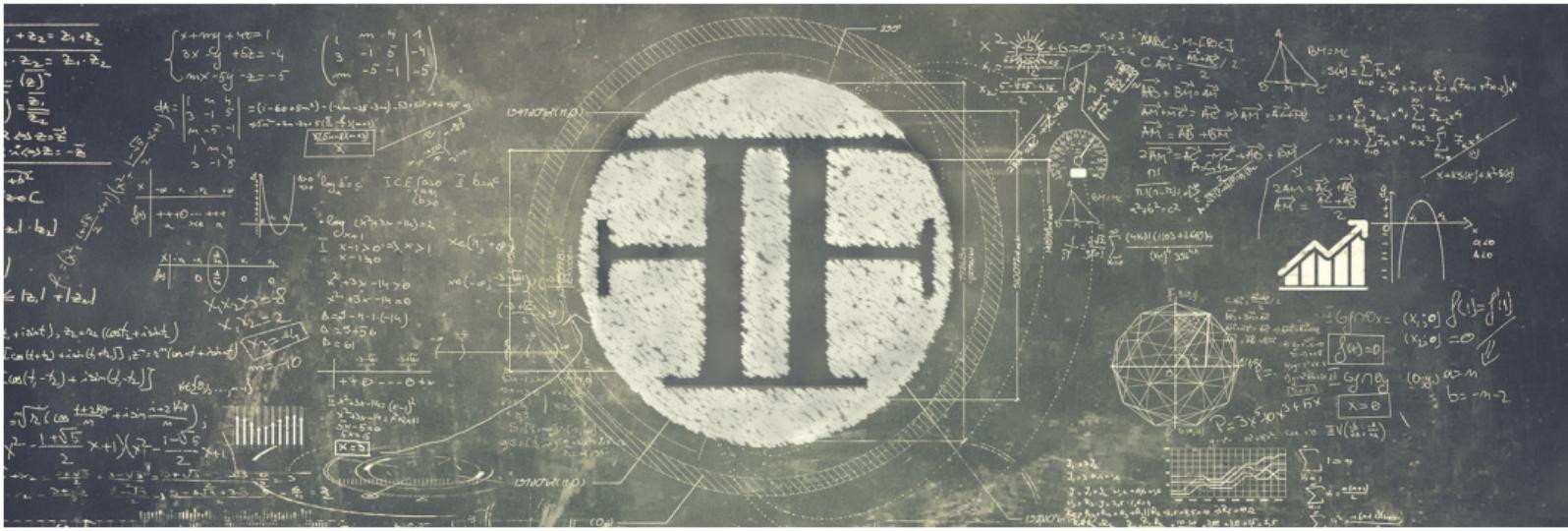
We should care about the predictability of social outcomes

- ▶ Scientific reasons

- ▶ Basic social fact
- ▶ Discovery

- ▶ Policy reasons



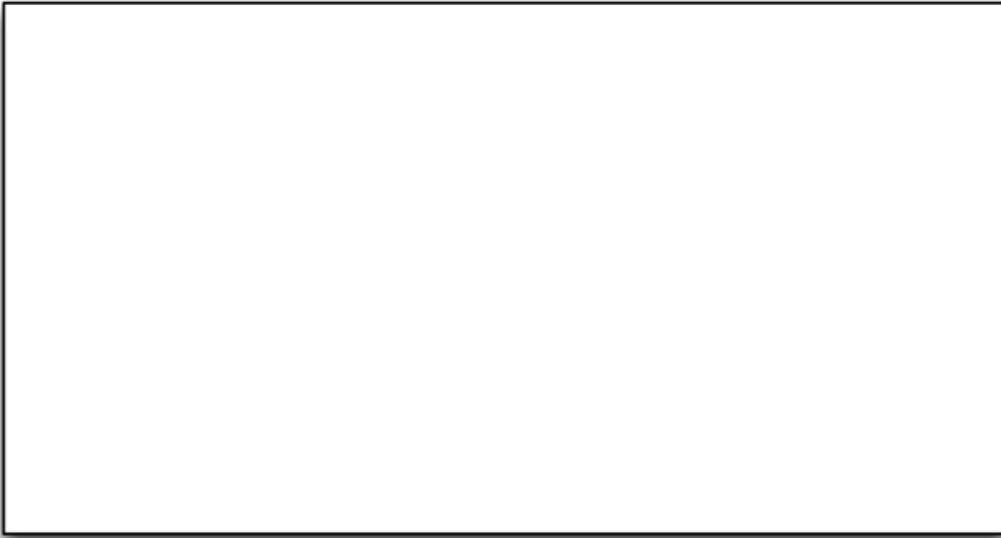


## Fragile Families Challenge

5,000 families

Birth to age 9  
12,000 features

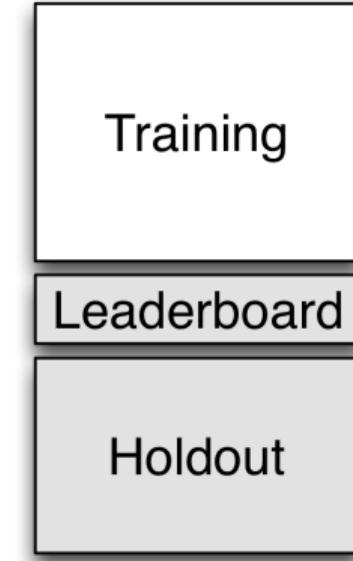
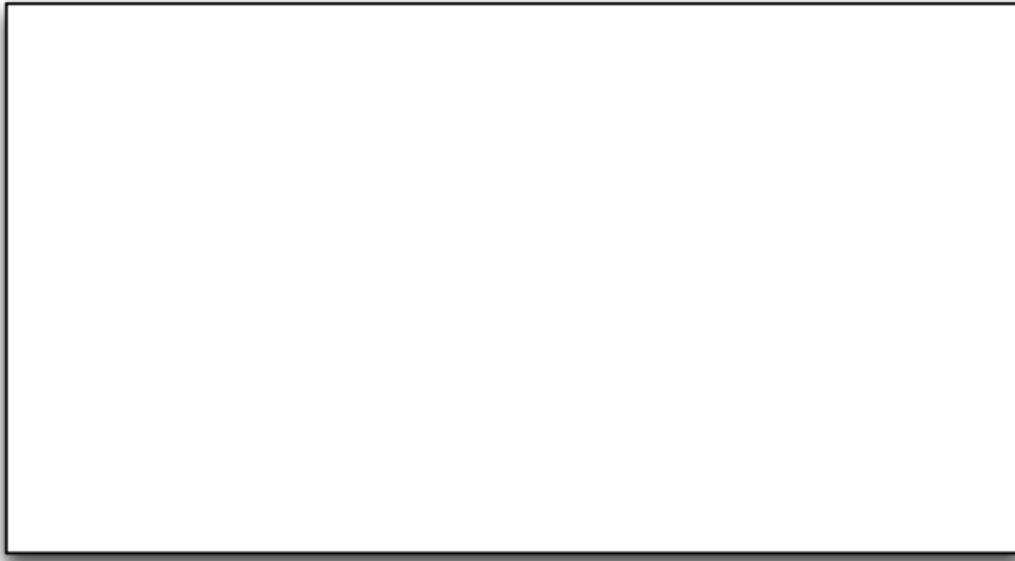
Age 15  
1,500 features



4,242 families

12,942 features  
birth to age 9

6 outcomes  
age 15

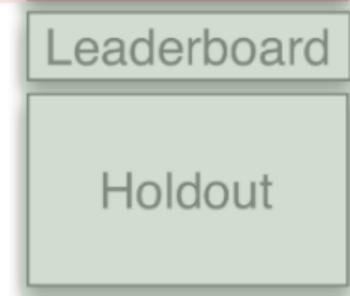
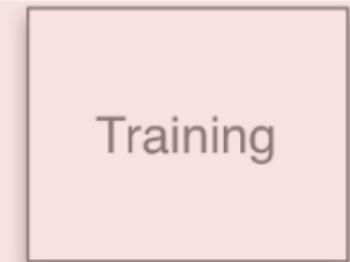


You have access  
 You don't have access

12,942 features  
birth to age 9

6 outcomes  
age 15

4,242 families



## Covariates

- ▶ ~ 13,000 features from questionnaires

## Outcomes

- ▶ GPA
- ▶ Grit
- ▶ Eviction
- ▶ Material hardship
- ▶ Job training
- ▶ Job loss

## Outcomes

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

# Measuring the predictability of life outcomes with a scientific mass collaboration

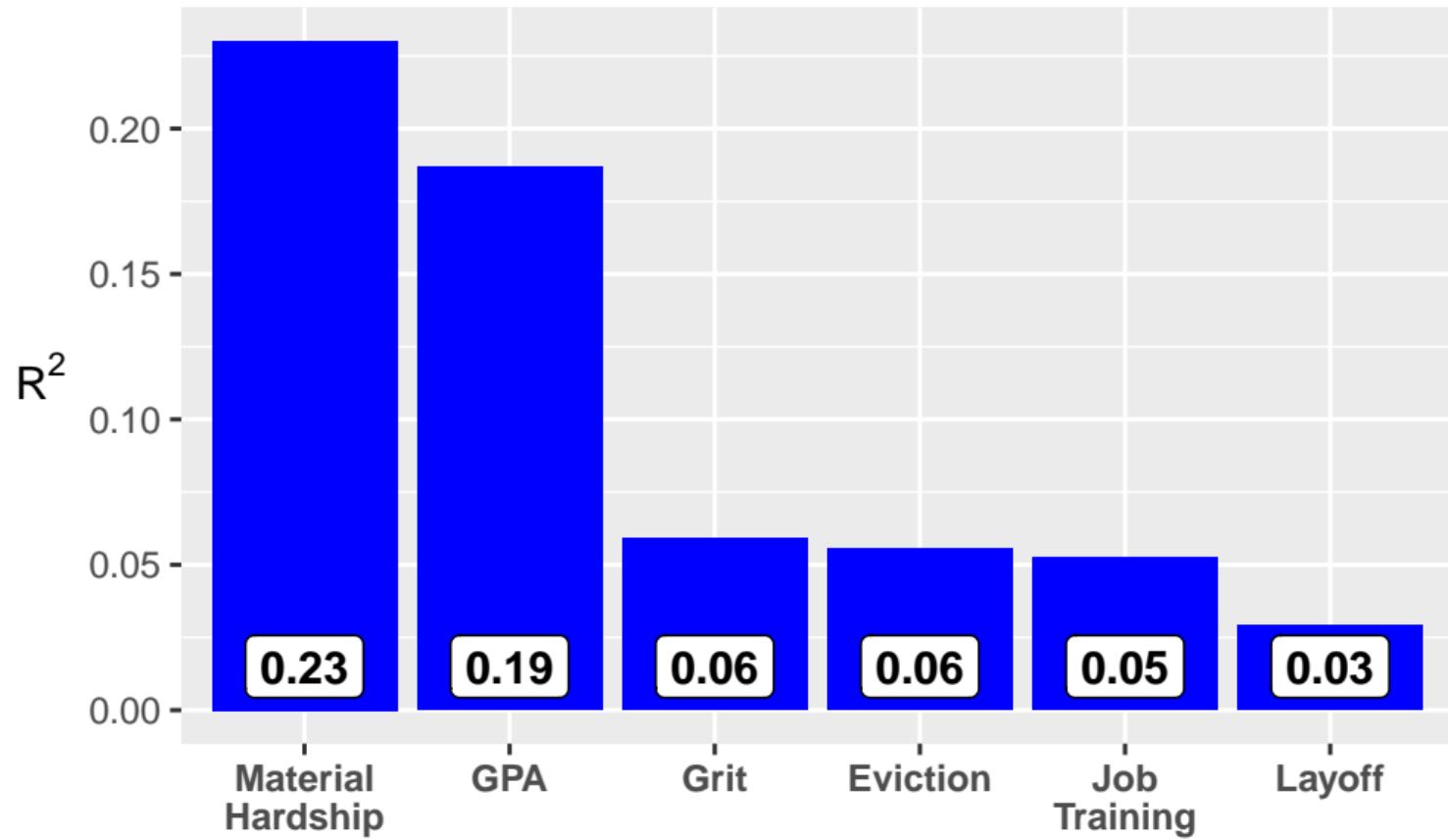
Matthew J. Salganik<sup>a,1</sup>, Ian Lundberg<sup>a</sup> , Alexander T. Kindel<sup>a</sup>, Caitlin E. Ahearn<sup>b</sup>, Khaled Al-Ghoneim<sup>c</sup>, Abdullah Almaatouq<sup>d,e</sup> , Drew M. Altschul<sup>f</sup> , Jennie E. Brand<sup>b,g</sup>, Nicole Bohme Carnegie<sup>h</sup> , Ryan James Compton<sup>i</sup>, Debanjan Datta<sup>j</sup>, Thomas Davidson<sup>k</sup>, Anna Filippova<sup>l</sup>, Connor Gilroy<sup>m</sup>, Brian J. Goode<sup>n</sup>, Eaman Jahani<sup>o</sup>, Ridhi Kashyap<sup>p,q,r</sup> , Antje Kirchner<sup>s</sup>, Stephen McKay<sup>t</sup> , Allison C. Morgan<sup>u</sup> , Alex Pentland<sup>e</sup>, Kivan Polimis<sup>v</sup>, Louis Raes<sup>w</sup> , Daniel E. Rigobon<sup>x</sup>, Claudia V. Roberts<sup>y</sup>, Diana M. Stanescu<sup>z</sup>, Yoshihiko Suhara<sup>e</sup>, Adaner Usmani<sup>aa</sup>, Erik H. Wang<sup>z</sup>, Muna Adem<sup>bb</sup>, Abdulla Alhajri<sup>cc</sup>, Bedoor AlShebli<sup>dd</sup>, Redwane Amin<sup>ee</sup>, Ryan B. Amos<sup>y</sup>, Lisa P. Argyle<sup>ff</sup> , Livia Baer-Bositis<sup>gg</sup>, Moritz Büchi<sup>hh</sup> , Bo-Ryehn Chung<sup>ii</sup>, William Eggert<sup>jj</sup>, Gregory Faletto<sup>kk</sup>, Zhilin Fan<sup>ll</sup>, Jeremy Freese<sup>gg</sup>, Tejomay Gadgil<sup>mm</sup>, Josh Gagné<sup>gg</sup>, Yue Gao<sup>nn</sup>, Andrew Halpern-Manners<sup>bb</sup>, Sonia P. Hashim<sup>y</sup>, Sonia Hausen<sup>gg</sup>, Guanhua He<sup>oo</sup>, Kimberly Higuera<sup>gg</sup>, Bernie Hogan<sup>pp</sup>, Ilana M. Horwitz<sup>qq</sup>, Lisa M. Hummel<sup>gg</sup>, Naman Jain<sup>x</sup>, Kun Jin<sup>rr</sup> , David Jurgens<sup>ss</sup>, Patrick Kaminski<sup>bb,tt</sup>, Areg Karapetyan<sup>uu,vv</sup>, E. H. Kim<sup>gg</sup>, Ben Leizman<sup>y</sup>, Naijia Liu<sup>z</sup>, Malte Möser<sup>y</sup>, Andrew E. Mack<sup>z</sup>, Mayank Mahajan<sup>y</sup>, Noah Mandell<sup>ww</sup>, Helge Marahrens<sup>bb</sup>, Diana Mercado-Garcia<sup>qq</sup>, Viola Mocz<sup>xx</sup>, Katarina Mueller-Gastell<sup>gg</sup>, Ahmed Musse<sup>yy</sup>, Qiankun Niu<sup>ee</sup>, William Nowak<sup>zz</sup>, Hamidreza Omidvar<sup>aaa</sup>, Andrew Or<sup>y</sup>, Karen Ouyang<sup>y</sup>, Katy M. Pinto<sup>bb,b</sup>, Ethan Porter<sup>cc</sup>, Kristin E. Porter<sup>ddd</sup>, Crystal Qian<sup>y</sup>, Tamkinat Rauf<sup>gg</sup>, Anahit Sargsyan<sup>ee</sup>, Thomas Schaffner<sup>y</sup>, Landon Schnabel<sup>gg</sup>, Bryan Schonfeld<sup>z</sup>, Ben Sender<sup>ff</sup>, Jonathan D. Tang<sup>y</sup>, Emma Tsurkov<sup>gg</sup>, Austin van Loon<sup>gg</sup>, Onur Varol<sup>ggg,hhh</sup> , Xiafei Wang<sup>ii</sup>, Zhi Wang<sup>hhh,jjj</sup>, Julia Wang<sup>y</sup>, Flora Wang<sup>ff</sup>, Samantha Weissman<sup>y</sup>, Kirstie Whitaker<sup>kkk,lll</sup>, Maria K. Wolters<sup>mmm</sup>, Wei Lee Woon<sup>nnn</sup>, James Wu<sup>ooo</sup>, Catherine Wu<sup>y</sup>, Kengran Yang<sup>aaa</sup>, Jingwen Yin<sup>ll</sup>, Bingyu Zhao<sup>ppp</sup>, Chenyun Zhu<sup>ll</sup>, Jeanne Brooks-Gunn<sup>qqq,rrr</sup>, Barbara E. Engelhardt<sup>y,ii</sup>, Moritz Hardt, Dean Knox<sup>z</sup>, Karen Levy<sup>ttt</sup>, Arvind Narayanan<sup>y</sup>, Brandon M. Stewart<sup>a</sup>, Duncan J. Watts<sup>uuu,vvv,wwww</sup> , and Sara McLanahan<sup>a,1</sup>

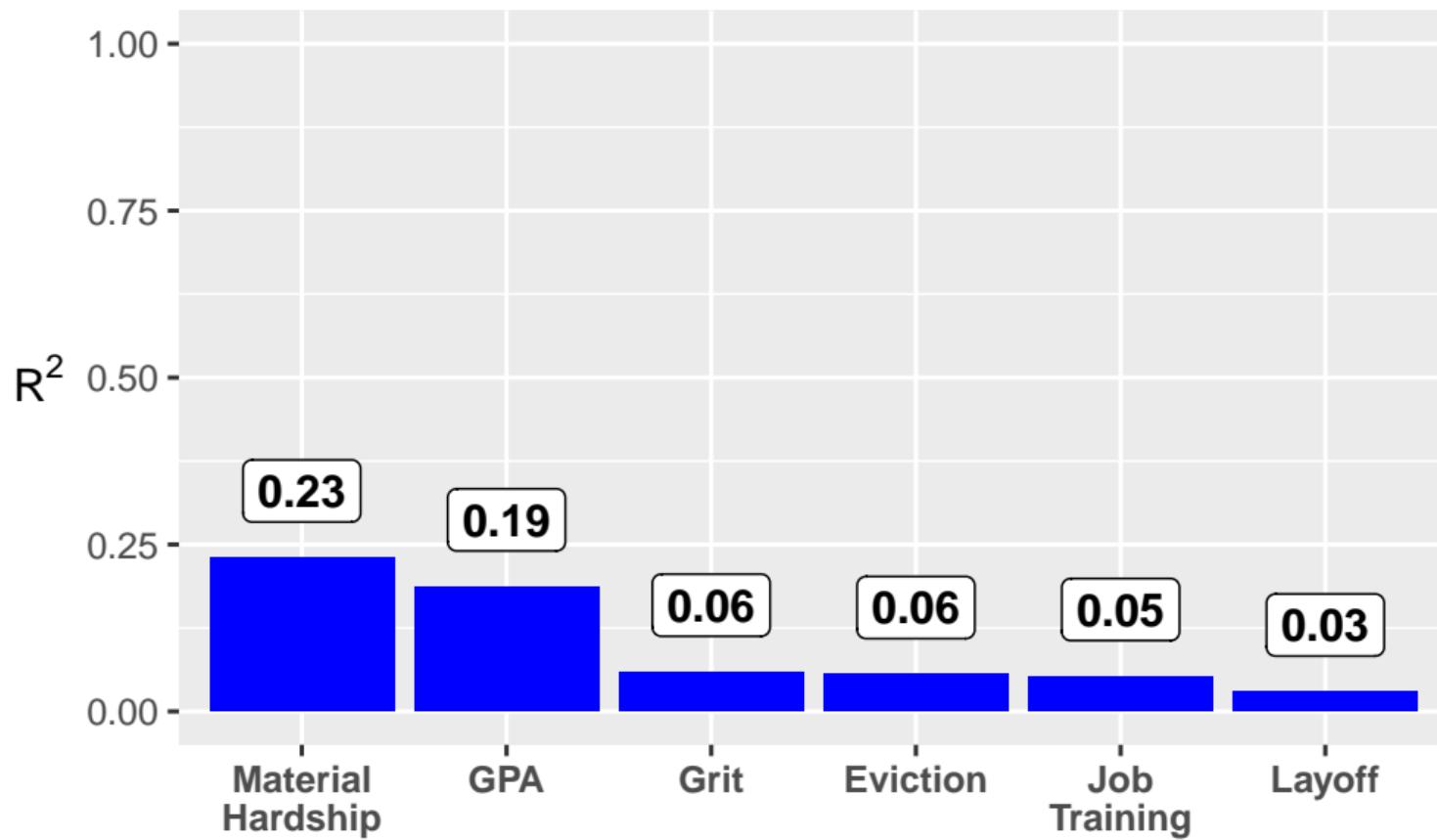
<https://doi.org/10.1073/pnas.1915006117>

457 researchers applied to participate. Many worked in interdisciplinary teams.

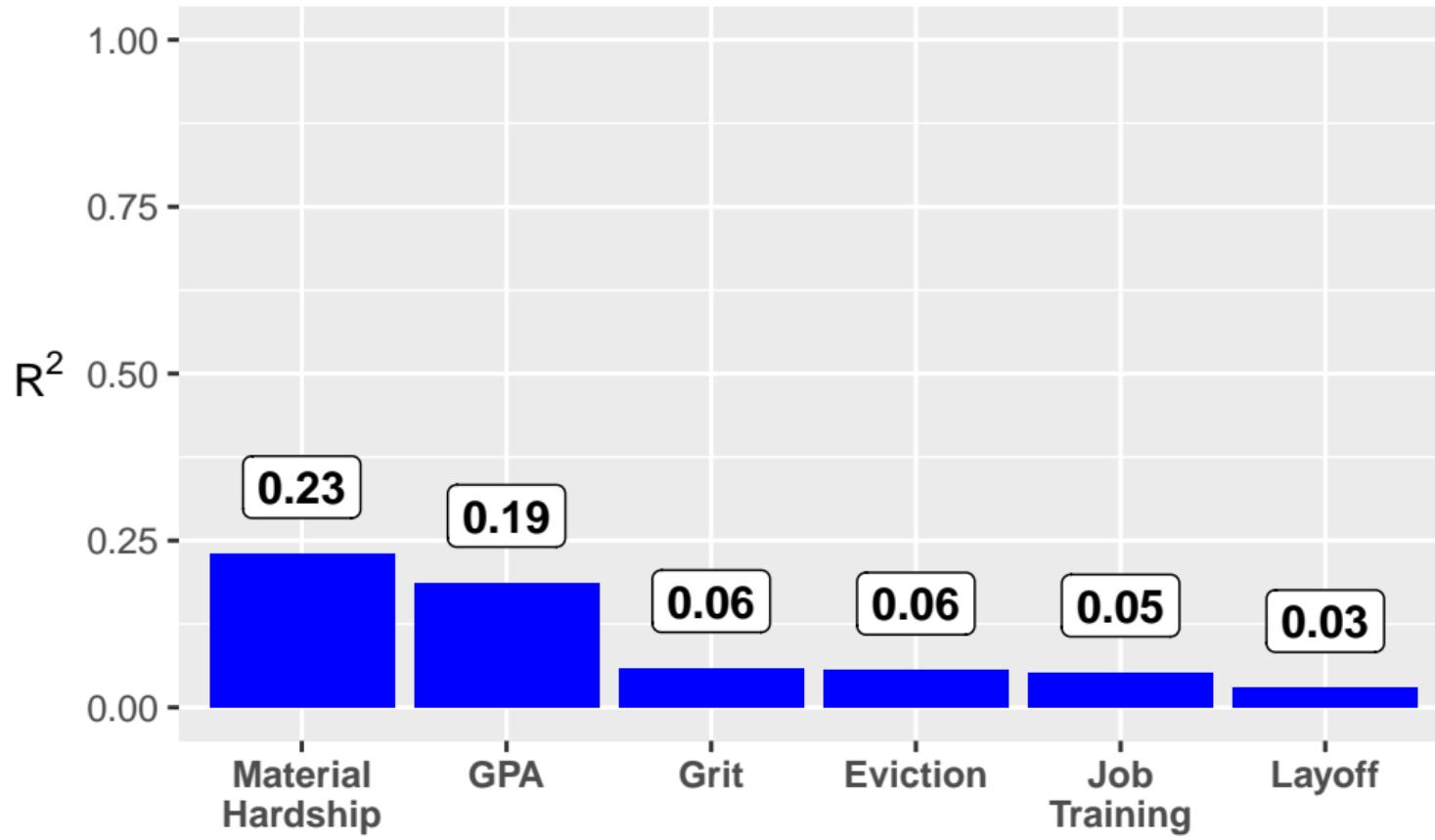
457 researchers applied to participate. Many worked in interdisciplinary teams. Using a large, high-quality social science dataset collected since birth and modern machine learning methods, how accurately can we predict outcomes from children, parents, and families?

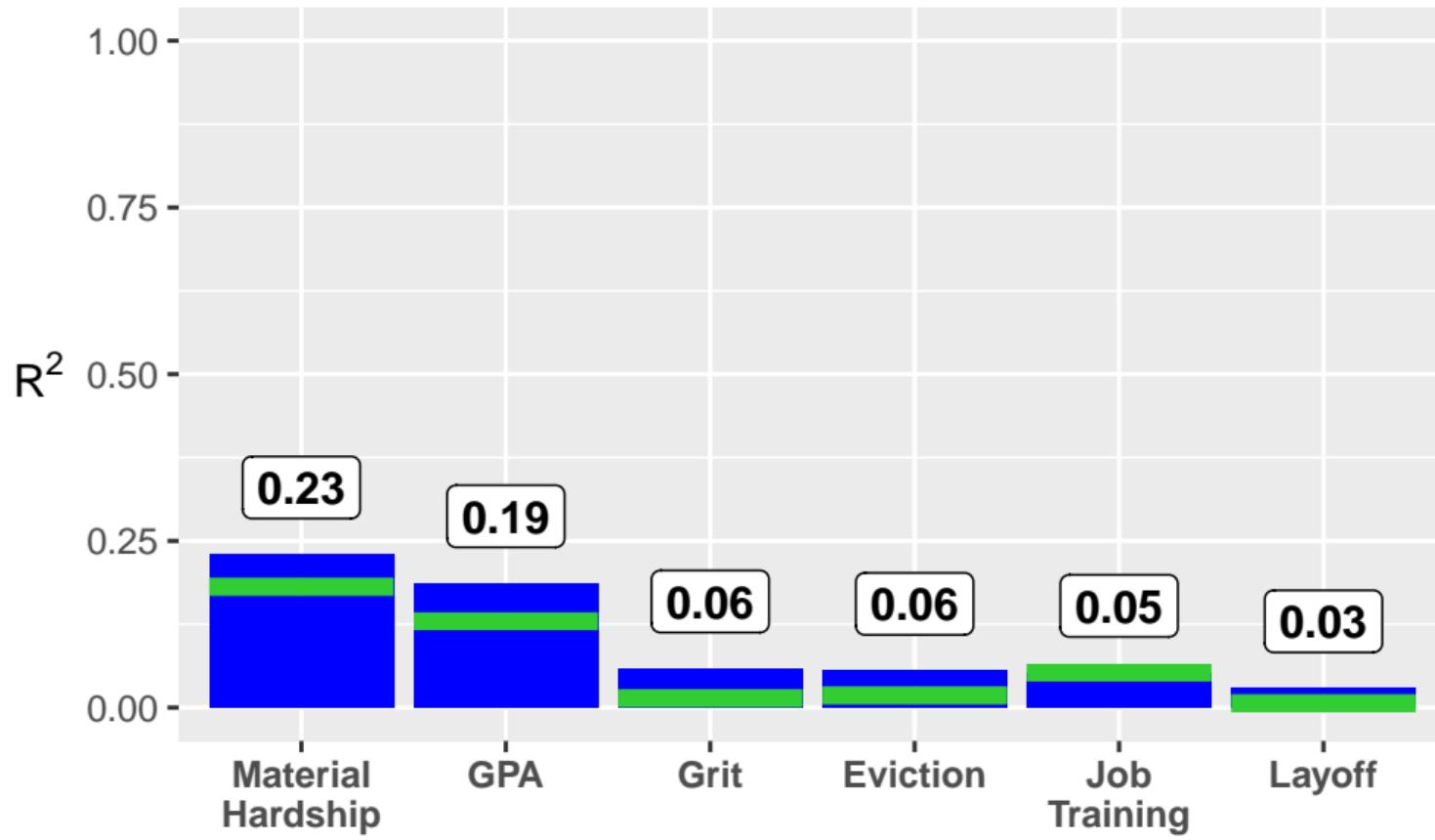
$$R_{holdout}^2 = 1 - \frac{\sum_{i \in holdout} (\hat{y}_i - y_i)^2}{\sum_{i \in holdout} (\bar{y}_{train} - y_i)^2}$$





Is this better than a simple benchmark model?





Green line: 4 variable regression model

**B**

## Material hardship

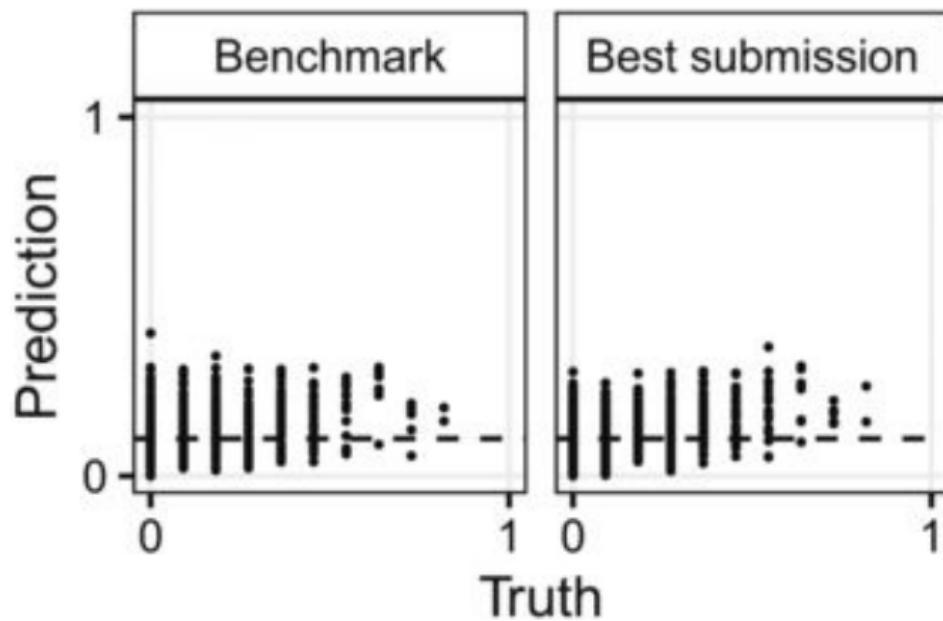


Fig 3, Salganik et al. (2020)

What can we learn looking at all the predictions?

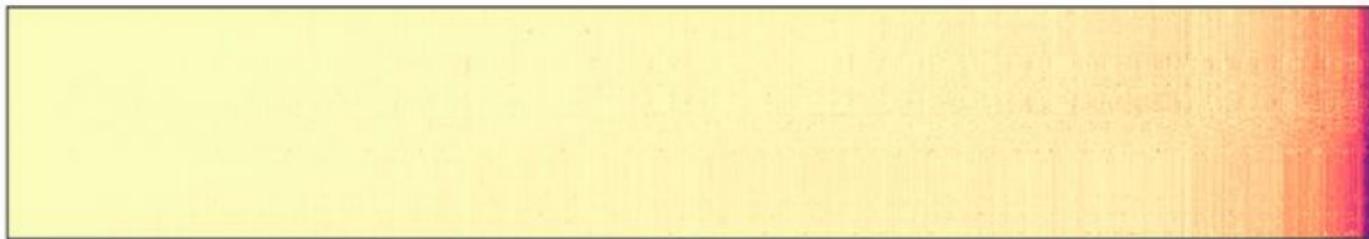
A

## Material hardship

Squared error

0.0 0.1 0.2 0.3 0.4 0.5

Team



Family

Fig 4, Salganik et al. (2020)

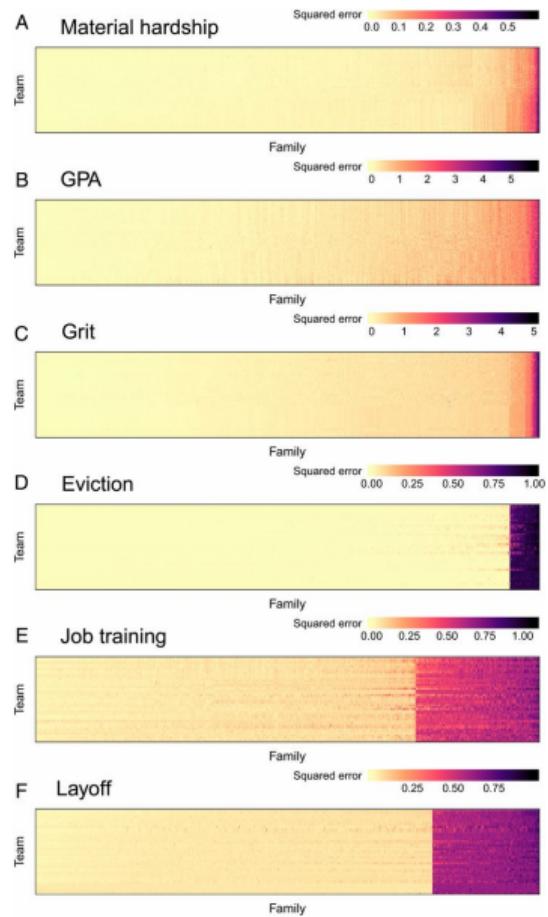


Fig 4, Salganik et al. (2020)

What do these results mean for policy makers?

- ▶ Machine learning is not magic

- ▶ Machine learning is not magic
- ▶ Transparent evaluation of any algorithm is needed

- ▶ Machine learning is not magic
- ▶ Transparent evaluation of any algorithm is needed
- ▶ Complex models may not outperform simple models

What do these results mean for researchers?

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

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

Mullainathan and Spiess (2017)

## Accessing the data

- ▶ You've already registered through Princeton's [OPR](#)
- ▶ Log in and download FFChallenge\_v5

Office of Population Research. [opr.princeton.edu/archive/restricted/Switchboard.aspx](#)

opr.princeton.edu

Paused

# Data Archive Switchboard for Andrew

Welcome to switchboard — your personal starting point in the restricted section of the data archive. Please note that some projects require signing up before you are granted access to data.

You are authorized to access:

Fragile Families and Child Wellbeing Study (FF): [Overview](#) [Download](#) [Update](#)

You may want to sign up for:

Addis Ababa Mortality Surveillance Project (AAMSP):	<a href="#">Overview</a>	<a href="#">SignUp</a>
Success and Failure in Cultural Markets (CM):	<a href="#">Overview</a>	<a href="#">SignUp</a>
The Game of Contacts (GC):	<a href="#">Overview</a>	<a href="#">SignUp</a>
Highly Skilled and Educated Immigrants (HSE):	<a href="#">Overview</a>	<a href="#">SignUp</a>
Immigrant Identity Project (IIP):	<a href="#">Overview</a>	<a href="#">SignUp</a>
Latin American Migration Project (LAMP):	<a href="#">Overview</a>	<a href="#">SignUp</a>
Mexican Migration Project (MMP):	<a href="#">Overview</a>	<a href="#">SignUp</a>
New Immigrant Survey (NIS):	<a href="#">Overview</a>	<a href="#">SignUp</a>
Survey of Unemployed Workers In New Jersey (NJUI):	<a href="#">Overview</a>	<a href="#">SignUp</a>
National Longitudinal Survey of Freshmen (NLSF):	<a href="#">Overview</a>	<a href="#">SignUp</a>
Network Scale-up Method for Heavy Drug Users in Curitiba, Brazil (NSUM):	<a href="#">Overview</a>	<a href="#">SignUp</a>
Project 90 (Partial Data) (P90):	<a href="#">Overview</a>	<a href="#">SignUp</a>
Texas Higher Education Opportunity Project (THEOP):	<a href="#">Overview</a>	<a href="#">SignUp</a>
Wiki Surveys: Open and Quantifiable Social Data Collection (WS):	<a href="#">Overview</a>	<a href="#">SignUp</a>
Wiki Surveys: Open and Quantifiable Social Data Collection (Restricted Sign-up Data) (WSR):	<a href="#">Overview</a>	<a href="#">SignUp</a>

Office of Population Research  opr.princeton.edu/archive/restricted/link.aspx?studyID=37

## Fragile Families and Child Wellbeing Study (FF)

Data for this study are organized by sub-studies as follows:

- All Waves (all FFCWS public data from Baseline - Year 15)
- Baseline
- Year 1
- Year 3
- Year 5
- Year 9
- Year 14
- **Fragile Families Challenge**

Visit project web site, <https://fragilefamilies.princeton.edu/>, for more information.

© 2021 The Trustees of Princeton University  
Office of Population Research | Princeton University, Wallace Hall, Princeton NJ 08544  
Phone: (609) 258-4870 | Fax: (609) 258-1039

PRINCETON UNIVERSITY 

Office of Population Research. Paused

[opr.princeton.edu/archive/restricted/Link.aspx?StudyID=37&SubStudy=substudy80\\_ftchallenge](http://opr.princeton.edu/archive/restricted/Link.aspx?StudyID=37&SubStudy=substudy80_ftchallenge)

 **OPR**  
OFFICE OF POPULATION RESEARCH

## Fragile Families and Child Wellbeing Study (FF)

### Fragile Families Challenge

Following 7 files are available for downloading:

File Name	Size (Zipped)	Last Updated	File Type	Note
ftchallenge_papers_replication_materials	71.8m (71.8m)	02/26/2020	Zip Archive	replication files for two paper published about the Challenge: Salganik et al 2020, PNAS and Salganik et al 2020, Socius
FFChallenge_v2	287.5m (287.5m)	02/26/2020	Zip Archive	files used by participants during the main part of the Challenge
FFChallenge_v5	271.0m (271.0m)	02/26/2020	Zip Archive	improved files first used for the Challenge during the Summer Institute in Computational Social Science (SICSS) in 2018
leaderboard	97.7k (13.7k)	02/26/2020	Zip Archive	outcomes for observations in the leaderboard set, with imputed values for missing outcomes
leaderboardUnfilled	97.7k (13.7k)	02/26/2020	Zip Archive	outcomes for observations in the leaderboard set, without imputed values for missing outcomes
README	5.9k (5.9k)	02/26/2020	Zip Archive	file descriptions
test	50.8k (8.8k)	02/26/2020	Zip Archive	outcome variables for test (holdout) observations

Visit project web site, <https://fragilefamilies.princeton.edu/>, for more information.

© 2021 The Trustees of Princeton University  
Office of Population Research | Princeton University, Wallace Hall, Princeton NJ 08544  
Phone: 609 258-4870 | Fax: 609 258-1059

 PRINCETON  
UNIVERSITY

# Your project

You'll be working with two files in FFChallenge\_v5/:

- ▶ `background.csv` — contains 13,026 variables for 4,242 families.
- ▶ `train.csv` — contains 6 later-life outcome variables for half (2,121) of the families.

# Your project

You'll be working with two files in FFChallenge\_v5/:

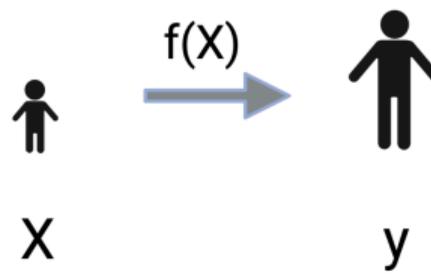
- ▶ `background.csv` — contains 13,026 variables for 4,242 families.  
**This is your  $X$ .**
- ▶ `train.csv` — contains 6 later-life outcome variables for half (2,121) of the families. **This is your  $y$ .**

# Your project

You'll be working with two files in FFChallenge\_v5/:

- ▶ `background.csv` — contains 13,026 variables for 4,242 families.  
**This is your  $X$ .**
- ▶ `train.csv` — contains 6 later-life outcome variables for half (2,121) of the families. **This is your  $y$ .**

**Your goal:** Using the training data, build a predictive model to predict  $y$  from  $X$ .



# Leaderboard

The leaderboard will track the performance of each person's/group's model on the test dataset.

The screenshot shows a web browser window for the "Fragile Families Challenge" on codalab.fragilefamilieschallenge.org. The title bar indicates the page is for the COS 424 Spring 2021 competition. The main content area displays the challenge details and a results table.

**Challenge Details:**

- Name:** COS 424 Spring 2021
- Organized by:** FFData
- Current server time:** Feb. 18, 2021, 5:51 p.m. UTC
- Phase:** Current
- End:** Dec. 21, 2016, midnight UTC
- Challenge Type:** Classification
- Status:** Never

**Results Tab:** This tab is selected, showing the following table:

#	User	GPI ▲	Grit ▲	Material hardship ▲	Eviction ▲	Layoff ▲	Job training ▲
1	baseline	0.39273 (1)	0.21997 (1)	0.02880 (1)	0.05341 (1)	0.17435 (1)	0.20224 (1)
2	ajones788	0.39273 (1)	0.21997 (1)	0.02880 (1)	0.05660 (2)	0.22453 (2)	0.27736 (2)

**Powered by Codalab v0.1.1**

# Imputing missing data

Use MissingDataScript.py to fill in missing values in background.csv.

```
1 import pandas as pd
2 import numpy as np
3
4 def fillMissing(inputcsv, outputcsv):
5
6     # read input csv - takes time
7     df = pd.read_csv(inputcsv, low_memory=False)
8     # Fix date bug
9     df.cf4fint = ((pd.to_datetime(df.cf4fint) - pd.to_datetime('1960-01-01')) / np.timedelta64(1, 'D')).astype(int)
10
11    # replace NA's with mode
12    df = df.fillna(df.mode().iloc[0])
13    # if still NA, replace with 1
14    df = df.fillna(value=1)
15    # replace negative values with 1
16    num = df._get_numeric_data()
17    num[num < 0] = 1
18    # write filled outputcsv
19    df.to_csv(outputcsv, index=False)
20
21    # Usage:
22    if __name__ == "__main__":
23        from os.path import join as pjoin
24
25        data_dir = "./FFChallenge_v5"
26
27        ## impute the background data
28        in_path = pjoin(data_dir, 'background.csv')
29        out_path = pjoin(data_dir, 'background_clean.csv')
30        fillMissing(in_path, out_path)
31
```

## Important notes

- ▶ Submitting the leaderboard requires a very specific format. Read the HW2 doc on Canvas carefully.
- ▶ You'll submit predictions for both the *train* and *test* data
- ▶ Each family has a unique challengeID. Use these to match the samples in `background.csv` with the samples in `train.csv`.

## Support-Vector Networks (Cortes, Vapnik 1995)

~ 15 minutes for discussion

# Wrap-up

- ▶ HW2: Fragile Families Challenge out today!

Good luck!