

# Studying Social Inequality with Data Science

INFO 3370 / 5371  
Spring 2024

## **Predicting life outcomes**

Results of the PSID Income Prediction Challenge

# Learning goals for today

By the end of class, you will be able to

- ▶ know who had the best predictions!
- ▶ reason about predictability of life outcomes

# Equality Opportunity and Prediction

## **Possible claim**

To the degree that we can predict life outcomes,  
people do not have equal opportunity

# Equality Opportunity and Prediction

Learning Set

	Respondent Income	Respondent Education	Parent Education	Grandparent Education	Sex	Race	Grandparent Income	Parent Income
Case 1								
Case 2								
Case 3								
Case 4								
Case 5								

Learn a  
prediction  
function  
→


Holdout Set

Case 6								
Case 7								
Case 8								

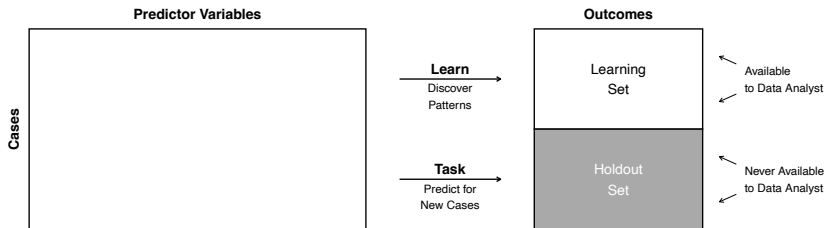
Predict for  
new cases  
→

?
?
?

# The model selection problem

In supervised machine learning, the goal is to

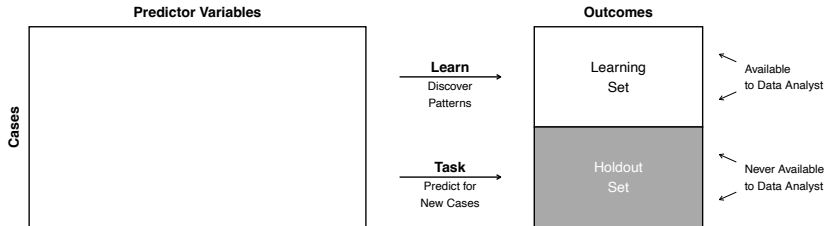
- ▶ learn patterns in the available data
- ▶ predict outcomes for previously unseen cases



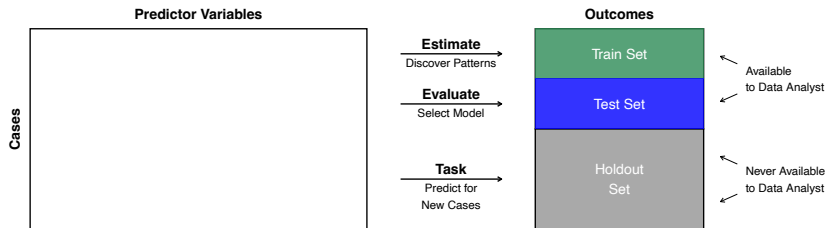
# The model selection problem

When a task involves unseen data,  
mimic the task with data we have

# The model selection problem



# The model selection problem





## Prepare environment

```
library(tidyverse)
library(rsample)
set.seed(14850)
```

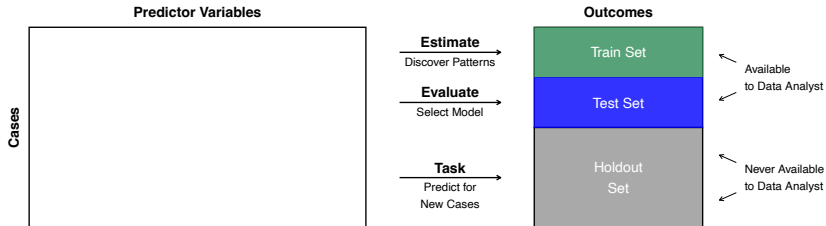
## Load data

```
learning <- read_csv("learning.csv")  
holdout_public <- read_csv("holdout_public.csv")
```

## Create a train-test split within learning

Using the rsample package,

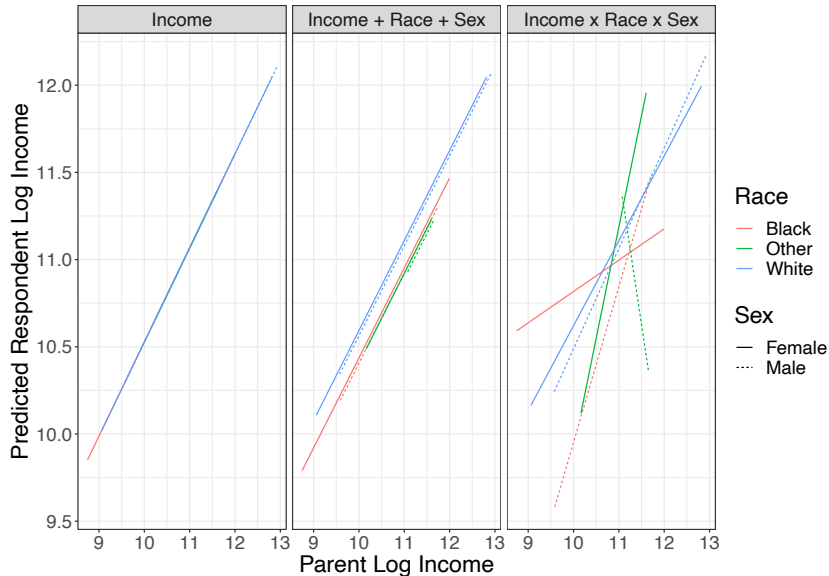
```
split <- learning |>  
  initial_split(prop = 0.5)
```



## Learn candidates in the train set

```
candidate_1 <- lm(  
  g3_log_income ~ g2_log_income,  
  data = training(split)  
)  
candidate_2 <- lm(  
  g3_log_income ~ g2_log_income + race + sex,  
  data = training(split)  
)  
candidate_3 <- lm(  
  g3_log_income ~ g2_log_income * race * sex,  
  data = training(split)  
)
```

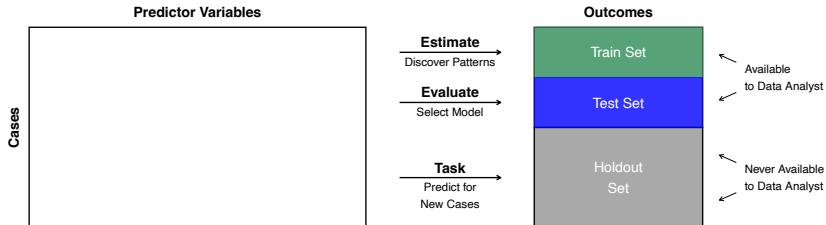
# Learn candidates in the train set



## Evaluate performance on the test set. Choose a model

```
fitted |>
  group_by(model) |>
  mutate(error = g3_log_income - yhat) |>
  mutate(squared_error = error ^ 2) |>
  summarize(mse = mean(squared_error))
```

```
## # A tibble: 3 x 2
##   model      mse
##   <chr>    <dbl>
## 1 candidate_1 0.439
## 2 candidate_2 0.437
## 3 candidate_3 0.477
```





## Apply your chosen model

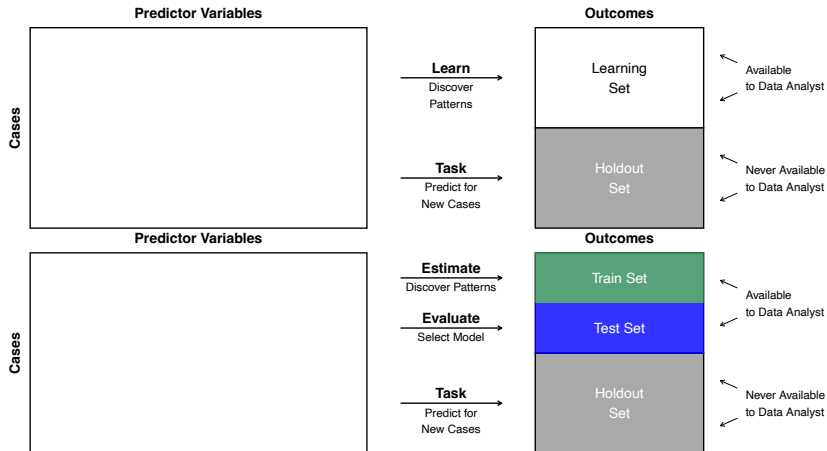
Learn in the full learning set

```
chosen <- lm(  
  g3_log_income ~ g2_log_income +  
    race + sex,  
  data = learning  
)
```

Predict for the holdout set

```
predicted <- holdout_public %>%  
  mutate(  
    predicted = predict(  
      chosen,  
      newdata = holdout_public  
    )  
  )
```

# Summary



## Your submissions








- ▶ 21 submissions
- ▶ 20 submissions predicting for all holdout cases
- ▶ 17 submissions with non-missing predictions
- ▶ 14 submissions by unique teams

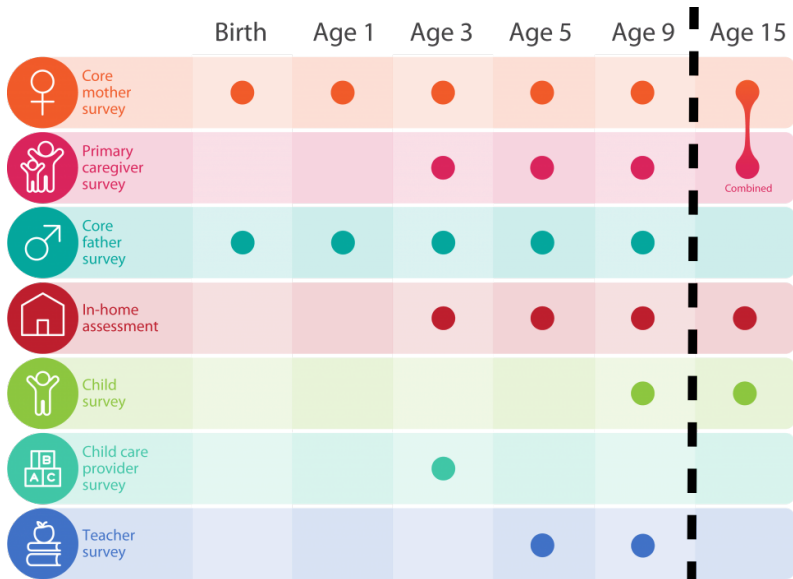
**[class submission results redacted for online posting]**

our exercise was a particular case  
of a broader research project

# 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>aa</sup>, Adaner Usmani<sup>ab</sup>, Erik H. Wang<sup>c</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>c</sup>, Malte Möser<sup>y</sup>, Andrew E. Mack<sup>c</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>bbb</sup>, Ethan Porter<sup>ccc</sup>, Kristin E. Porter<sup>ddd</sup>, Crystal Qian<sup>y</sup>, Tamkinat Rauf<sup>gg</sup>, Anahit Sargsyan<sup>eee</sup>, Thomas Schaffner<sup>y</sup>, Landon Schnabel<sup>gg</sup>, Bryan Schonfeld<sup>z</sup>, Ben Sender<sup>fff</sup>, Jonathan D. Tang<sup>y</sup>, Emma Tsurkov<sup>gg</sup>, Austin van Loon<sup>gg</sup>, Onur Varo<sup>ggg,hhh</sup>, Xiafei Wang<sup>iii</sup>, Zhi Wang<sup>hhh,jjj</sup>, Julia Wang<sup>y</sup>, Flora Wang<sup>fff</sup>, Samantha Weissman<sup>y</sup>, Kirstie Whitaker<sup>kkk,lll</sup>, Maria K. Wolters<sup>mmmm</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<sup>sss</sup>, 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>

	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				●	●



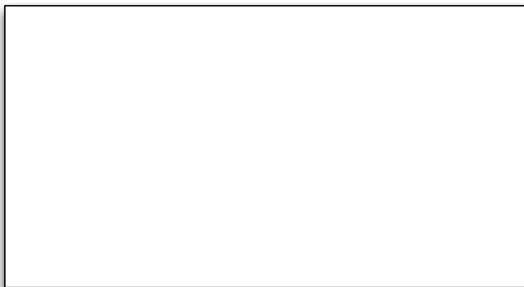


Six age 15 outcomes:

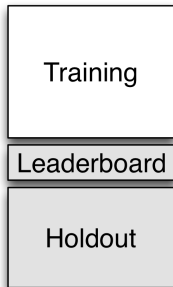
- ▶ GPA
- ▶ Material Hardship
- ▶ Grit
- ▶ Evicted
- ▶ Job training
- ▶ Job loss

4,200 families

12,000 features  
birth to age 9



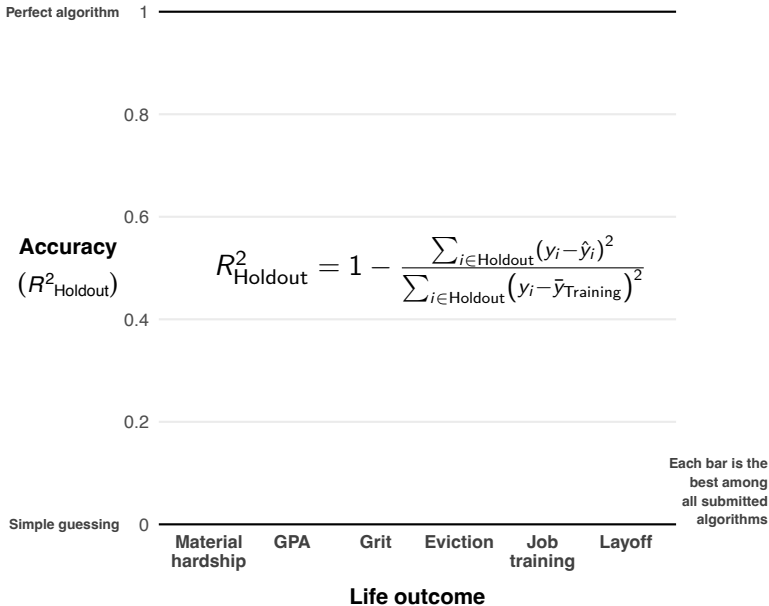
6 outcomes  
age 15



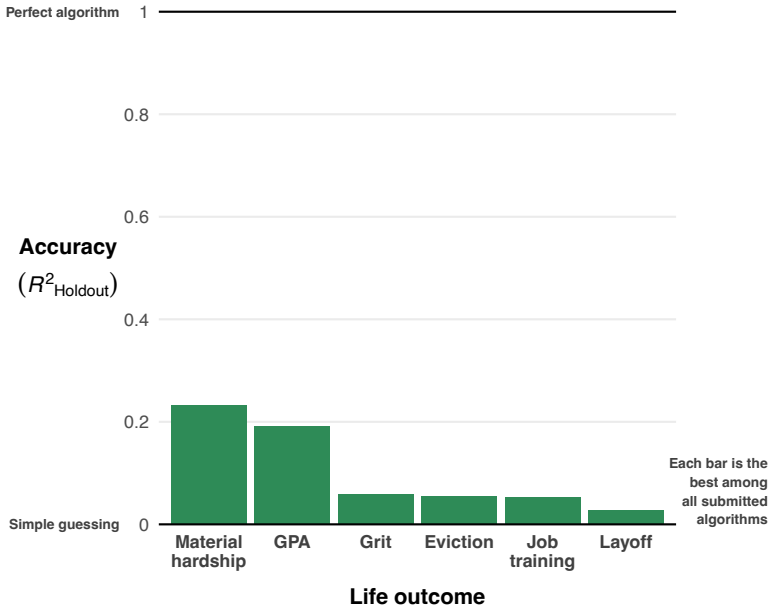
441 registered participants

- ▶ social scientists and data scientists
- ▶ undergraduates, grad students, and professionals
- ▶ many working in teams

How did they do?

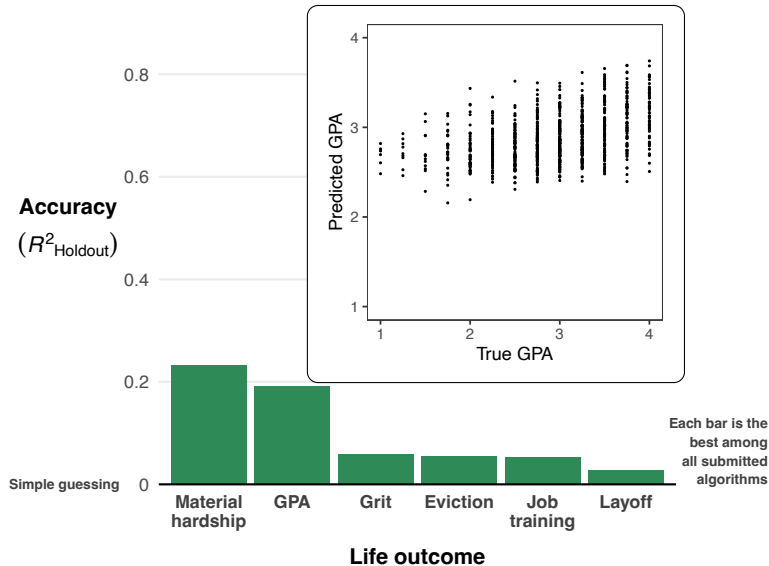


## Best algorithms were not very accurate



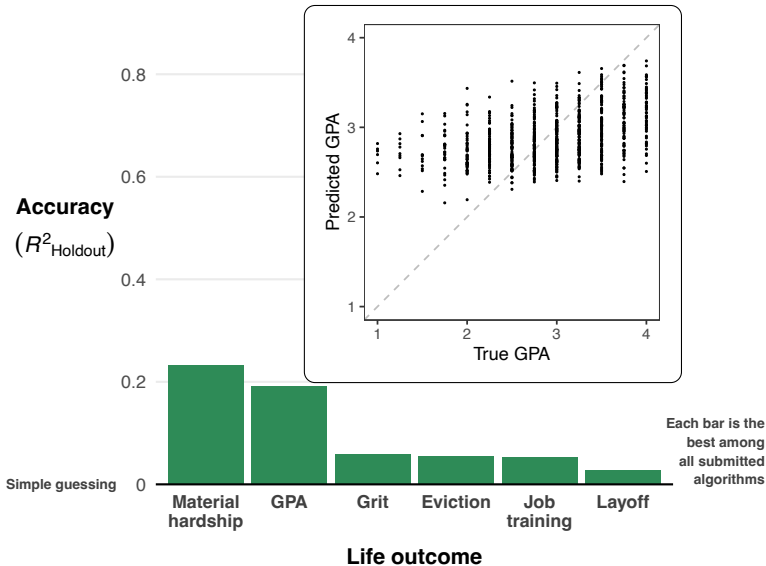
# Best algorithms were not very accurate

Perfect algorithm 1



## Best algorithms were not very accurate

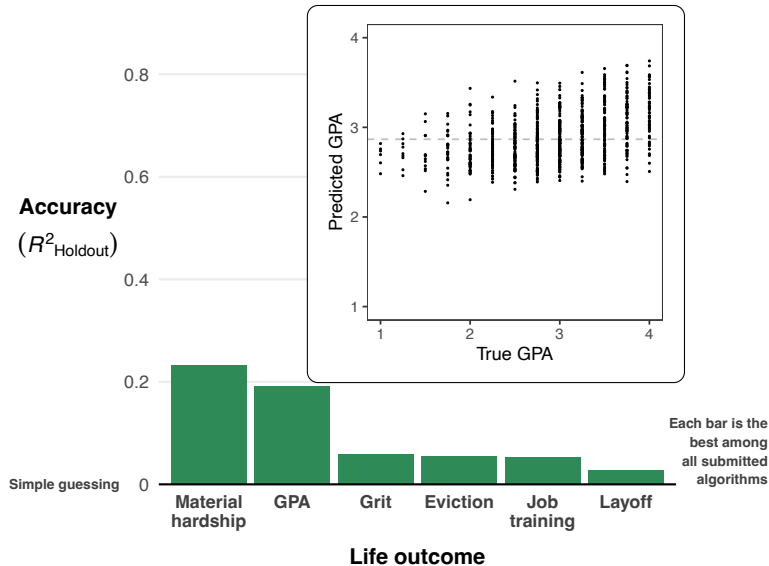
Perfect algorithm 1



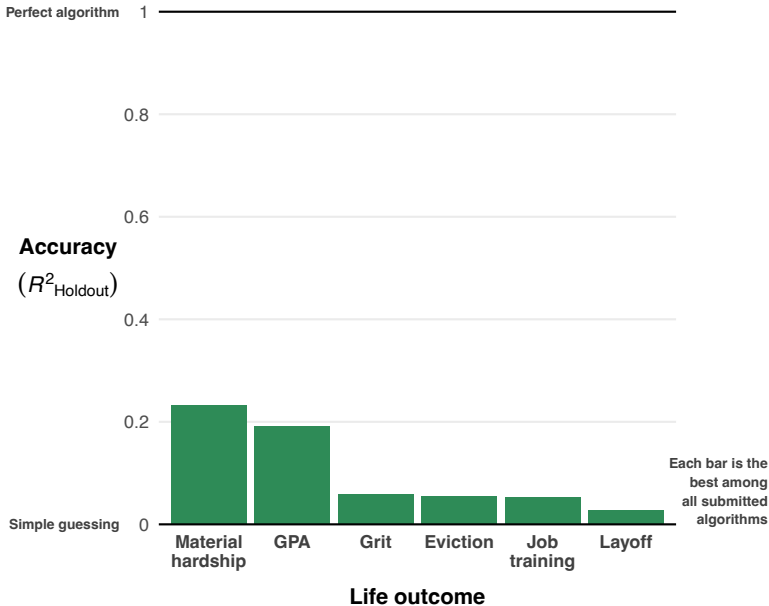


## Best algorithms were not very accurate

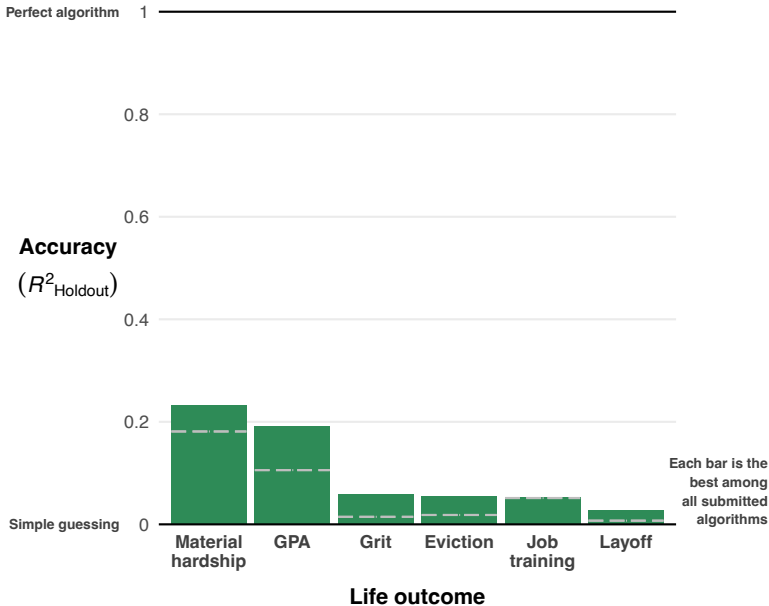
Perfect algorithm 1



## Best algorithms were not very accurate

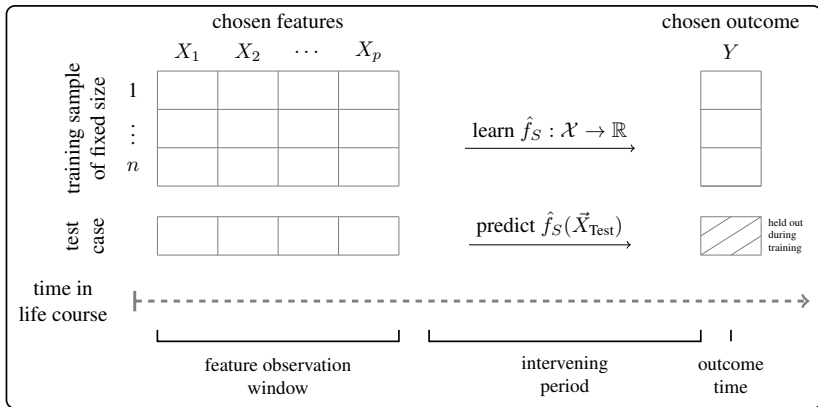


## Best algorithms were not very accurate



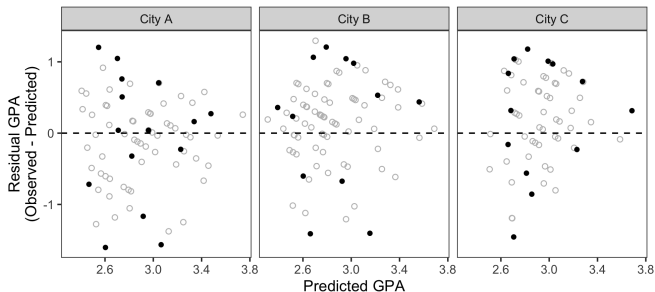
Lundberg et al. 2024.

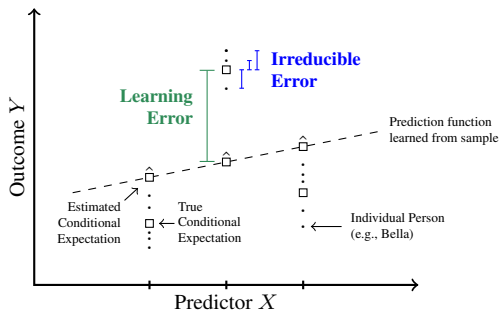
The origins of unpredictability in life outcome prediction tasks



## In-depth, qualitative interviews

- ▶ 73 respondents in 40 families
- ▶ Separate interviews with the youth and primary caregiver
- ▶ Life history of the youth from birth to the interview ( $\approx$  age 18)

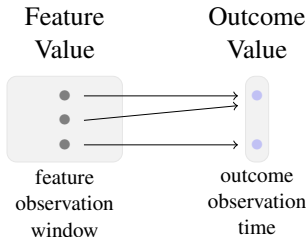




# Irreducible error

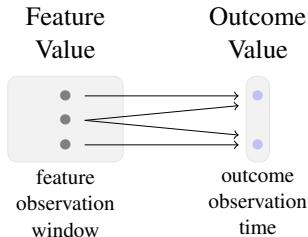
## Zero Irreducible Error

Irreducible error is zero if  
each feature value  
maps to **one** outcome value



## Non-Zero Irreducible Error

Irreducible error is non-zero if  
at least one feature value  
maps to **multiple** outcome values





## Irreducible error: Unmeasurable features

Unmeasurable features occur after the feature observation window

## Irreducible error: Unmeasurable features

Unmeasurable features occur after the feature observation window

- ▶ Bella: A lasting event

# Irreducible error: Unmeasurable features

Unmeasurable features occur after the feature observation window

- ▶ Bella: A lasting event
  - ▶ after age 9, her father died

# Irreducible error: Unmeasurable features

Unmeasurable features occur after the feature observation window

- ▶ Bella: A lasting event
  - ▶ after age 9, her father died
  - ▶ high school went off course

## Irreducible error: Unmeasurable features

Unmeasurable features occur after the feature observation window

- ▶ Bella: A lasting event
  - ▶ after age 9, her father died
  - ▶ high school went off course
- ▶ Charles: A fleeting event

# Irreducible error: Unmeasurable features

Unmeasurable features occur after the feature observation window

- ▶ Bella: A lasting event
  - ▶ after age 9, her father died
  - ▶ high school went off course
- ▶ Charles: A fleeting event
  - ▶ online high school

# Irreducible error: Unmeasurable features

Unmeasurable features occur after the feature observation window

- ▶ Bella: A lasting event
  - ▶ after age 9, her father died
  - ▶ high school went off course
- ▶ Charles: A fleeting event
  - ▶ online high school
  - ▶ worked in the basement for one semester

# Irreducible error: Unmeasurable features

Unmeasurable features occur after the feature observation window

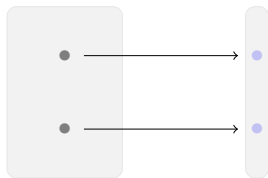
- ▶ Bella: A lasting event
  - ▶ after age 9, her father died
  - ▶ high school went off course
- ▶ Charles: A fleeting event
  - ▶ online high school
  - ▶ worked in the basement for one semester
  - ▶ video games = bad grades that semester



# Irreducible error: Unmeasurable features

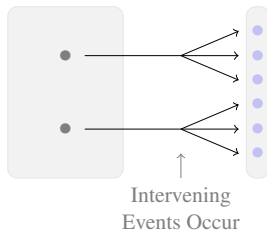
## Zero Irreducible Error

Without intervening events,



## Non-Zero Irreducible Error

With intervening events,



Irreducible error: Unmeasured features

## Irreducible error: Unmeasured features

Lola's social network

## Irreducible error: Unmeasured features

Lola's social network

- ▶ elderly neighbor got Lola ready for school each day

## Irreducible error: Unmeasured features

Lola's social network

- ▶ elderly neighbor got Lola ready for school each day
- ▶ grandparents remodeled the basement to house Lola

## Irreducible error: Unmeasured features

Lola's social network

- ▶ elderly neighbor got Lola ready for school each day
- ▶ grandparents remodeled the basement to house Lola
- ▶ aunt employed Lola's mother in a family business

## Irreducible error: Unmeasured features

Lola's social network

- ▶ elderly neighbor got Lola ready for school each day
- ▶ grandparents remodeled the basement to house Lola
- ▶ aunt employed Lola's mother in a family business

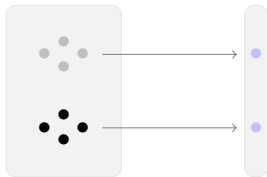
Predicted GPA: 3.04

Actual GPA: 3.75

# Irreducible error: Unmeasured features

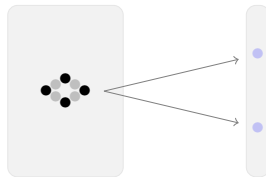
## Zero Irreducible Error

Feature is measured,



## Non-Zero Irreducible Error

Feature is unmeasured,





Irreducible error: Imperfectly measured features

## Irreducible error: Imperfectly measured features

How close do you feel to your mom? Would you say...

Extremely close, .....	1
Quite close,.....	2
Fairly close, or, .....	3
Not very close? .....	4
REFUSED .....	-1
DON'T KNOW .....	-2

## Irreducible error: Imperfectly measured features

How close do you feel to your mom? Would you say...

Extremely close, .....	1
Quite close,.....	2
Fairly close, or, .....	3
Not very close? .....	4
REFUSED .....	-1
DON'T KNOW .....	-2

A daughter told us about her “not very close” mother

# Irreducible error: Imperfectly measured features

How close do you feel to your mom? Would you say...

Extremely close, .....	1
Quite close,.....	2
Fairly close, or, .....	3
Not very close? .....	4
REFUSED .....	-1
DON'T KNOW .....	-2

A daughter told us about her “not very close” mother

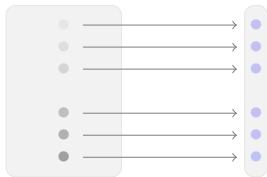
- ▶ kicked her out of the house and called police
- ▶ mother: “you better start treating me better, because I might not live that long.’ ’
- ▶ daughter: “I couldn’t even focus in class. . . I was shaking.’ ’

Outcome: Failed 8th grade. Low GPA. Dropped out.

# Irreducible error: Imperfectly measured features

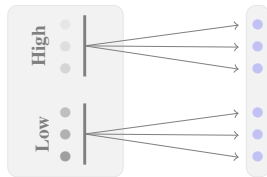
## Zero Irreducible Error

Granular measurement,



## Non-Zero Irreducible Error

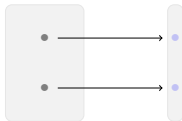
Coarse measurement,



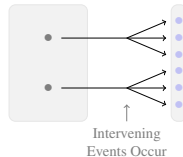
**Unmeasurable features**

Events after the feature observation window create outcome variance

Without intervening events,

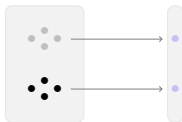


With intervening events,

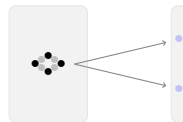
**Unmeasured features**

A measurable feature could distinguish units with highly disparate outcomes

Feature is measured,



Feature is unmeasured,

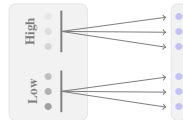
**Imperfectly-measured features**

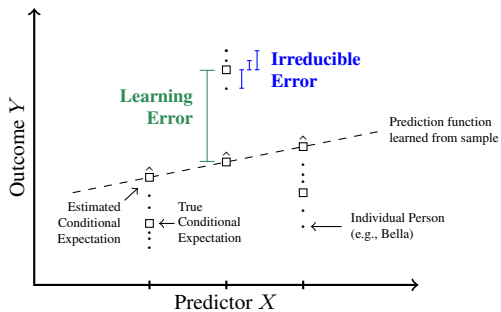
A feature is measured in coarse categories

Granular measurement,



Coarse measurement,

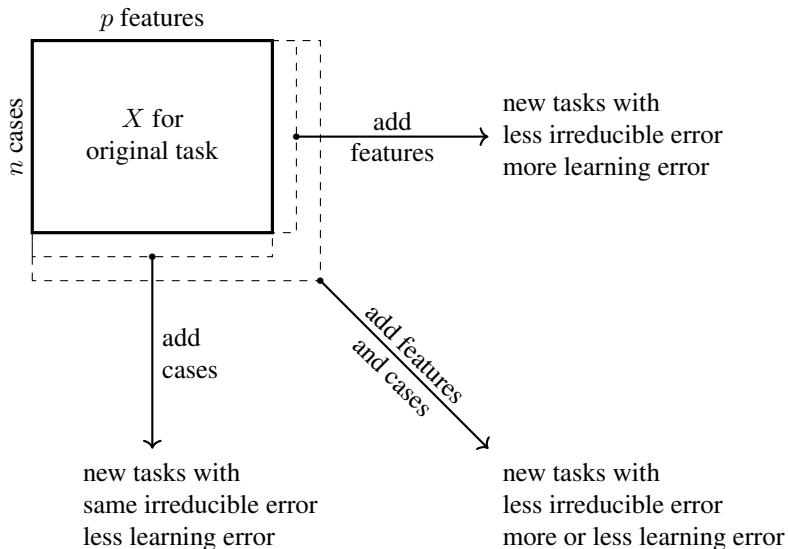




## DISCUSSION



# Generalizing to other life outcome prediction tasks



## Implications for policy

## Implications for policy

- ▶ life outcome predictions may be inaccurate

# Implications for policy

- ▶ life outcome predictions may be inaccurate
  - ▶ if generated by algorithms
  - ▶ if generated by humans

# Implications for policy

- ▶ life outcome predictions may be inaccurate
  - ▶ if generated by algorithms
  - ▶ if generated by humans
- ▶ from accuracy to impact evaluations

## Implications for science

# Implications for science

- ▶ old goal: between-group variability
  - ▶ how means vary across groups

# Implications for science

- ▶ old goal: between-group variability
  - ▶ how means vary across groups
- ▶ new goal: within-group variability
  - ▶ how variances vary across groups



# Implications for science

- ▶ old goal: between-group variability
  - ▶ how means vary across groups
- ▶ new goal: within-group variability
  - ▶ how variances vary across groups
- ▶ more work to better understand unpredictability
  - ▶ empirical estimates
  - ▶ formal models

# Learning goals for today

By the end of class, you will be able to

- ▶ know who had the best predictions!
- ▶ reason about predictability of life outcomes