

Machine learning for population description

Ian Lundberg
UCLA

Including past work with
Rebecca Johnson (Georgetown)
Brandon Stewart (Princeton)

and current work with
Kristin Liao (UCLA)

ilundberg.github.io/description

We acknowledge support through facilities and resources provided by the California Center for Population Research at UCLA (CCPR), which receives core support (P2C-HD041022) from the Eunice Kennedy Shriver National Institute of Child Health and Human Development (NICHD). The content is solely the responsibility of the authors and does not necessarily represent the official views of the Eunice Kennedy Shriver National Institute of Child Health & Human Development or the National Institutes of Health.

Plan for today

- ▶ Estimands in quantitative social science
- ▶ Descriptive estimands: A \hat{Y} view
- ▶ Intro to tomorrow's computer tutorial

What Is Your Estimand? Defining the Target Quantity Connects Statistical Evidence to Theory

Ian Lundberg,^a  Rebecca Johnson,^b  and
Brandon M. Stewart^a 

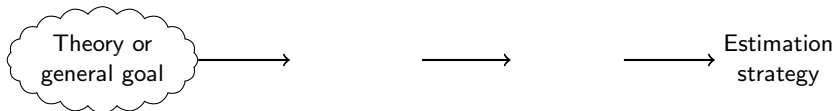
American Sociological Review
1–34

© American Sociological
Association 2021

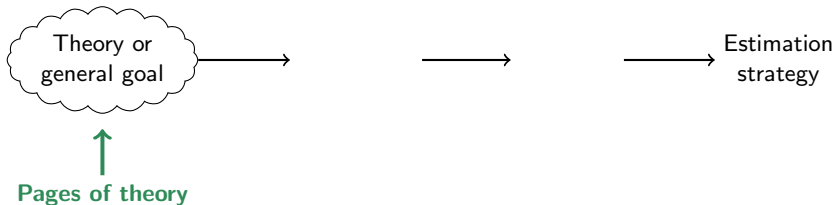
DOI:10.1177/00031224211004187
journals.sagepub.com/home/asr



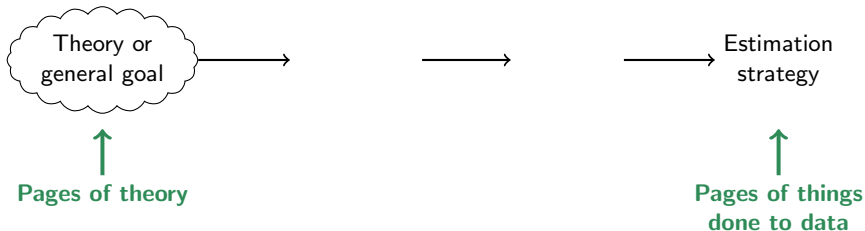
Research framework: Estimands connect theory to evidence



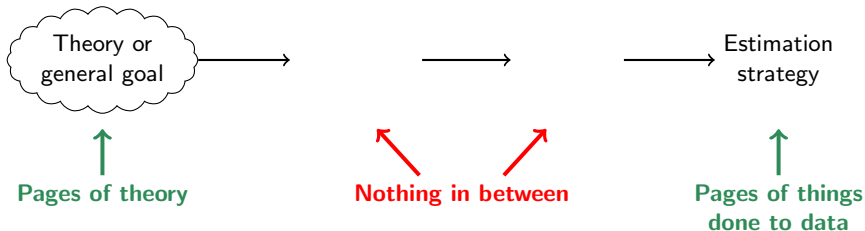
Research framework: Estimands connect theory to evidence



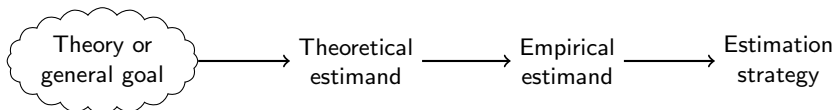
Research framework: Estimands connect theory to evidence



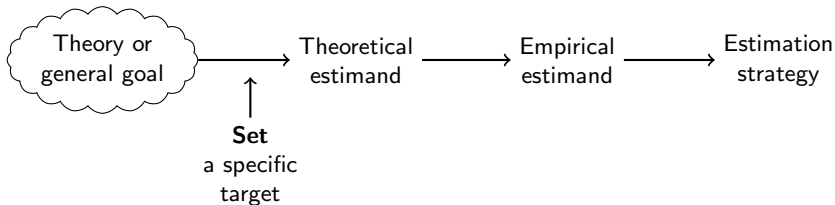
Research framework: Estimands connect theory to evidence



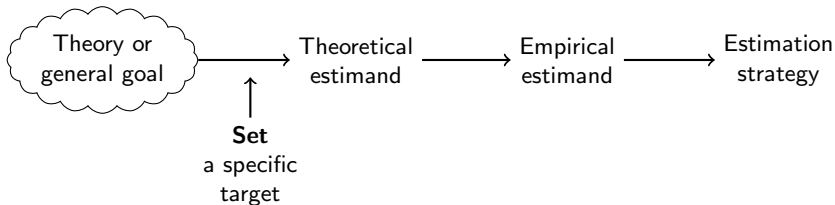
Research framework: Estimands connect theory to evidence



Research framework: Estimands connect theory to evidence



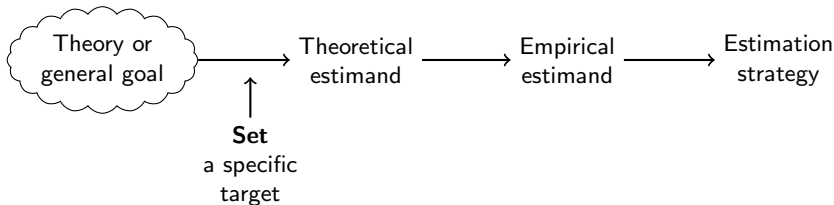
Research framework: Estimands connect theory to evidence



Definition

A **unit-specific quantity**
aggregated over a
target population

Research framework: Estimands connect theory to evidence



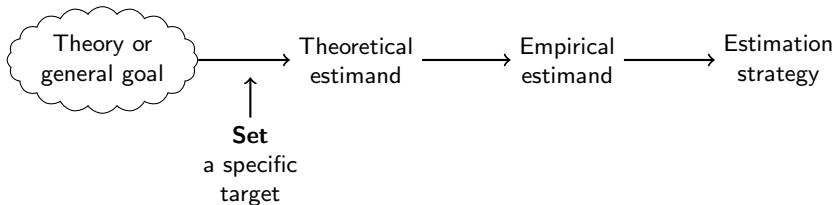
Definition

A **unit-specific quantity**
aggregated over a
target population

Example

$$\frac{1}{\text{Size of U.S. adult population}} \sum_{i \text{ in U.S. adult population}} \left(\text{Employed}_i \right)$$

Research framework: Estimands connect theory to evidence



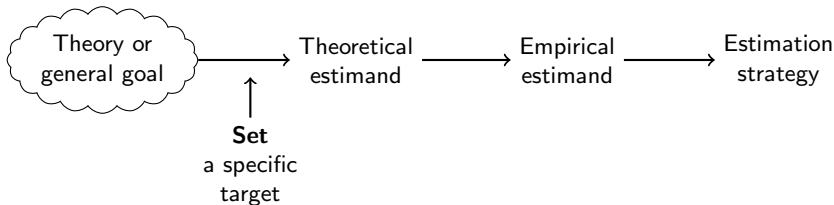
Definition

A **unit-specific quantity**
aggregated over a
target population

Example

$$\frac{1}{\text{Size of U.S. adult population}} \sum_{i \text{ in U.S. adult population}} \left(\underbrace{\text{Employed}_i(\text{Job training})}_{\text{Employment if received job training}} - \underbrace{\text{Employed}_i(\text{No job training})}_{\text{Employment if did not receive job training}} \right)$$

Research framework: Estimands connect theory to evidence



Definition

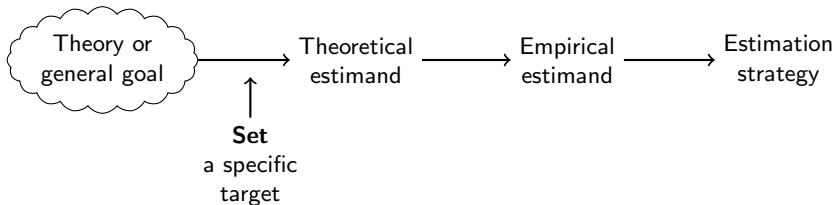
A **unit-specific quantity**
aggregated over a
target population

Example

$$\frac{1}{\text{Size of U.S. adult population}} \sum_{i \text{ in U.S. adult population}} \left(\underbrace{\text{Employed}_i(\text{Job training})}_{\text{Employment if received job training}} - \underbrace{\text{Employed}_i(\text{No job training})}_{\text{Employment if did not receive job training}} \right)$$

Liebersen 1987, Abbott 1988, Freedman 1991, Xie 2013, Hernán 2018

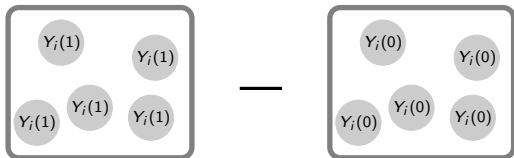
Research framework: Estimands connect theory to evidence



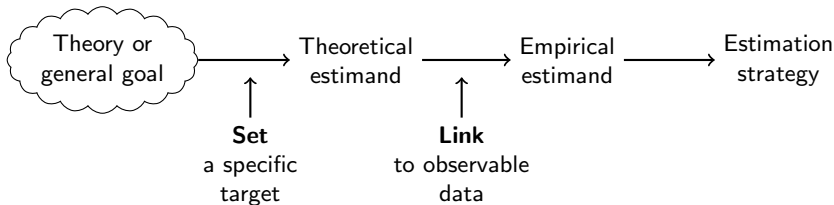
Definition

A **unit-specific quantity**
aggregated over a
target population

Example



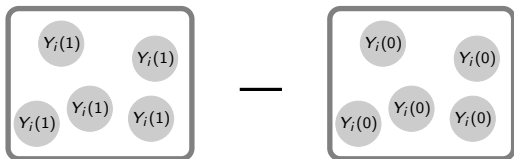
Research framework: Estimands connect theory to evidence



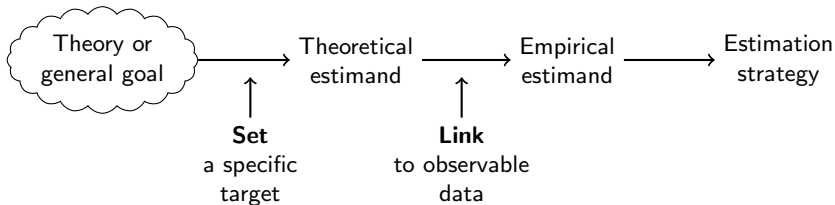
Definition

A quantity involving
observable data

Example



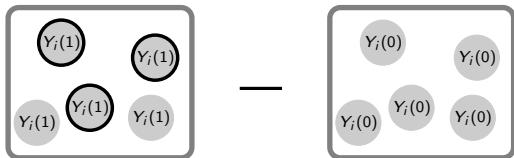
Research framework: Estimands connect theory to evidence



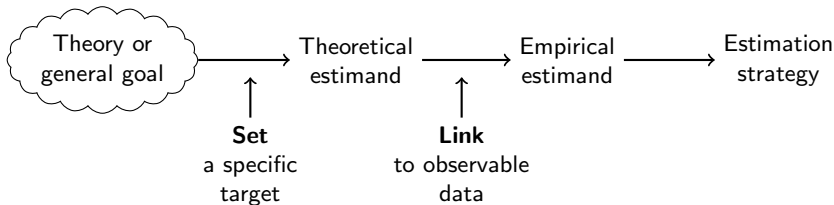
Definition

A quantity involving
observable data

Example



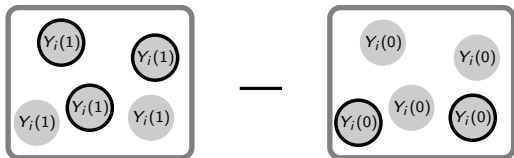
Research framework: Estimands connect theory to evidence



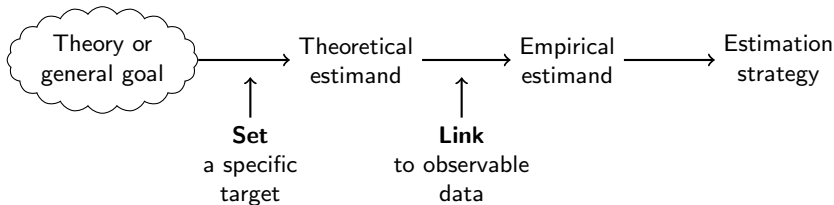
Definition

A quantity involving
observable data

Example



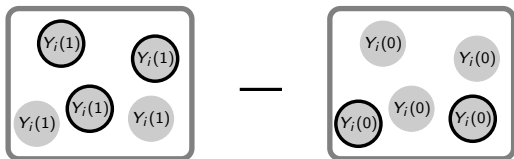
Research framework: Estimands connect theory to evidence



Definition

A quantity involving
observable data

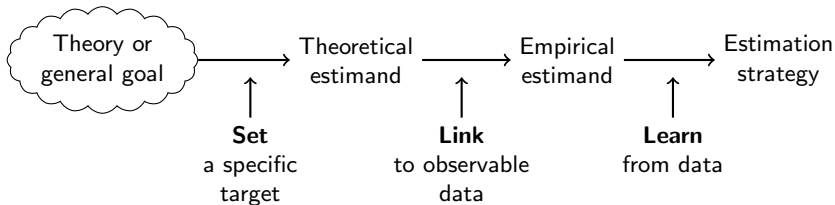
Example



$$\vec{X} \begin{matrix} \xrightarrow{\quad} T \xrightarrow{\quad} Y \\ \quad \quad \quad \searrow \quad \quad \quad \nearrow \end{matrix}$$

Pearl 2009, Imbens and Rubin 2015,
Morgan and Winship 2015, Elwert and Winship 2014

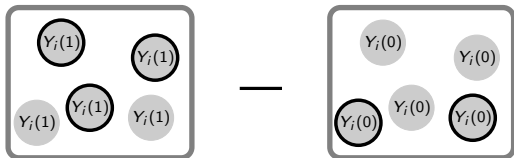
Research framework: Estimands connect theory to evidence



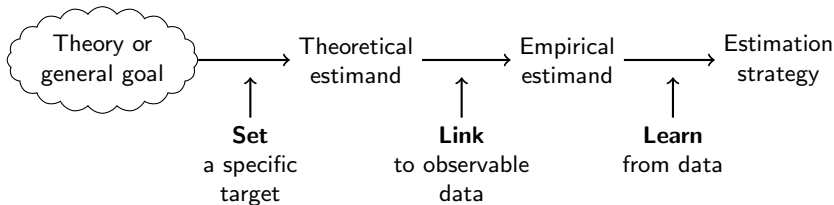
Definition

An algorithm applied to data

Example



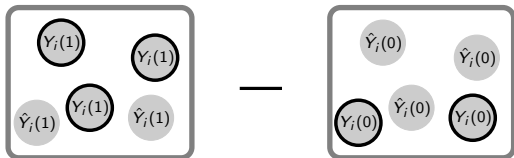
Research framework: Estimands connect theory to evidence



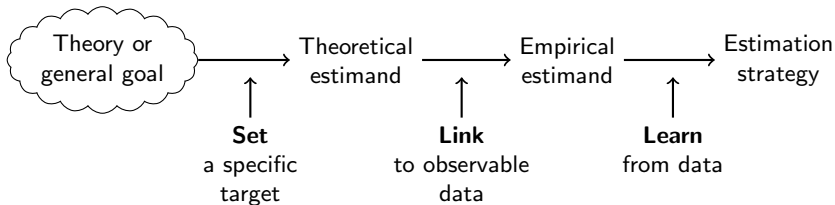
Definition

An algorithm applied to data

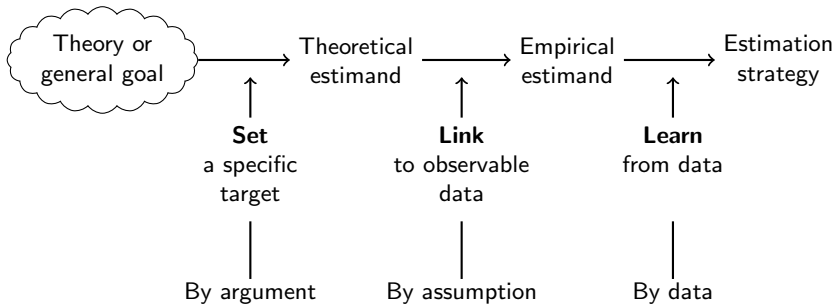
Example

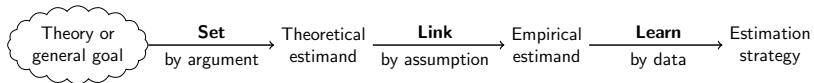


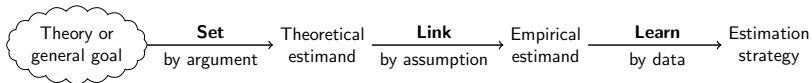
Research framework: Estimands connect theory to evidence



Research framework: Estimands connect theory to evidence







Effect of motherhood
on employment



Effect of motherhood
on employment

First two births
are the same sex

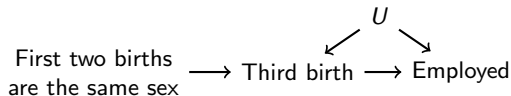


Effect of motherhood
on employment

First two births are the same sex → Third birth



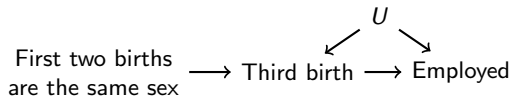
Effect of motherhood
on employment

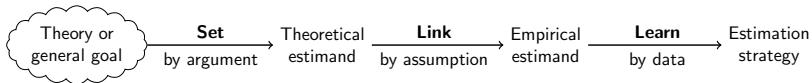




Vague estimand

Effect of motherhood
on employment

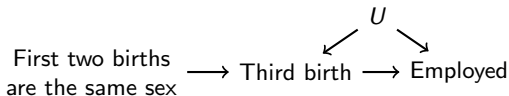




Vague estimand

Effect of motherhood
on employment

Precise estimand





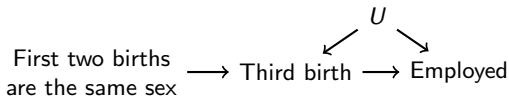
Vague estimand

Effect of motherhood
on employment

Precise estimand

Effect of having **3 vs. 2 children**

**unit-specific
quantity**





Vague estimand

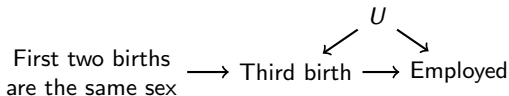
Effect of motherhood on employment

target population



Precise estimand

Effect of having 3 vs. 2 children among those with at least two children who would have a third birth if and only if the first two were of the same sex





Precise estimand

Effect of having 3 vs. 2 children
among those with at least two children who
would have a third birth if and only if the
first two were of the same sex

$\approx 4\%$ of all mothers



Precise estimand

Effect of having 3 vs. 2 children
among those with at least two children who
would have a third birth if and only if the
first two were of the same sex

$\approx 4\%$ of all mothers

You have to argue either:

- 1)
- 2)



Precise estimand

Effect of having 3 vs. 2 children
among those with at least two children who
would have a third birth if and only if the
first two were of the same sex

$\approx 4\%$ of all mothers

You have to argue either:

- 1) That estimand matters for theory, or
- 2)



Precise estimand

Effect of having 3 vs. 2 children
among those with at least two children who
would have a third birth if and only if the
first two were of the same sex

$\approx 4\%$ of all mothers

You have to argue either:

- 1) That estimand matters for theory, or
- 2) It speaks to some broader estimand



1. Set the target quantity.



Describe a population

What is the proportion employed
among U.S. resident women ages 21–35?



Describe a population

What is the proportion employed
among U.S. resident women ages 21–35?

Woman 1

Woman 2

Woman 3

Woman 4



Describe a population

What is the proportion employed
among U.S. resident women ages 21–35?

	<u>Employed?</u>
Woman 1	1
Woman 2	0
Woman 3	1
Woman 4	1



Describe population subgroups

What is the proportion employed among U.S. resident women ages 21–35, comparing mothers to non-mothers?



Describe population subgroups

What is the proportion employed among U.S. resident women ages 21–35, comparing mothers to non-mothers?

	<u>Employed?</u>		<u>Employed?</u>
Mother 1	0	Non-Mother 1	1
Mother 2	0	Non-Mother 2	0
Mother 3	0	Non-Mother 3	1
Mother 4	1	Non-Mother 4	1



Causal effect in a population

What is the causal effect of motherhood on employment among U.S. resident women ages 21–35?



Causal effect in a population

What is the causal effect of motherhood on employment among U.S. resident women ages 21–35?

Woman 1

Woman 2

Woman 3

Woman 4



Causal effect in a population

What is the causal effect of motherhood on employment among U.S. resident women ages 21–35?

	Would be employed if a mother? $Y(1)$
Woman 1	0
Woman 2	0
Woman 3	0
Woman 4	1



Causal effect in a population

What is the causal effect of motherhood on employment among U.S. resident women ages 21–35?

	Would be employed if a mother? $Y(1)$	Would be employed if a non-mother? $Y(0)$
Woman 1	0	1
Woman 2	0	0
Woman 3	0	1
Woman 4	1	1



Causal effect in a population

What is the causal effect of motherhood on employment among U.S. resident women ages 21–35?

	Would be employed if a mother? $Y(1)$	Would be employed if a non-mother? $Y(0)$	Causal effect $Y(1) - Y(0)$
Woman 1	0	1	-1
Woman 2	0	0	0
Woman 3	0	1	-1
Woman 4	1	1	0

Why model?

Causal effect in a population

What is the causal effect of motherhood on employment among U.S. resident women ages 21–35?

	Would be employed if a mother? $Y(1)$	Would be employed if a non-mother? $Y(0)$	Causal effect $Y(1) - Y(0)$
Woman 1	0	1	-1
Woman 2	0	0	0
Woman 3	0	1	-1
Woman 4	1	1	0

Why model?

Causal effect in a population

What is the causal effect of motherhood on employment among U.S. resident women ages 21–35?

	Would be employed if a mother? $Y(1)$	Would be employed if a non-mother? $Y(0)$	Causal effect $Y(1) - Y(0)$
Woman 1	?	1	?
Woman 2	?	0	?
Woman 3	0	?	?
Woman 4	1	?	?

Why model?

Describe population subgroups

What is the proportion employed among U.S. resident women ages 21–35, comparing mothers to non-mothers?

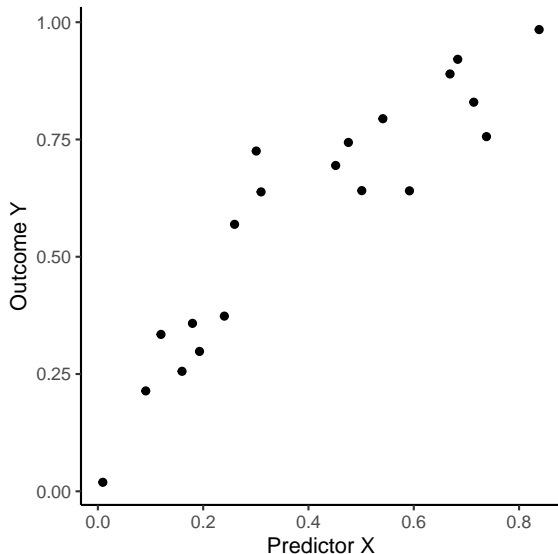
	<u>Employed?</u>		<u>Employed?</u>
Mother 1	0	Non-Mother 1	1
Mother 2	0	Non-Mother 2	0
Mother 3	0	Non-Mother 3	1
Mother 4	1	Non-Mother 4	1

A \hat{Y} view of description

With Kristin Liao, UCLA

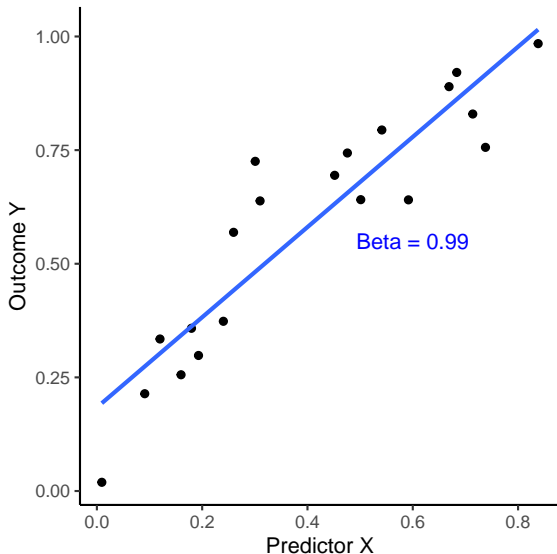
A \hat{Y} view of description

With Kristin Liao, UCLA



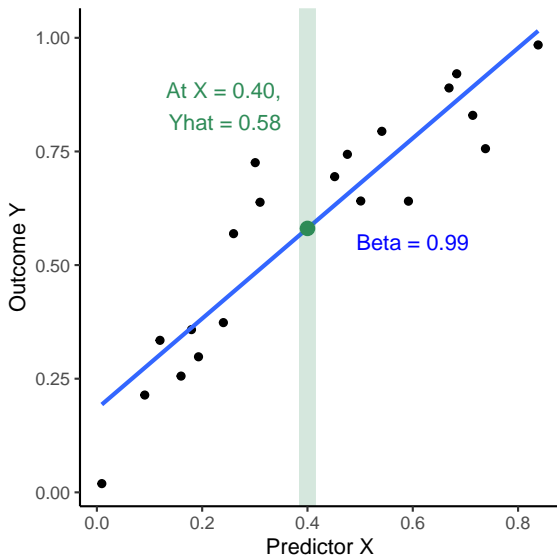
A \hat{Y} view of description

With Kristin Liao, UCLA



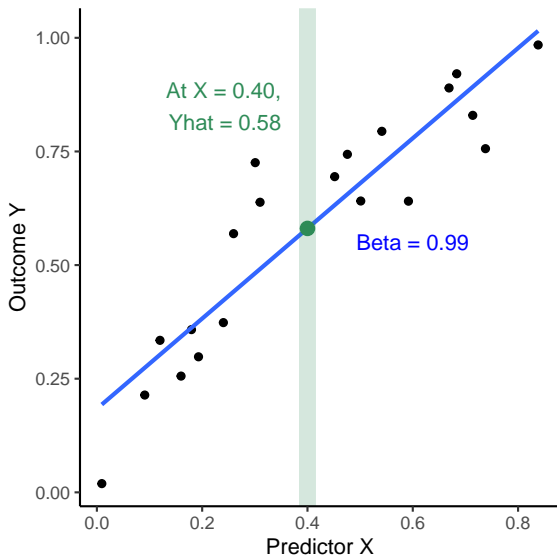
A \hat{Y} view of description

With Kristin Liao, UCLA



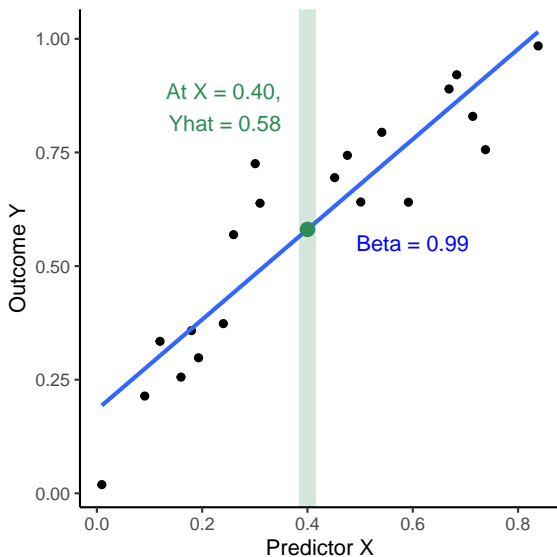
A \hat{Y} view of description

With Kristin Liao, UCLA



A \hat{Y} view of description

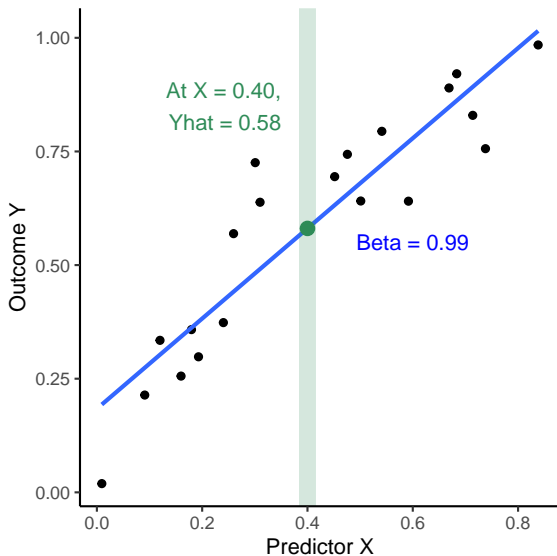
With Kristin Liao, UCLA



Why model?

A \hat{Y} view of description

With Kristin Liao, UCLA

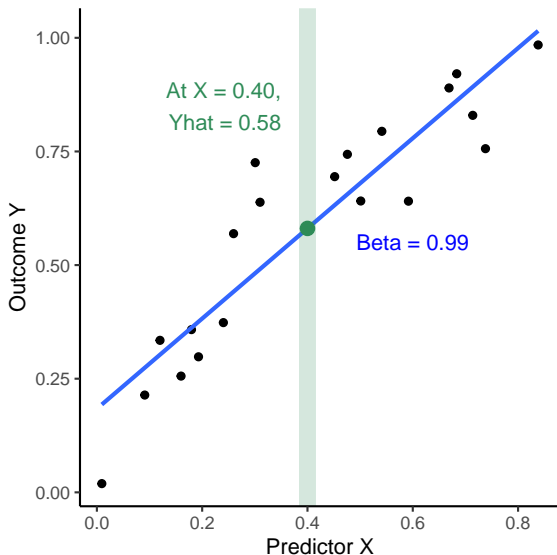


Why model?

A subgroups may have few units

A \hat{Y} view of description

With Kristin Liao, UCLA



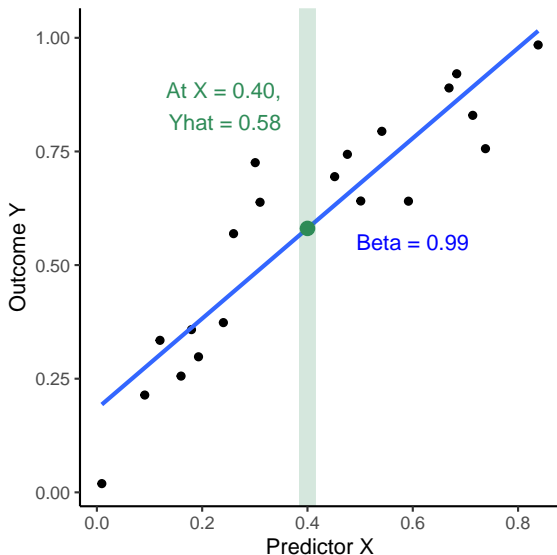
Why model?

A subgroups may have few units

Model pools information across subgroups

A \hat{Y} view of description

With Kristin Liao, UCLA



Why model?

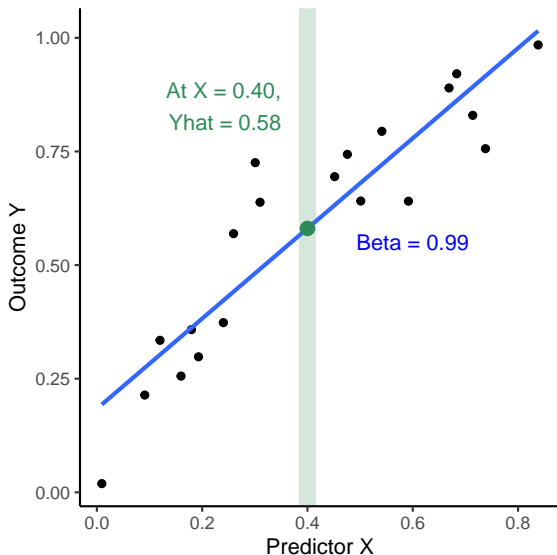
A subgroups may have few units

Model pools information across subgroups

Report \hat{Y} , not $\hat{\beta}$

A \hat{Y} view of description

With Kristin Liao, UCLA



Why model?

A subgroups may have few units

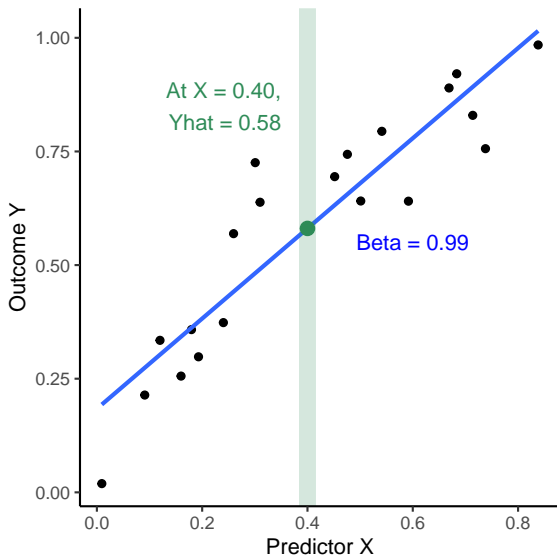
Model pools information across subgroups

Report \hat{Y} , not $\hat{\beta}$

Benefits

A \hat{Y} view of description

With Kristin Liao, UCLA



Why model?

A subgroups may have few units

Model pools information across subgroups

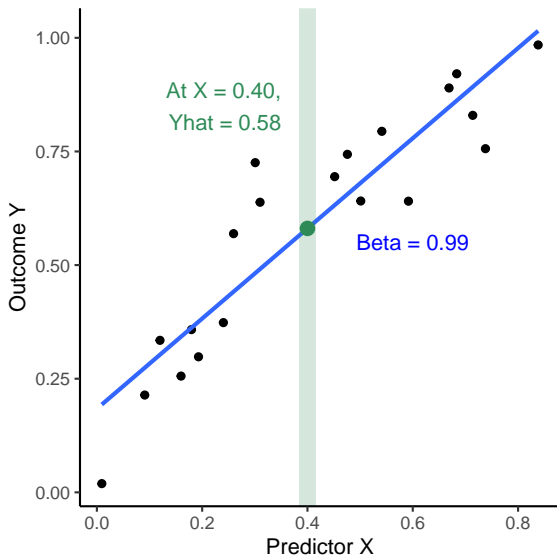
Report \hat{Y} , not $\hat{\beta}$

Benefits

Jargon-free results

A \hat{Y} view of description

With Kristin Liao, UCLA



Why model?

A subgroups may have few units

Model pools information across subgroups

Report \hat{Y} , not $\hat{\beta}$

Benefits

Jargon-free results

Plug in machine learning

Concrete exercise: Sex gap in pay

Sample of 5 million cases (true nonparametric estimates)

Simulate a sample of 100 (evaluate sample-based estimators)

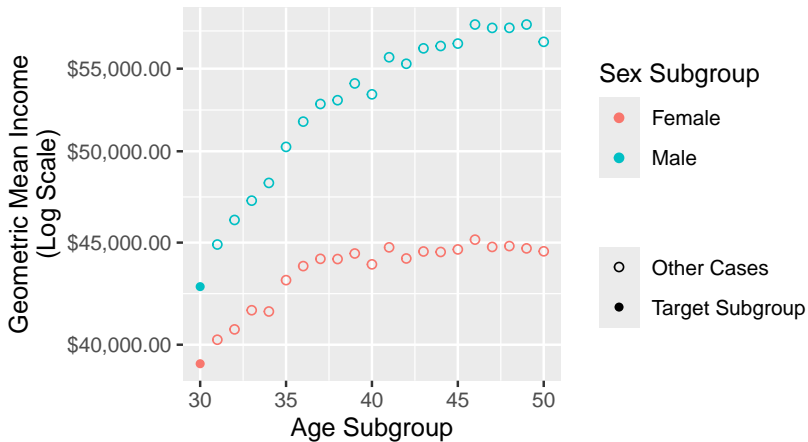
Concrete exercise: Sex gap in pay

Source of 5 million cases

- ▶ American Community Survey (ACS) 2010–2019
- ▶ Adults age 30–50
- ▶ Worked 35+ hours per week in 50+ weeks last year
- ▶ Outcome: Annual wage and salary income

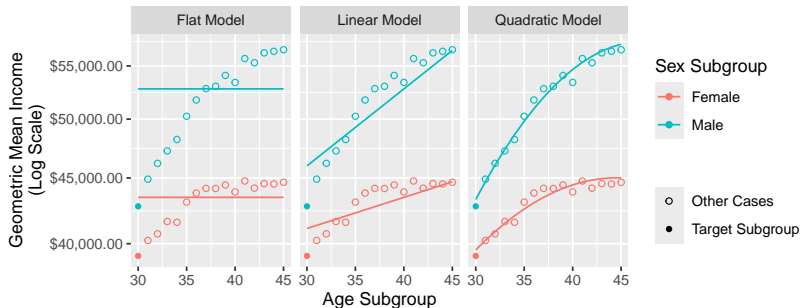
Concrete exercise: Sex gap in pay

Illustrated on full data: 5 million cases



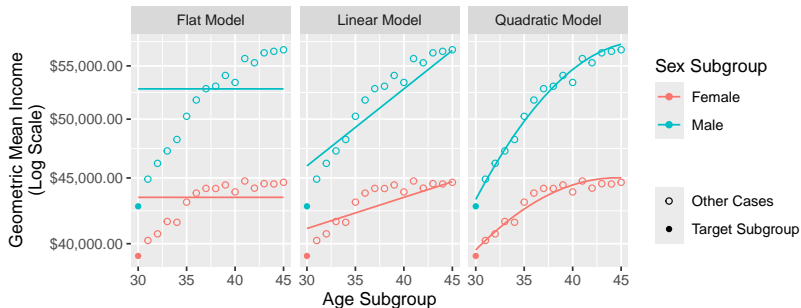
Concrete exercise: Sex gap in pay

Illustrated on full data: 5 million cases



Concrete exercise: Sex gap in pay

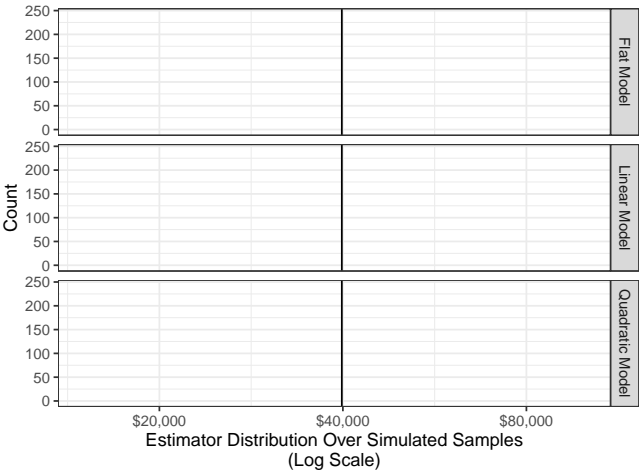
Illustrated on full data: 5 million cases



Evaluate models

Estimator Performance: Histogram Over Simulations

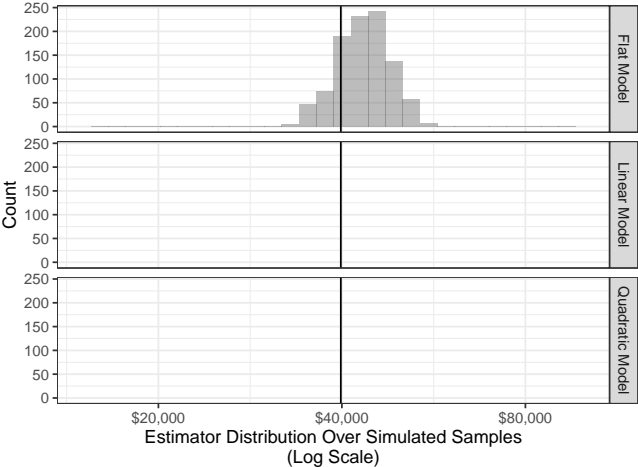
Estimand: Geometric mean income among 30-year-old female subgroup



Evaluate models

Estimator Performance: Histogram Over Simulations

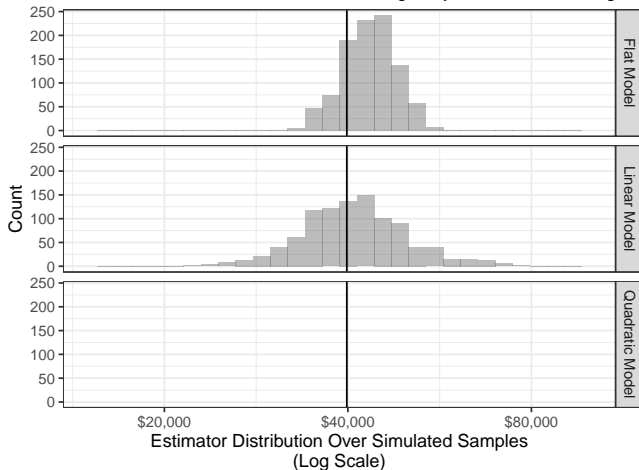
Estimand: Geometric mean income among 30-year-old female subgroup



Evaluate models

Estimator Performance: Histogram Over Simulations

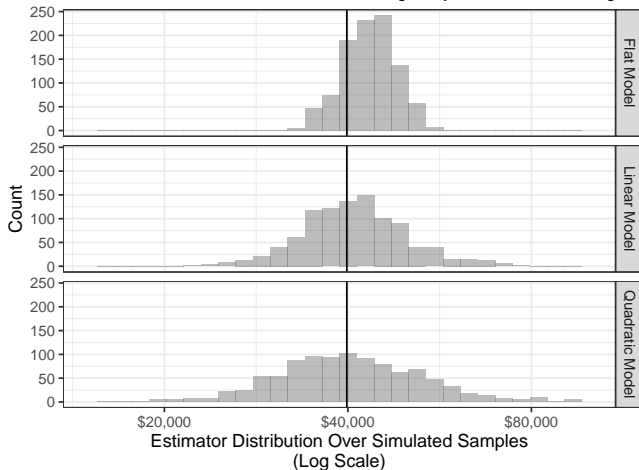
Estimand: Geometric mean income among 30-year-old female subgroup



Evaluate models

Estimator Performance: Histogram Over Simulations

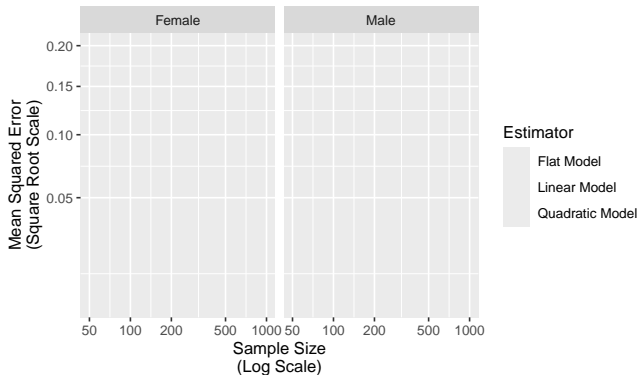
Estimand: Geometric mean income among 30-year-old female subgroup



Evaluate models

Best Estimator Depends on Estimand and Sample Size

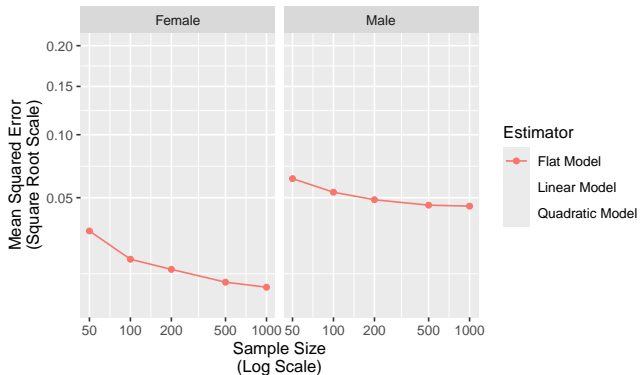
Estimand: Geometric mean income among U.S. adults age 30



Evaluate models

Best Estimator Depends on Estimand and Sample Size

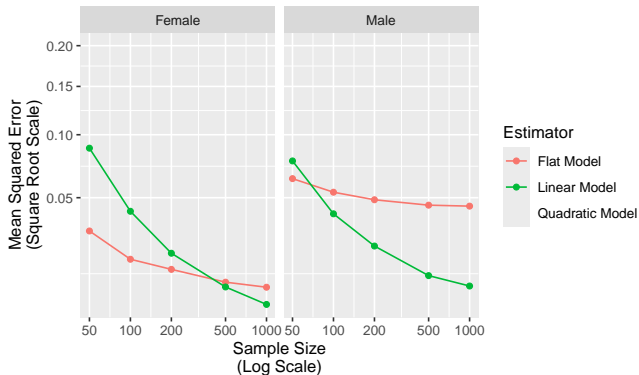
Estimand: Geometric mean income among U.S. adults age 30



Evaluate models

Best Estimator Depends on Estimand and Sample Size

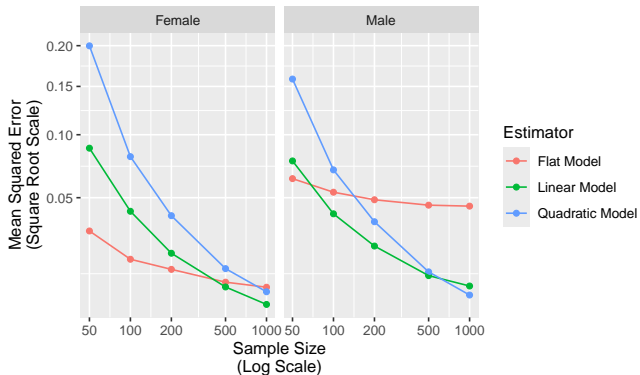
Estimand: Geometric mean income among U.S. adults age 30



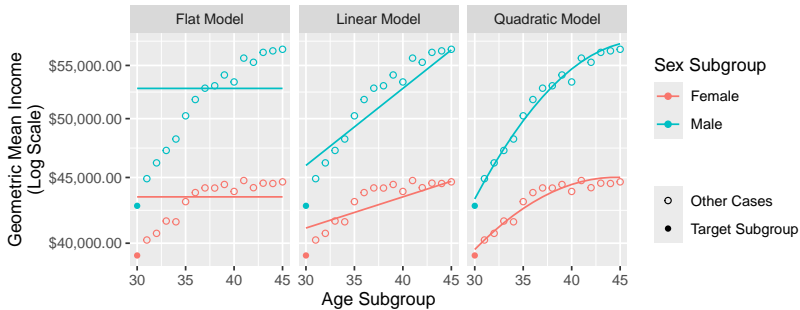
Evaluate models

Best Estimator Depends on Estimand and Sample Size

Estimand: Geometric mean income among U.S. adults age 30



Illustrated on full data: 5 million cases



Implications of a \hat{Y} view of description

Implications of a \hat{Y} view of description

- ▶ a model as a means to an end
 - ▶ we would rather not model
 - ▶ model only when you lack data

Implications of a \hat{Y} view of description

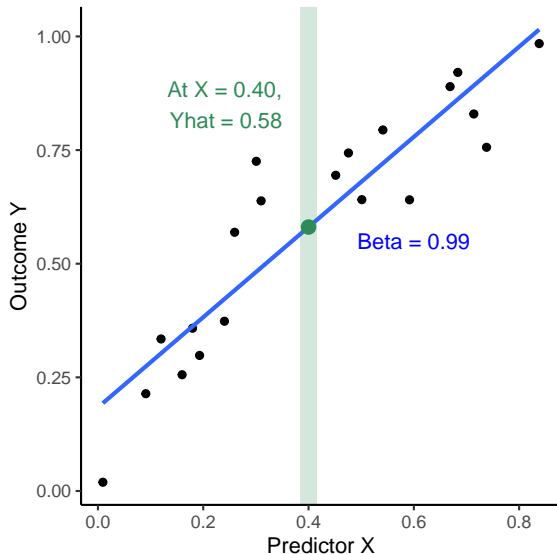
- ▶ a model as a means to an end
 - ▶ we would rather not model
 - ▶ model only when you lack data
- ▶ misspecified models are ok
 - ▶ flat model was wrong
 - ▶ flat model was best

(lower variance)

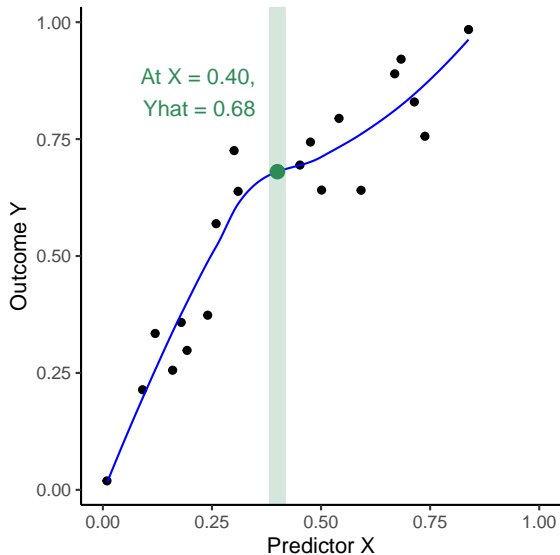
Implications of a \hat{Y} view of description

- ▶ a model as a means to an end
 - ▶ we would rather not model
 - ▶ model only when you lack data
- ▶ misspecified models are ok
 - ▶ flat model was wrong
 - ▶ flat model was best (lower variance)
- ▶ machine learning becomes a plug-in

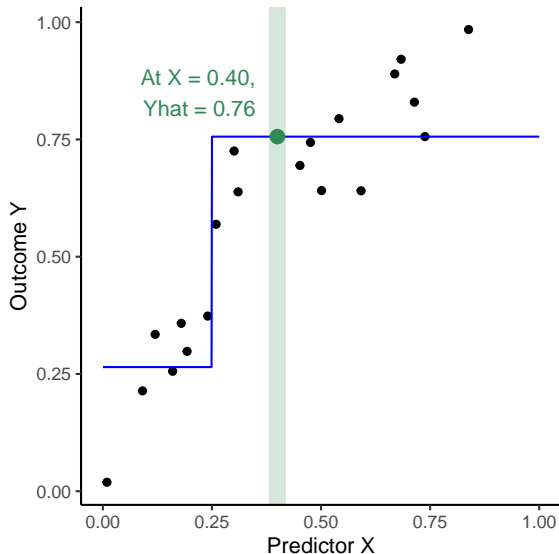
With \hat{Y} description,
machine learning becomes a plug-in



With \hat{Y} description,
machine learning becomes a plug-in



With \hat{Y} description,
machine learning becomes a plug-in



Computer tutorial: Introduction

ilundberg.github.io/description

Computer tutorial: Introduction

ilundberg.github.io/description

We will give you data:

- ▶ male and female incomes at age 30–50 in 2010–2019

You will make a forecast:

- ▶ male and female geometric mean income at age 30–50 in 2022

Computer tutorial: Introduction

ilundberg.github.io/description

Prepare the environment by loading the `tidyverse` package.

```
library(tidyverse)
```

The function below simulates a sample of 100 cases.

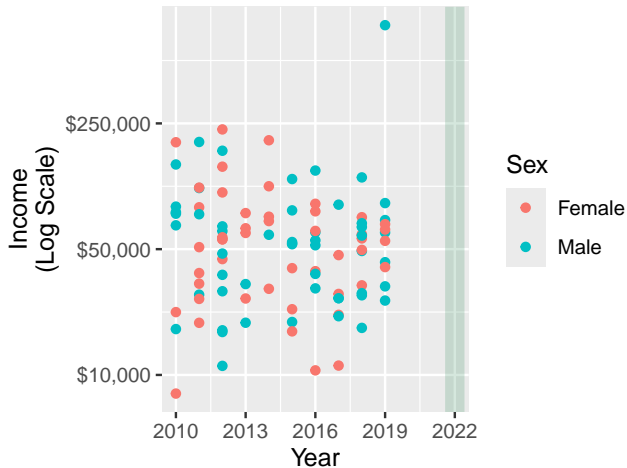
```
simulate <- function(n = 100) {  
  read_csv("https://ilundberg.github.io/description/assets/truth.csv") |>  
  slice_sample(n = n, weight_by = weight, replace = T) |>  
  mutate(income = exp(rnorm(n(), meanlog, sdlog))) |>  
  select(year, age, sex, income)  
}
```

We can see how it works below.

```
simulated <- simulate(n = 100)
```

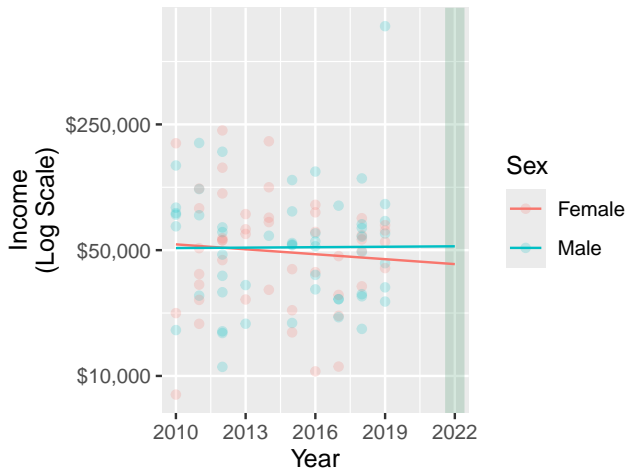
Computer tutorial: Introduction

ilundberg.github.io/description



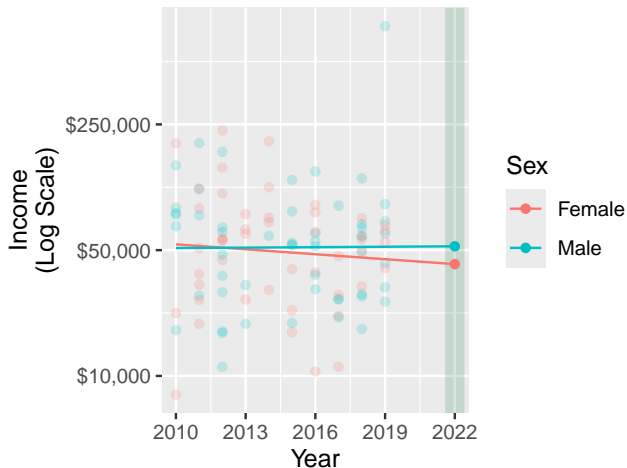
Computer tutorial: Introduction

ilundberg.github.io/description



Computer tutorial: Introduction

ilundberg.github.io/description



Computer tutorial: Introduction

ilundberg.github.io/description

We will give you data:

- ▶ male and female incomes at age 30–50 in 2010–2019

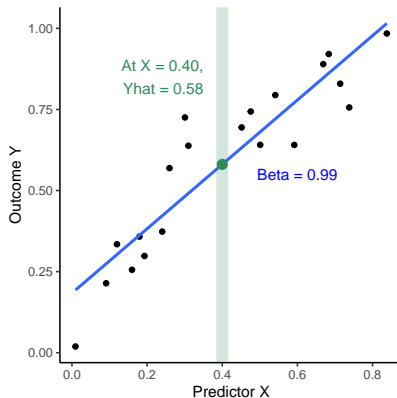
You will make a forecast:

- ▶ male and female geometric mean income at age 30–50 in 2022

We will see who comes closest

- ▶ to gold-standard truth from ACS 2022

Thanks!



Ian Lundberg
ianlundberg.org
ianlundberg@ucla.edu

Kristin Liao
kristinliao.com
ktliao@ucla.edu