

# Machine Learning for Healthcare

HST.956, 6.S897

## Lecture 14: Causal Inference Part 1

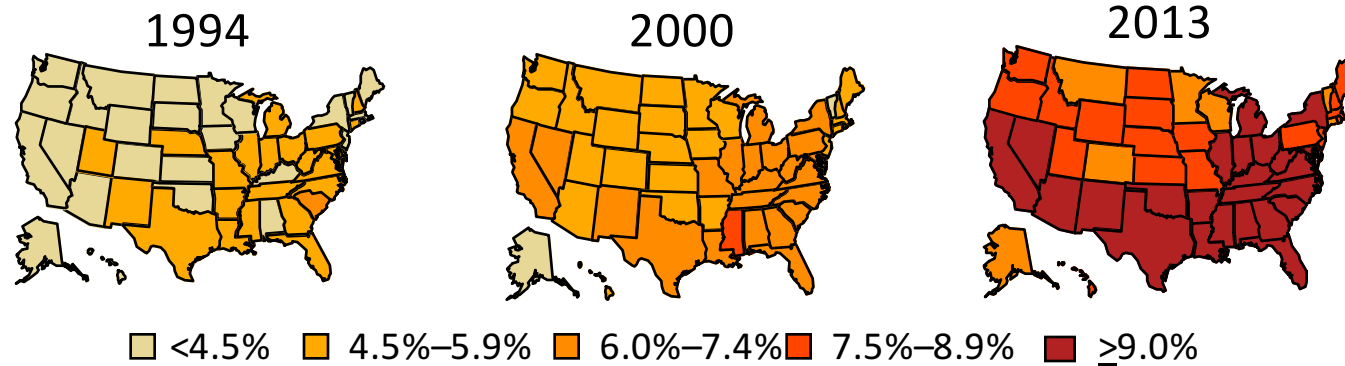
David Sontag



# Course announcements

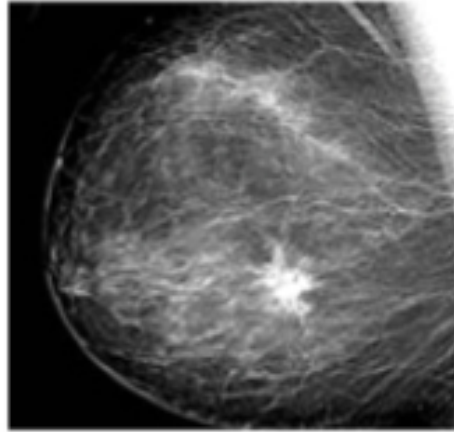
- **Please fill out mid-semester survey**
- **Project proposals**
  - You will receive e-mail feedback this week
  - Office hours next Tuesday, 10-11:30am
- **Problem sets**
  - PS1-4 graded (see Stellar)
  - PS5 out tonight, due next Tuesday, April 9
  - Last problem set, PS6, released in ~2 weeks
- **Recitation this week will be a discussion of**
  - Brat et al., Postsurgical prescriptions for opioid naïve patients and association with overdose and misuse, BMJ 2018
  - Bertsimas et al., Personalized diabetes management using electronic medical records, Diabetes Care 2017

# Does gastric bypass surgery prevent onset of diabetes?

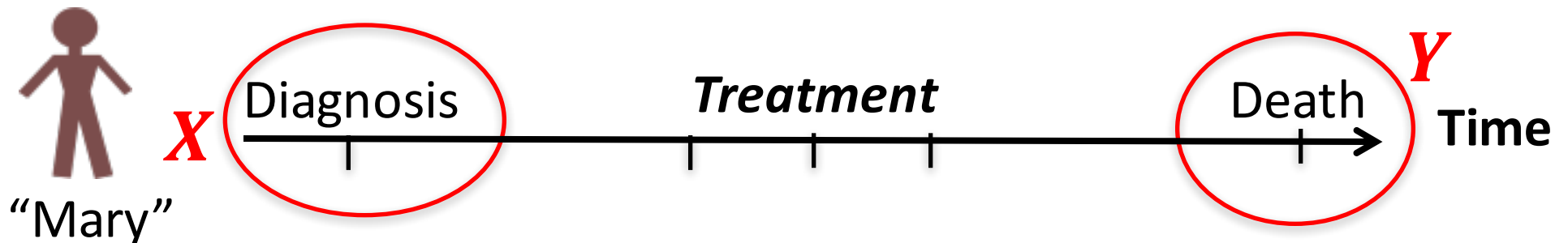


- In Lecture 4 & PS2 we used machine learning for early detection of Type 2 diabetes
- Health system doesn't want to know how to predict diabetes – they want to know how to *prevent it*
- Gastric bypass surgery is the highest negative weight (9th most predictive feature)
  - Does this mean it would be a good intervention?

What is the likelihood this patient, with breast cancer, will survive 5 years?

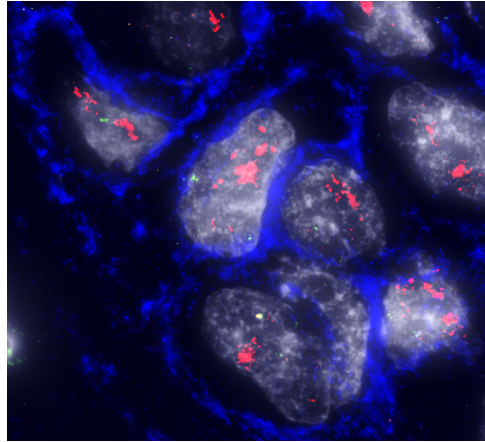


- Such predictive models widely used to stage patients. Should we initiate treatment? How aggressive?
- What could go wrong if we trained to predict survival, and then used to guide patient care?



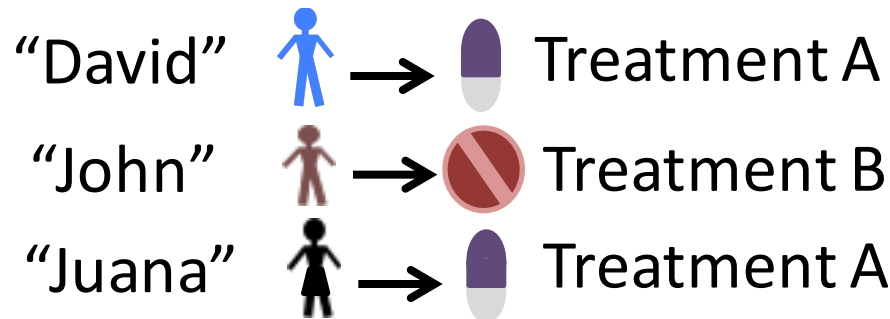
**A long survival time may be because of treatment!**

# What treatment should we give this patient?



Expansion pathology  
(image from Andy Beck)

- People respond differently to treatment
- Goal: use data from other patients and their journeys to guide future treatment decisions
- What could go wrong if we trained to predict (past) treatment decisions?



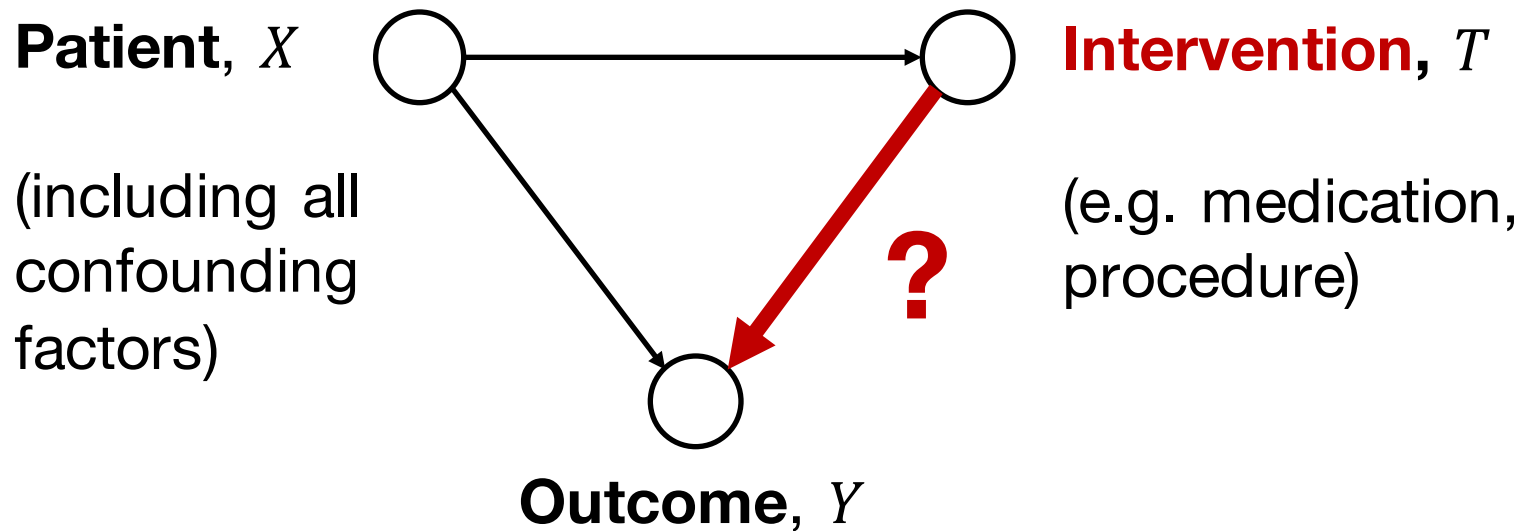
**Best this can do is  
match current  
medical practice!**

# Does smoking cause lung cancer?



- Doing a randomized control trial is unethical
- Could we simply answer this question by comparing  $\Pr(\text{lung cancer} \mid \text{smoker})$  vs  $\Pr(\text{lung cancer} \mid \text{nonsmoker})$ ?
- **No! Answering such questions from observational data is difficult because of *confounding***

To properly answer, need to formulate as *causal* questions:



*High dimensional*

*Observational data*

# Potential Outcomes Framework (Rubin-Neyman Causal Model)

- Each unit (individual)  $x_i$  has two potential outcomes:
  - $Y_0(x_i)$  is the potential outcome had the unit not been treated:  
“**control outcome**”
  - $Y_1(x_i)$  is the potential outcome had the unit been treated:  
“**treated outcome**”
- Conditional average treatment effect for unit  $i$ :  
$$CATE(x_i) = \mathbb{E}_{Y_1 \sim p(Y_1|x_i)} [Y_1|x_i] - \mathbb{E}_{Y_0 \sim p(Y_0|x_i)} [Y_0|x_i]$$
- Average Treatment Effect:  
$$ATE := \mathbb{E}[Y_1 - Y_0] = \mathbb{E}_{x \sim p(x)} [CATE(x)]$$



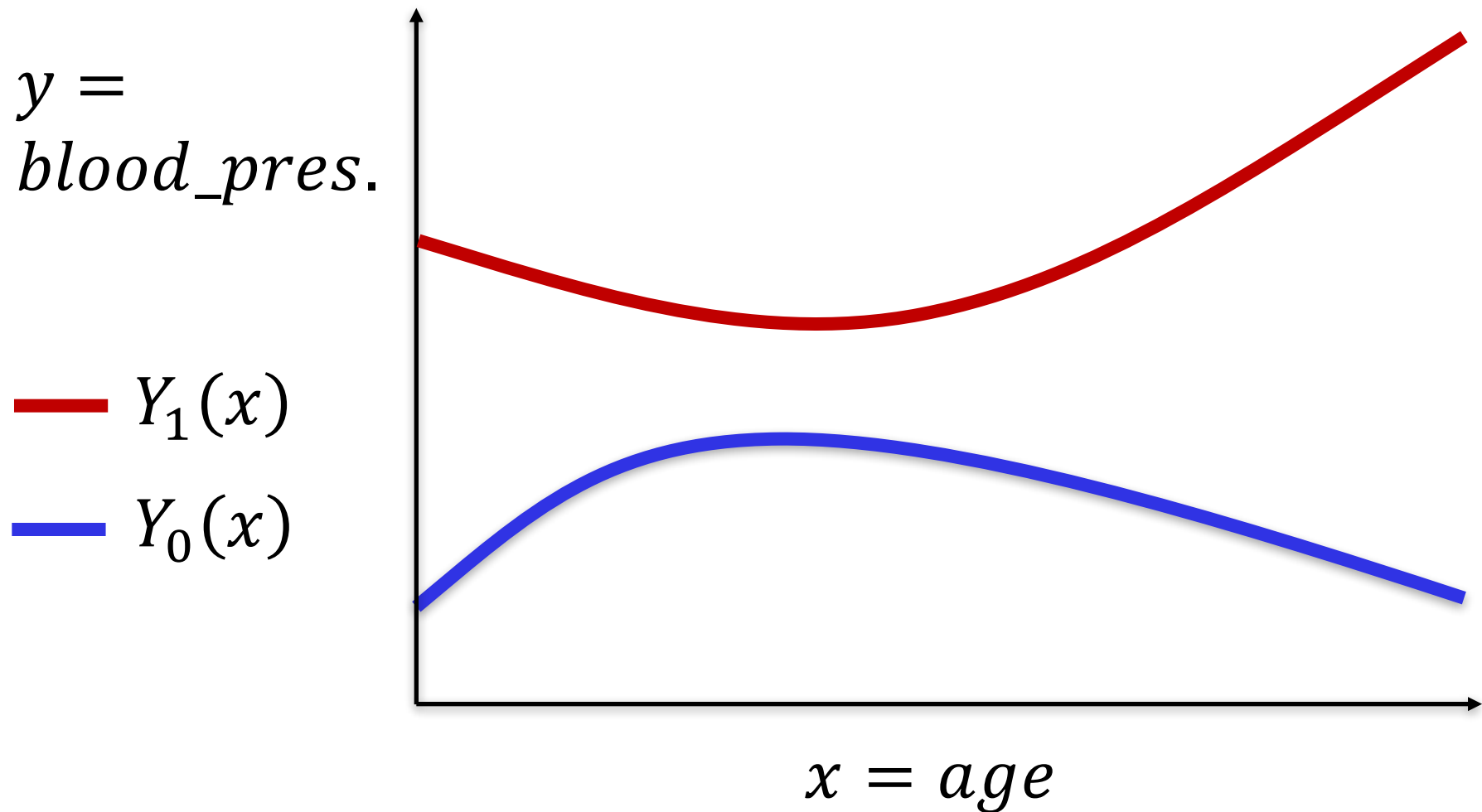
# Potential Outcomes Framework (Rubin-Neyman Causal Model)

- Each unit (individual)  $x_i$  has two potential outcomes:
  - $Y_0(x_i)$  is the potential outcome had the unit not been treated:  
“**control outcome**”
  - $Y_1(x_i)$  is the potential outcome had the unit been treated:  
“**treated outcome**”
- Observed factual outcome:  
$$y_i = t_i Y_1(x_i) + (1 - t_i) Y_0(x_i)$$
- Unobserved counterfactual outcome:  
$$y_i^{CF} = (1 - t_i) Y_1(x_i) + t_i Y_0(x_i)$$

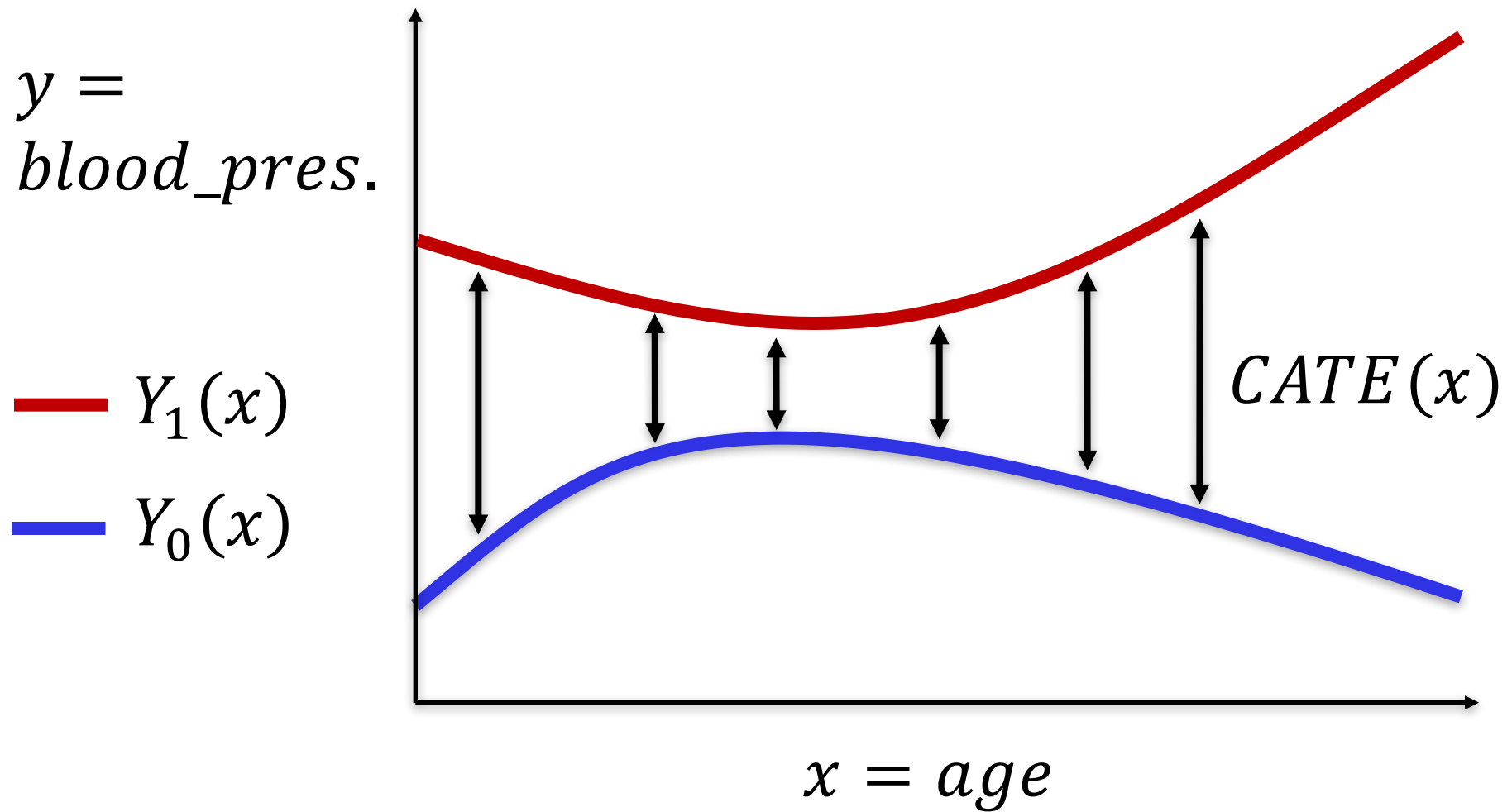
“The fundamental problem of  
causal inference”

We only ever observe one of the  
two outcomes

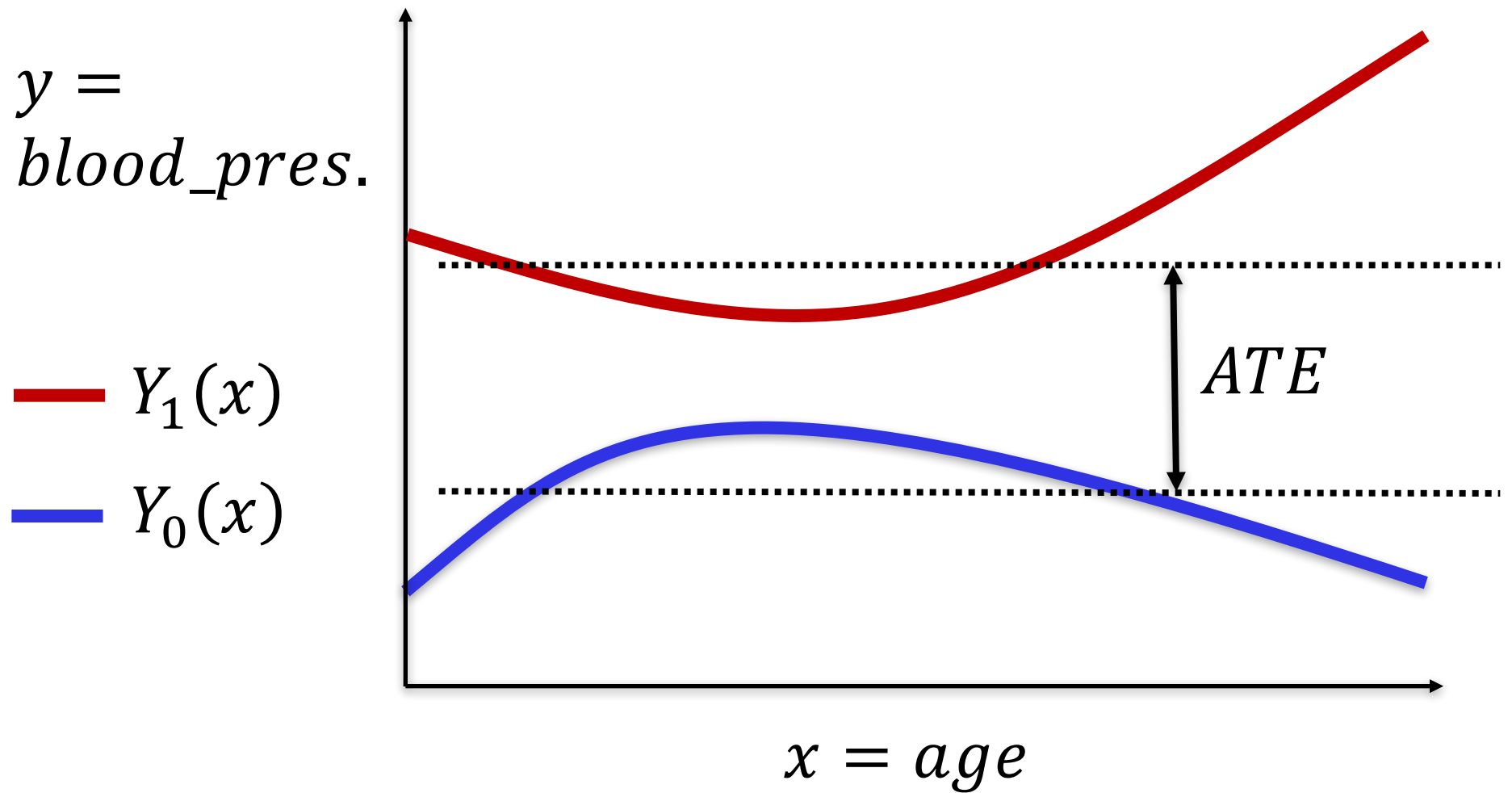
# Example – Blood pressure and age



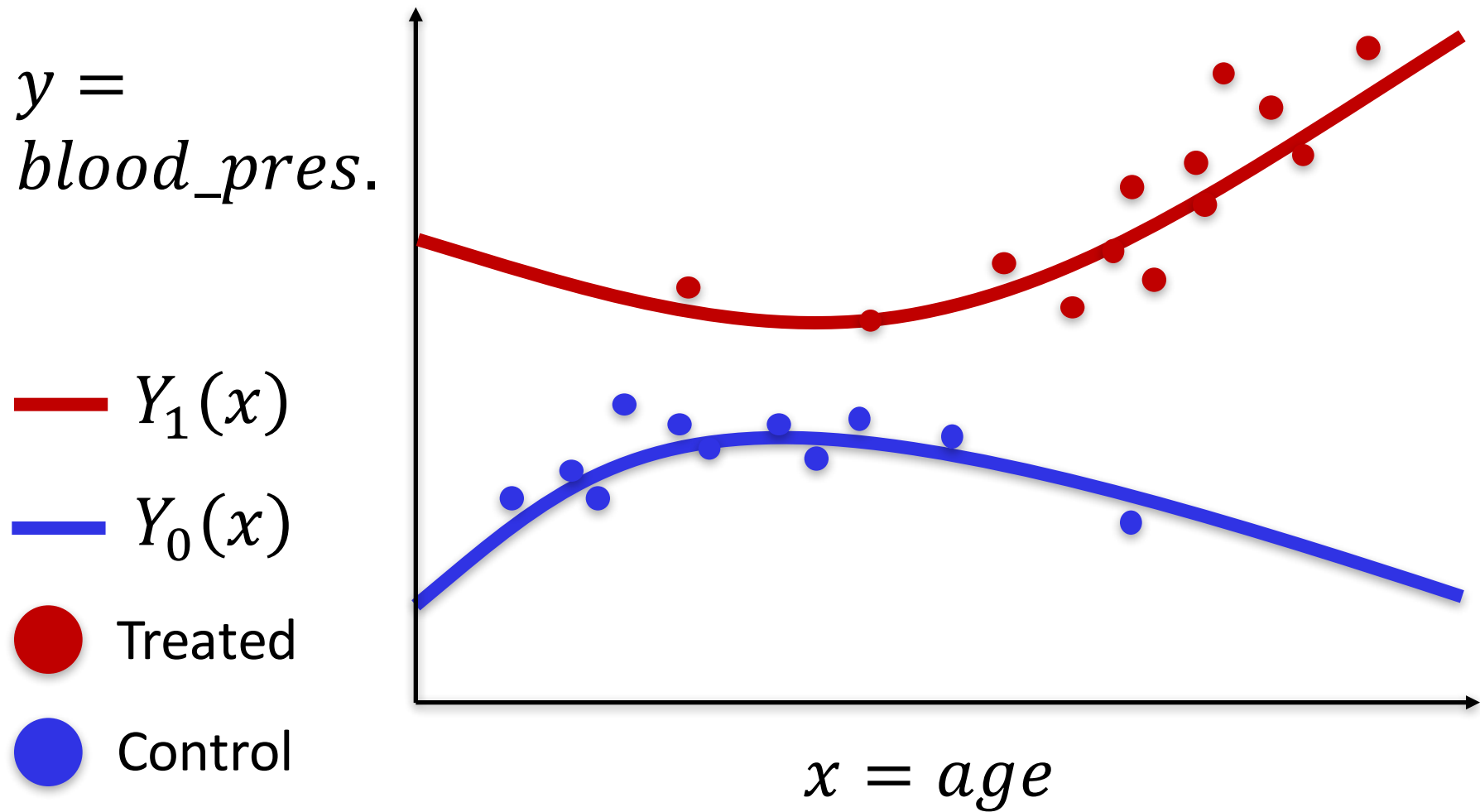
# Blood pressure and age



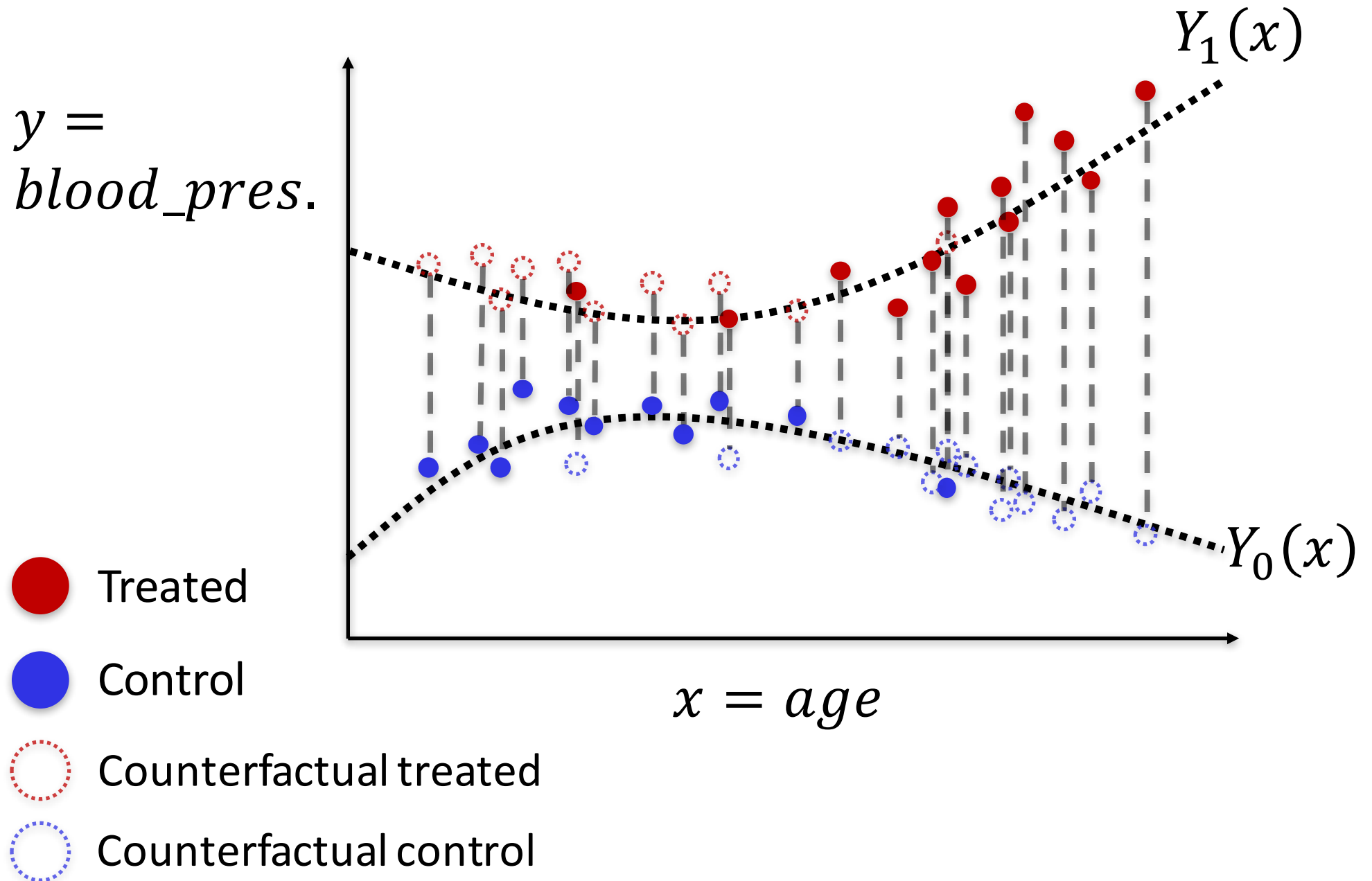
# Blood pressure and age



# Blood pressure and age



# Blood pressure and age



(age, gender, exercise, treatment)			Observed sugar levels
(45, F, 0, <b>A</b> )			6
(45, F, 1, <b>B</b> )			6.5
(55, M, 0, <b>A</b> )			7
(55, M, 1, <b>B</b> )			8
(65, F, 0, <b>B</b> )			8
(65, F, 1, <b>A</b> )			7.5
(75, M, 0, <b>B</b> )			9
(75, M, 1, <b>A</b> )			8

(Example from Uri Shalit)



(age, gender, exercise)			Observed sugar levels
(45, F, 0)			6
(45, F, 1)			6.5
(55, M, 0)			7
(55, M, 1)			8
(65, F, 0)			8
(65, F, 1)			7.5
(75, M, 0)			9
(75, M, 1)			8

(Example from Uri Shalit)

(age, gender, exercise)	$Y_0$ : Sugar levels <i>had they received medication A</i>	$Y_1$ : Sugar levels <i>had they received medication B</i>	Observed sugar levels
(45, F, 0)	<b>6</b>	5.5	6
(45, F, 1)	7	<b>6.5</b>	6.5
(55, M, 0)	<b>7</b>	6	7
(55, M, 1)	9	<b>8</b>	8
(65, F, 0)	8.5	<b>8</b>	8
(65, F, 1)	<b>7.5</b>	7	7.5
(75, M, 0)	10	<b>9</b>	9
(75, M, 1)	<b>8</b>	7	8

(Example from Uri Shalit)

(age,gender, exercise)	Sugar levels <i>had they received medication A</i>	Sugar levels <i>had they received medication B</i>	Observed sugar levels
(45, F, 0)	<b>6</b>	5.5	6
(45, F, 1)	7	<b>6.5</b>	6.5
(55, M, 0)	<b>7</b>	6	7
(55, M, 1)	9	<b>8</b>	8
(65, F, 0)	8.5	<b>8</b>	8
(65,F, 1)	<b>7.5</b>	7	7.5
(75,M, 0)	10	<b>9</b>	9
(75,M, 1)	<b>8</b>	7	8

mean(sugar | medication B) –  
mean(sugar | medicaton A) =  
?

mean(sugar | *had they received* B) –  
mean(sugar | *had they received* A) =  
?

(Example from Uri Shalit)

(age,gender, exercise)	Sugar levels <i>had they received medication A</i>	Sugar levels <i>had they received medication B</i>	Observed sugar levels
(45, F, 0)	<b>6</b>	5.5	6
(45, F, 1)	7	<b>6.5</b>	6.5
(55, M, 0)	<b>7</b>	6	7
(55, M, 1)	9	<b>8</b>	8
(65, F, 0)	8.5	<b>8</b>	8
(65,F, 1)	<b>7.5</b>	7	7.5
(75,M, 0)	10	<b>9</b>	9
(75,M, 1)	<b>8</b>	7	8

$$\text{mean}(\text{sugar} \mid \text{medication B}) - \text{mean}(\text{sugar} \mid \text{medication A}) = 7.875 - 7.125 = 0.75$$

$$\text{mean}(\text{sugar} \mid \textit{had they received B}) - \text{mean}(\text{sugar} \mid \textit{had they received A}) = 7.125 - 7.875 = -0.75$$

(Example from Uri Shalit)

# Typical assumption – no unmeasured confounders

$Y_0, Y_1$ : potential outcomes for control and treated

$x$ : unit covariates (features)

$T$ : treatment assignment

We assume:

$$(Y_0, Y_1) \perp\!\!\!\perp T \mid x$$

The potential outcomes are independent of treatment assignment, conditioned on covariates  $x$

# Typical assumption – no unmeasured confounders

$Y_0, Y_1$ : potential outcomes for control and treated

$x$ : unit covariates (features)

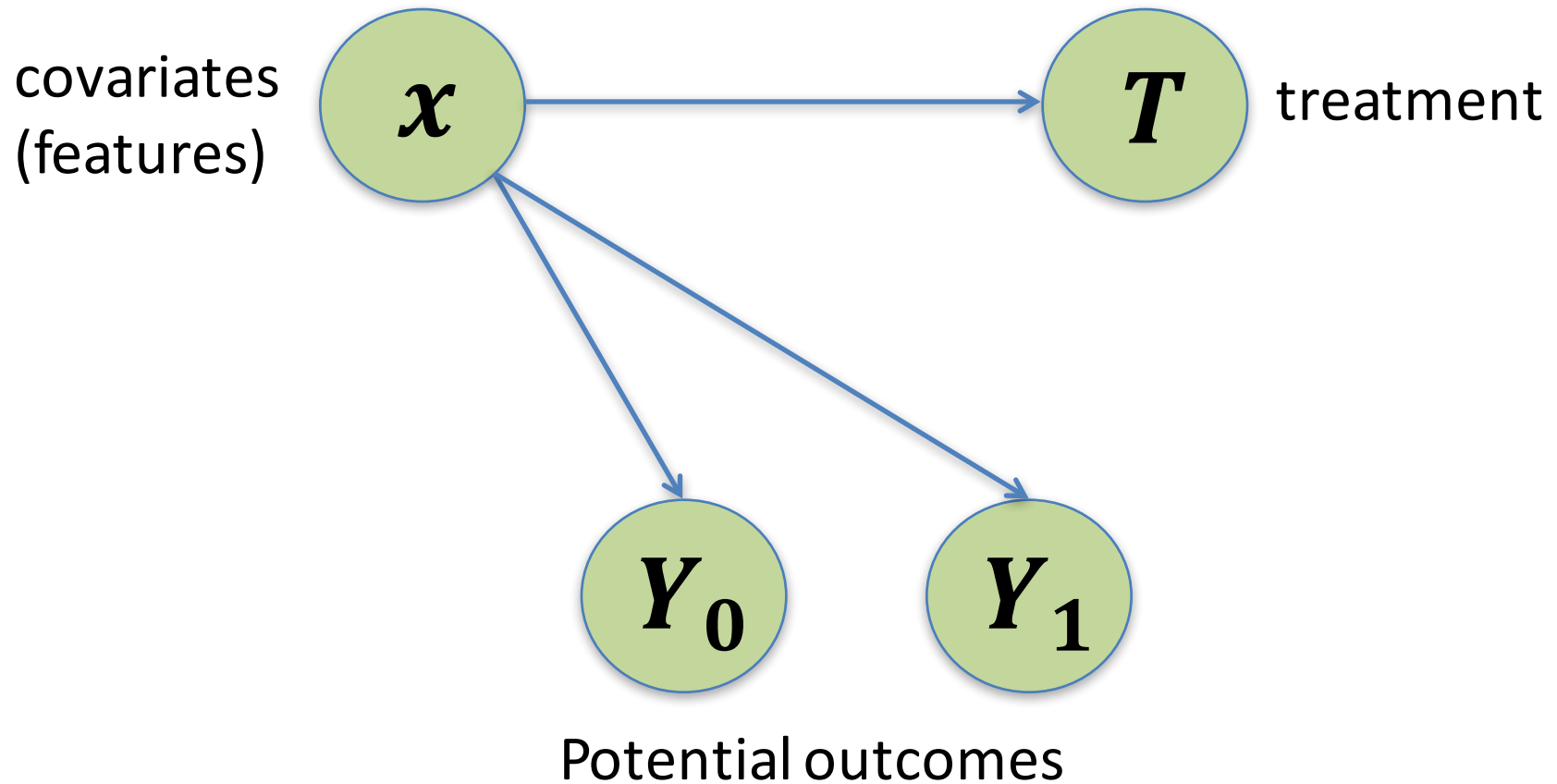
$T$ : treatment assignment

We assume:

$$(Y_0, Y_1) \perp\!\!\!\perp T \mid x$$

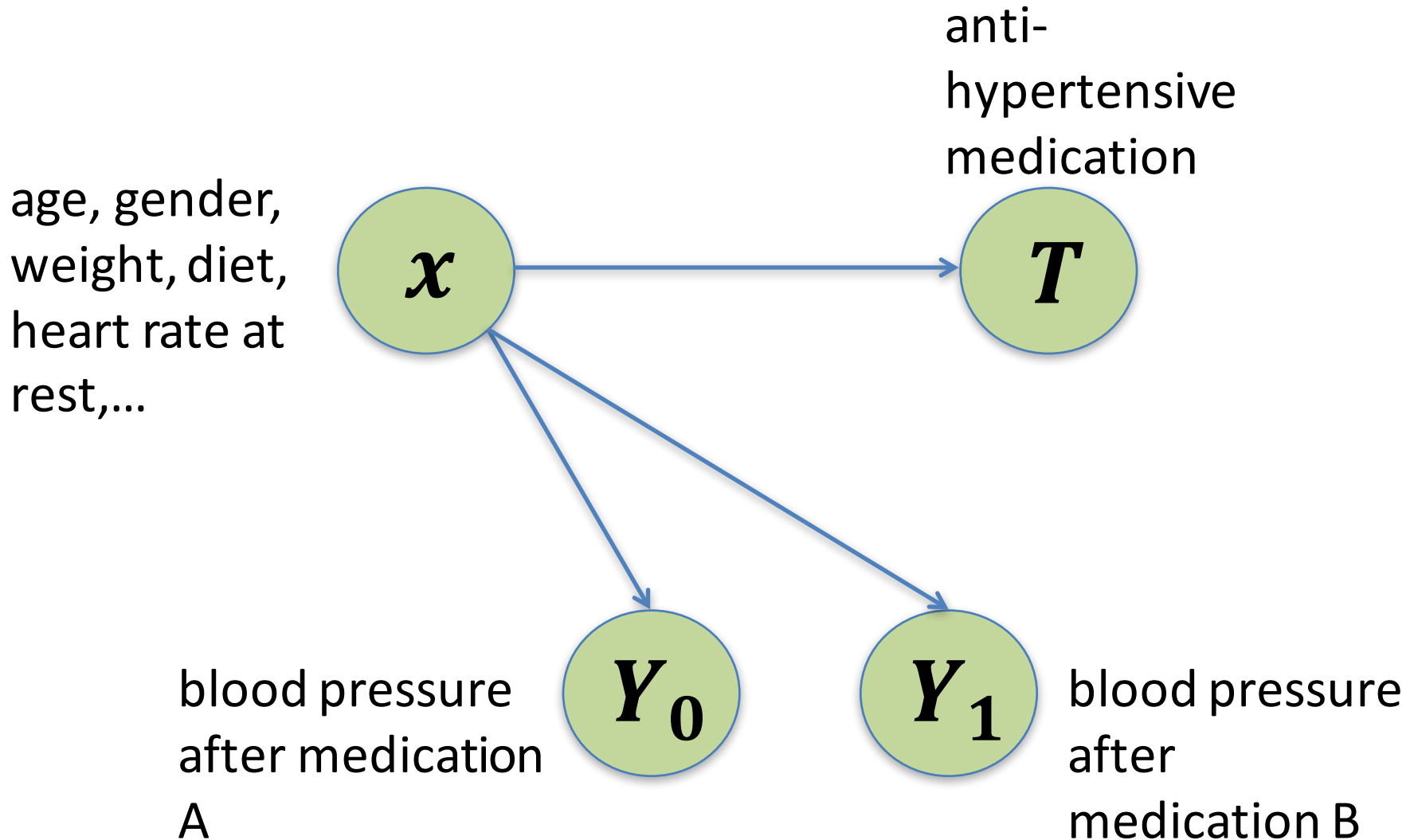
*Ignorability*

# Ignorability



$$(Y_0, Y_1) \perp\!\!\!\perp T \mid x$$

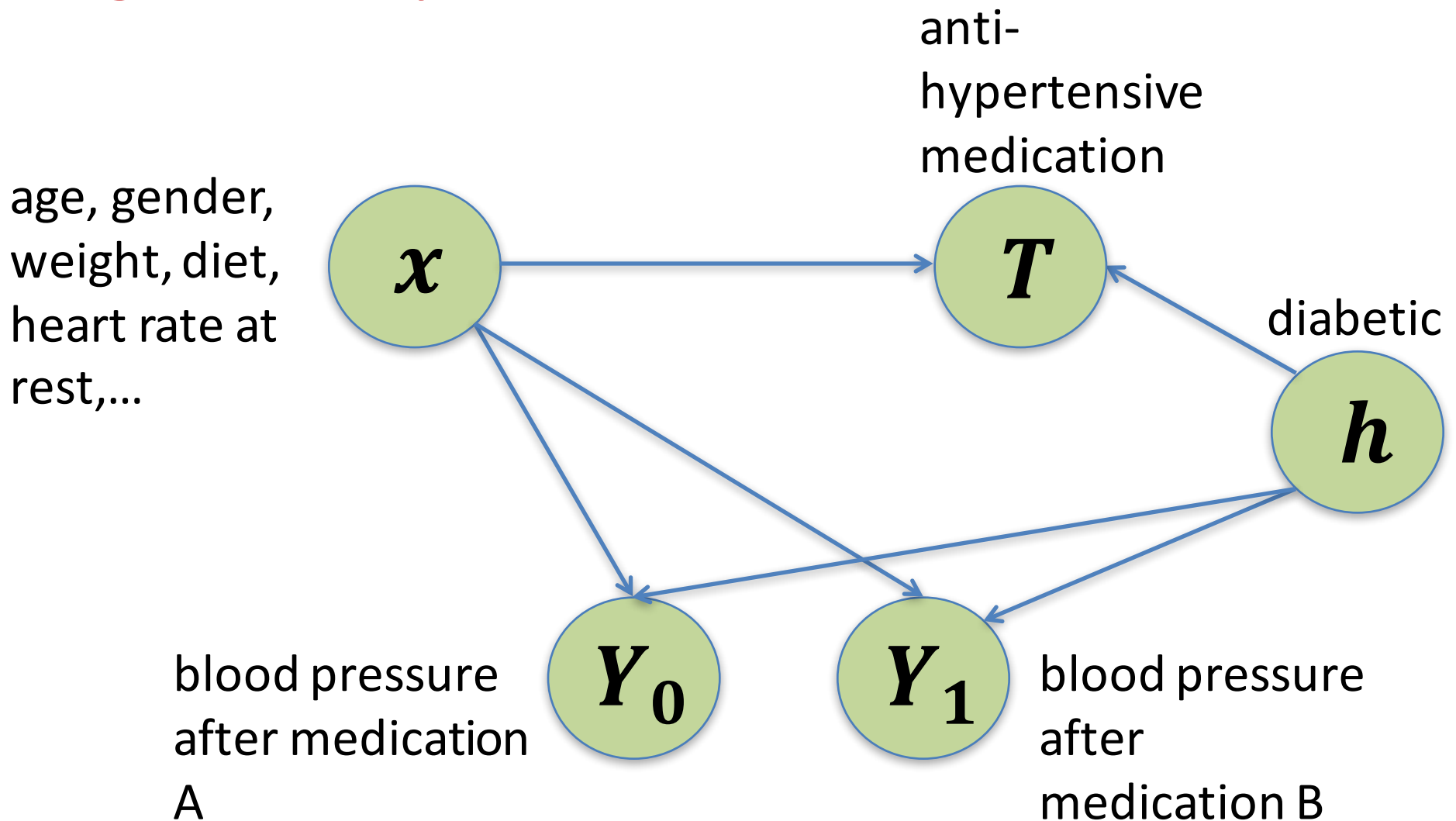
# Ignorability



$$(Y_0, Y_1) \perp\!\!\!\perp T \mid X$$



# No Ignorability



$$(Y_0, Y_1) \not\perp T \mid x$$

## Typical assumption – common support

$Y_0, Y_1$ : potential outcomes for control and treated

$x$ : unit covariates (features)

$T$ : treatment assignment

We assume:

$$p(T = t | X = x) > 0 \quad \forall t, x$$

# Framing the question

1. Where could we go to for data to answer these questions?
2. What should  $\mathbf{X}$ ,  $T$ , and  $Y$  be to satisfy ignorability?
3. What is the specific causal inference question that we are interested in?
4. Are you worried about common support?

# Outline for lecture

- How to recognize a causal inference problem
- Potential outcomes framework
  - Average treatment effect (ATE)
  - Conditional average treatment effect (CATE)
- **Algorithms for estimating ATE and CATE**

# Average Treatment Effect

The expected causal effect of  $T$  on  $Y$ :

$$ATE := \mathbb{E} [Y_1 - Y_0]$$

# Average Treatment Effect – the adjustment formula

- Assuming ignorability, we will derive the *adjustment formula* (Hernán & Robins 2010, Pearl 2009)
- The adjustment formula is extremely useful in causal inference
- Also called *G-formula*

# Average Treatment Effect

The expected causal effect of  $T$  on  $Y$ :

$$ATE := \mathbb{E} [Y_1 - Y_0]$$

# Average Treatment Effect

The expected causal effect of  $T$  on  $Y$ :

$$ATE := \mathbb{E} [Y_1 - Y_0]$$

$$\mathbb{E} [Y_1] =$$

law of total  
expectation

$$\mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{Y_1 \sim p(Y_1|x)} [Y_1 | x] \right] =$$



# Average Treatment Effect

The expected causal effect of  $T$  on  $Y$ :

$$ATE := \mathbb{E} [Y_1 - Y_0]$$

$$\mathbb{E} [Y_1] =$$

$$\mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{Y_1 \sim p(Y_1 | x)} [Y_1 | x] \right] = \text{ignorability} \\ (Y_0, Y_1) \perp\!\!\!\perp T | x$$

$$\mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{Y_1 \sim p(Y_1 | x)} [Y_1 | x, T = 1] \right] =$$

# Average Treatment Effect

The expected causal effect of  $T$  on  $Y$ :

$$ATE := \mathbb{E} [Y_1 - Y_0]$$

$$\mathbb{E} [Y_1] =$$

$$\mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{Y_1 \sim p(Y_1 | x)} [Y_1 | x] \right] =$$

$$\mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{Y_1 \sim p(Y_1 | x)} [Y_1 | x, T = 1] \right] =$$

$$\mathbb{E}_{x \sim p(x)} \left[ \mathbb{E} [Y_1 | x, T = 1] \right] \quad \text{shorter notation}$$

# Average Treatment Effect

The expected causal effect of  $T$  on  $Y$ :

$$ATE := \mathbb{E} [Y_1 - Y_0]$$

$$\mathbb{E} [Y_0] =$$

$$\mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{Y_0 \sim p(Y_0|x)} [Y_0|x] \right] =$$

$$\mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{Y_0 \sim p(Y_0|x)} [Y_0|x, T = 1] \right] =$$

$$\mathbb{E}_{x \sim p(x)} \left[ \mathbb{E} [Y_0|x, T = 0] \right]$$

# The adjustment formula

Under the assumption of ignorability, we have that:

$$ATE = \mathbb{E} [Y_1 - Y_0] =$$

$$\mathbb{E}_{x \sim p(x)} [ \mathbb{E} [Y_1 | x, T = 1] - \mathbb{E} [Y_0 | x, T = 0] ]$$

$\mathbb{E} [Y_1 | x, T = 1]$   
 $\mathbb{E} [Y_0 | x, T = 0]$  } Quantities we  
can estimate  
from data

# The adjustment formula

Under the assumption of ignorability, we have that:

$$ATE = \mathbb{E} [Y_1 - Y_0] = \mathbb{E}_{x \sim p(x)} [ \mathbb{E} [Y_1 | x, T = 1] - \mathbb{E} [Y_0 | x, T = 0] ]$$

$$\mathbb{E} [Y_0 | x, T = 1]$$

$$\mathbb{E} [Y_1 | x, T = 0]$$

$$\mathbb{E} [Y_0 | x]$$

$$\mathbb{E} [Y_1 | x]$$

Quantities we  
*cannot* directly  
estimate from data

# The adjustment formula

Under the assumption of ignorability, we have that:

$$ATE = \mathbb{E} [Y_1 - Y_0] = \mathbb{E}_{x \sim p(x)} [ \underbrace{\mathbb{E} [Y_1 | x, T = 1] - \mathbb{E} [Y_0 | x, T = 0]}_{\text{Quantities we can estimate from data}} ]$$

$$\mathbb{E} [Y_1 | x, T = 1]$$

$$\mathbb{E} [Y_0 | x, T = 0]$$

Quantities we  
can estimate  
from data

Empirically we have samples from  $p(x|T = 1)$  or  $p(x|T = 0)$ .

*Extrapolate to  $p(x)$*

# Many methods!

Covariate adjustment

Propensity score re-weighting

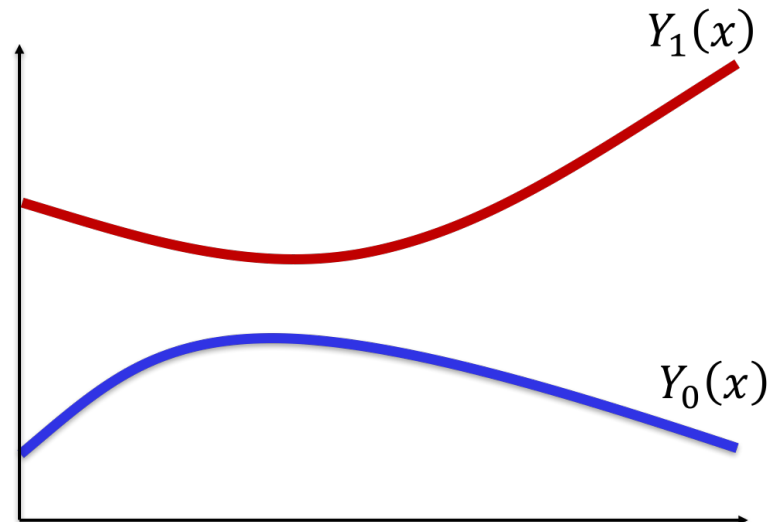
Doubly robust estimators

Matching

...

# Covariate adjustment

- Explicitly model the relationship between treatment, confounders, and outcome
- Also called “Response Surface Modeling”
- Used for both ITE and ATE
- A regression problem





Covariates  
(Features)

$x_1$

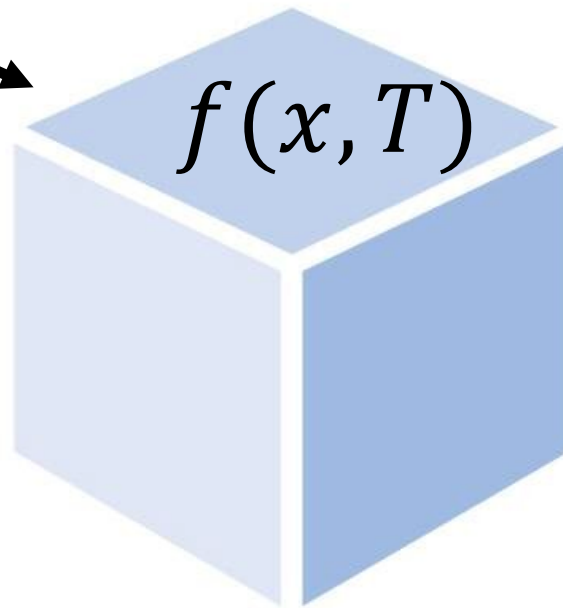
$x_2$

$\vdots$

$x_d$

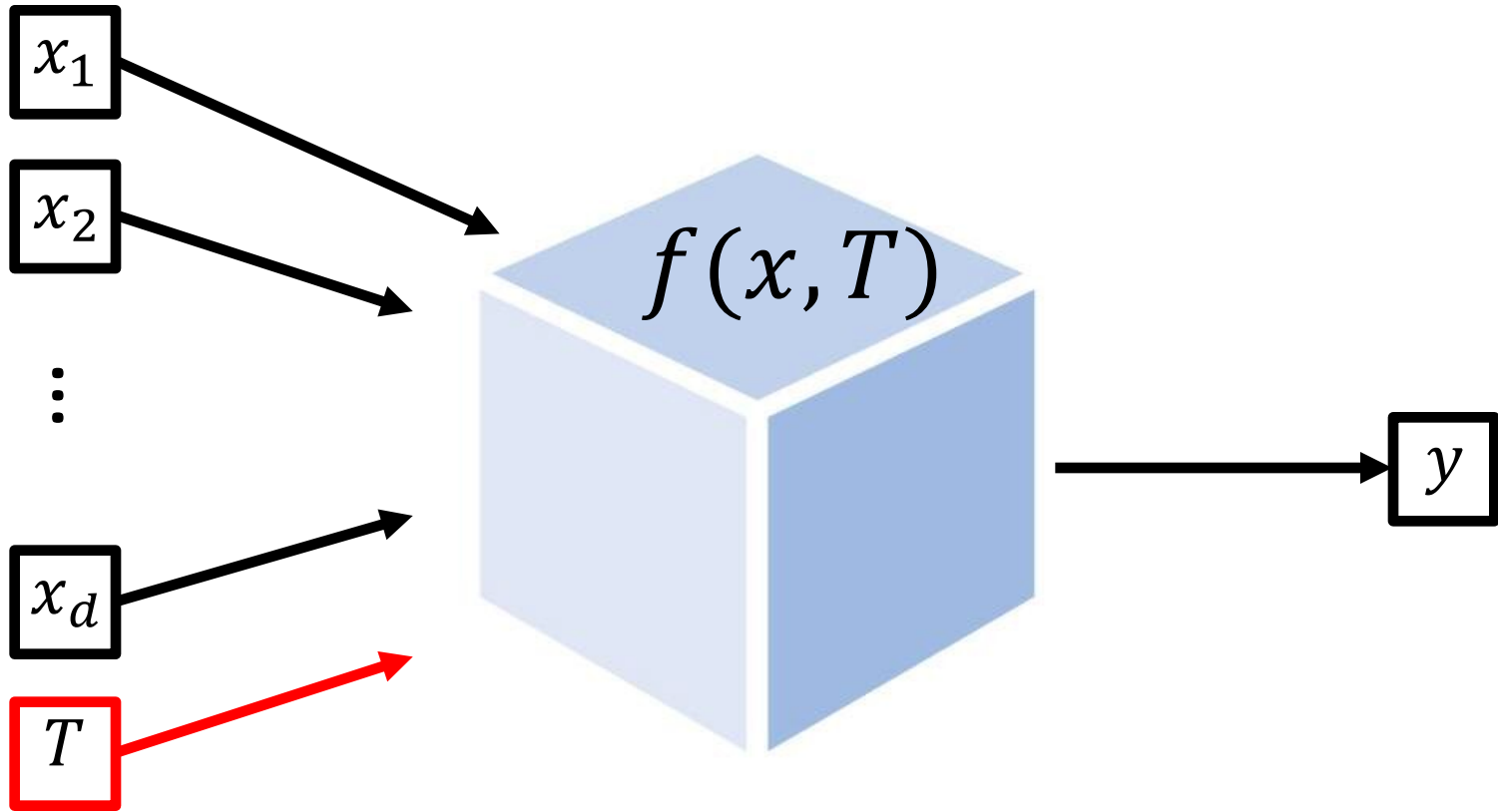
$T$

Regression  
model



Outcome

$y$



Nuisance  
Parameters

$x_1$

$x_2$

$\vdots$

$x_d$

Regression  
model

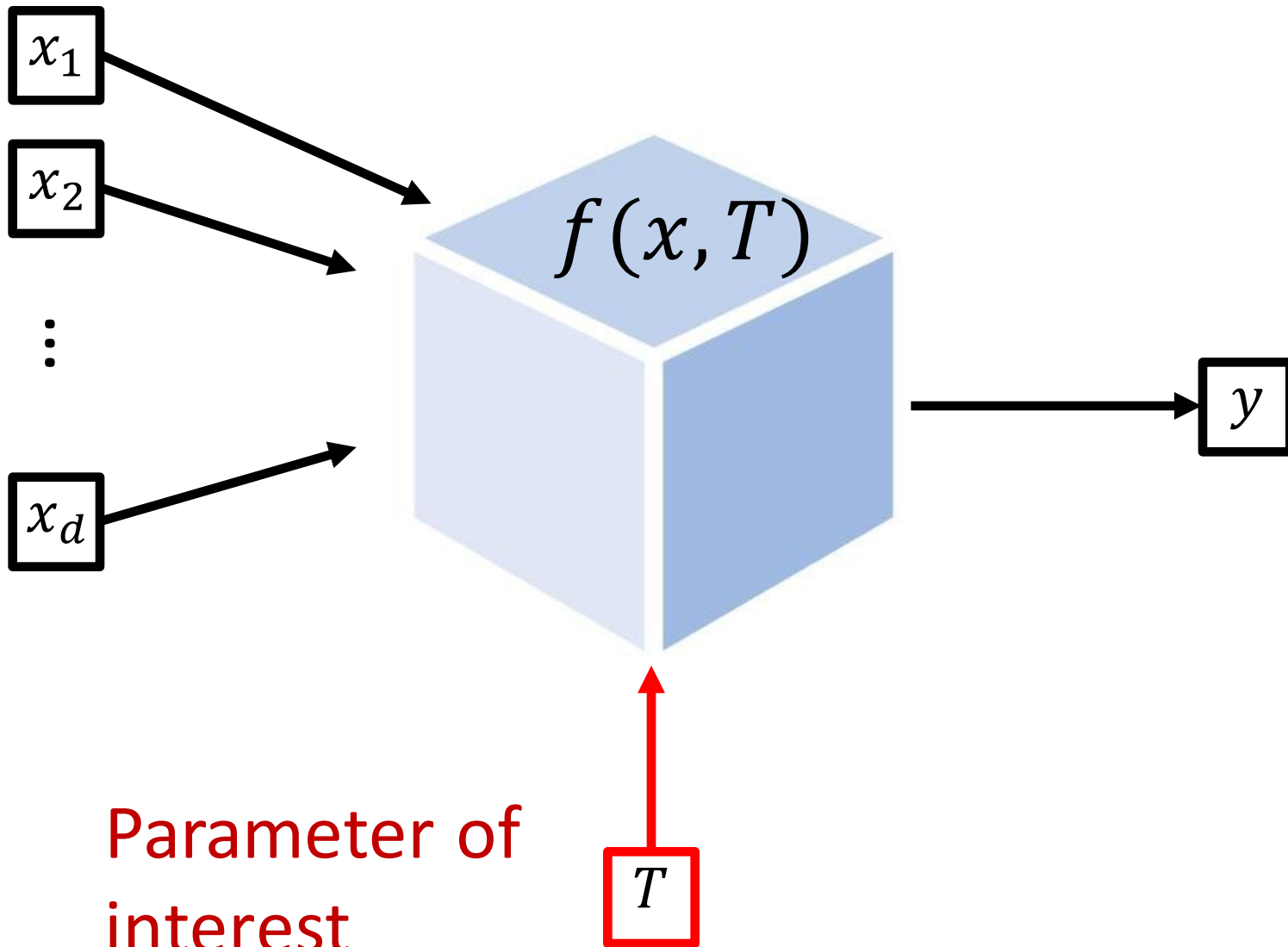
$f(x, T)$

Outcome

$y$

Parameter of  
interest

$T$



# Covariate adjustment (parametric g-formula)

- Explicitly model the relationship between treatment, confounders, and outcome
- Under ignorability, the expected causal effect of  $T$  on  $Y$ :

$$\mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}[Y_1 | T = 1, x] - \mathbb{E}[Y_0 | T = 0, x] \right]$$

- Fit a model  $f(x, t) \approx \mathbb{E}[Y_t | T = t, x]$

$$\widehat{ATE} = \frac{1}{n} \sum_{i=1}^n f(x_i, 1) - f(x_i, 0)$$

# Covariate adjustment (parametric g-formula)

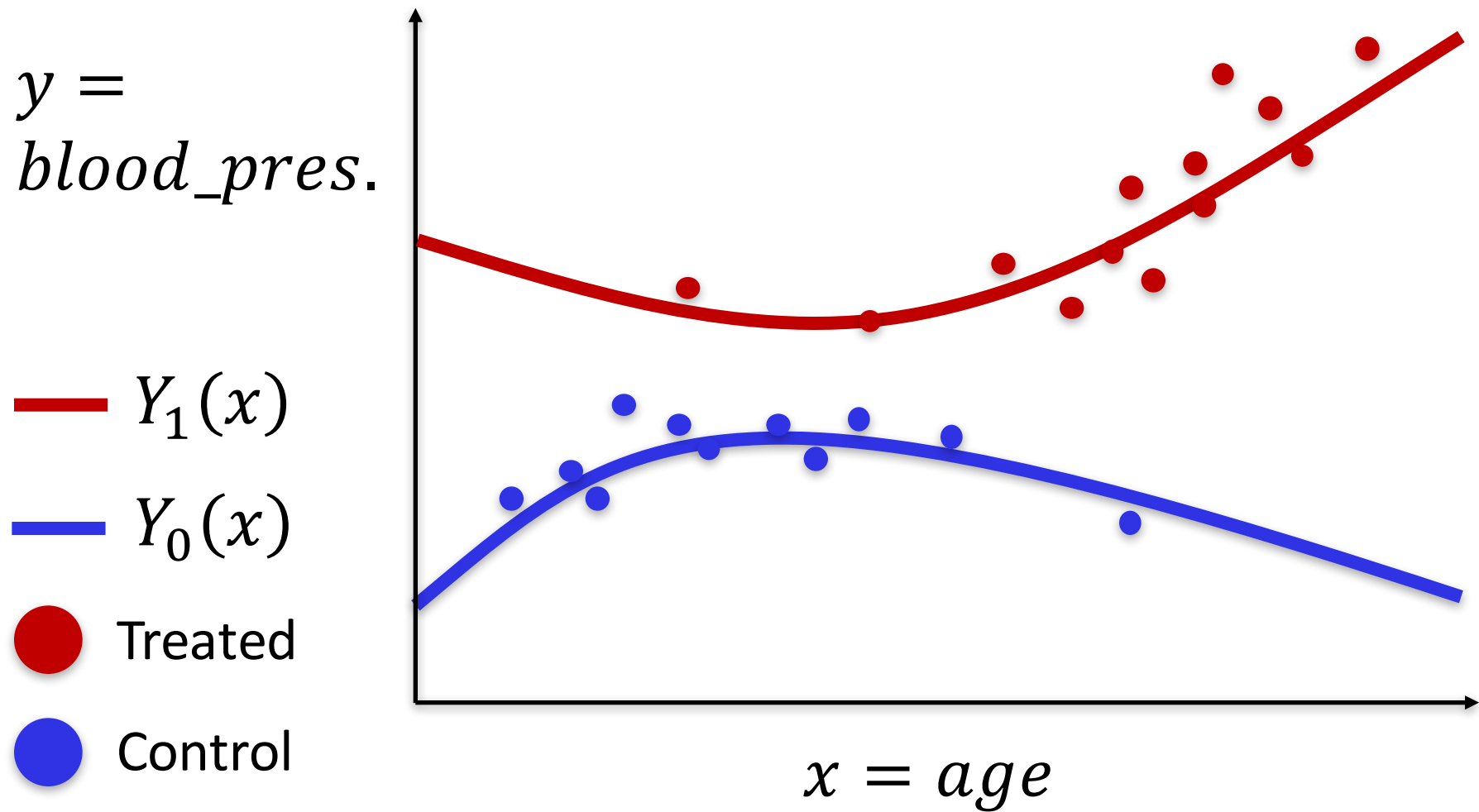
- Explicitly model the relationship between treatment, confounders, and outcome
- Under ignorability, the expected causal effect of  $T$  on  $Y$ :

$$\mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}[Y_1 | T = 1, x] - \mathbb{E}[Y_0 | T = 0, x] \right]$$

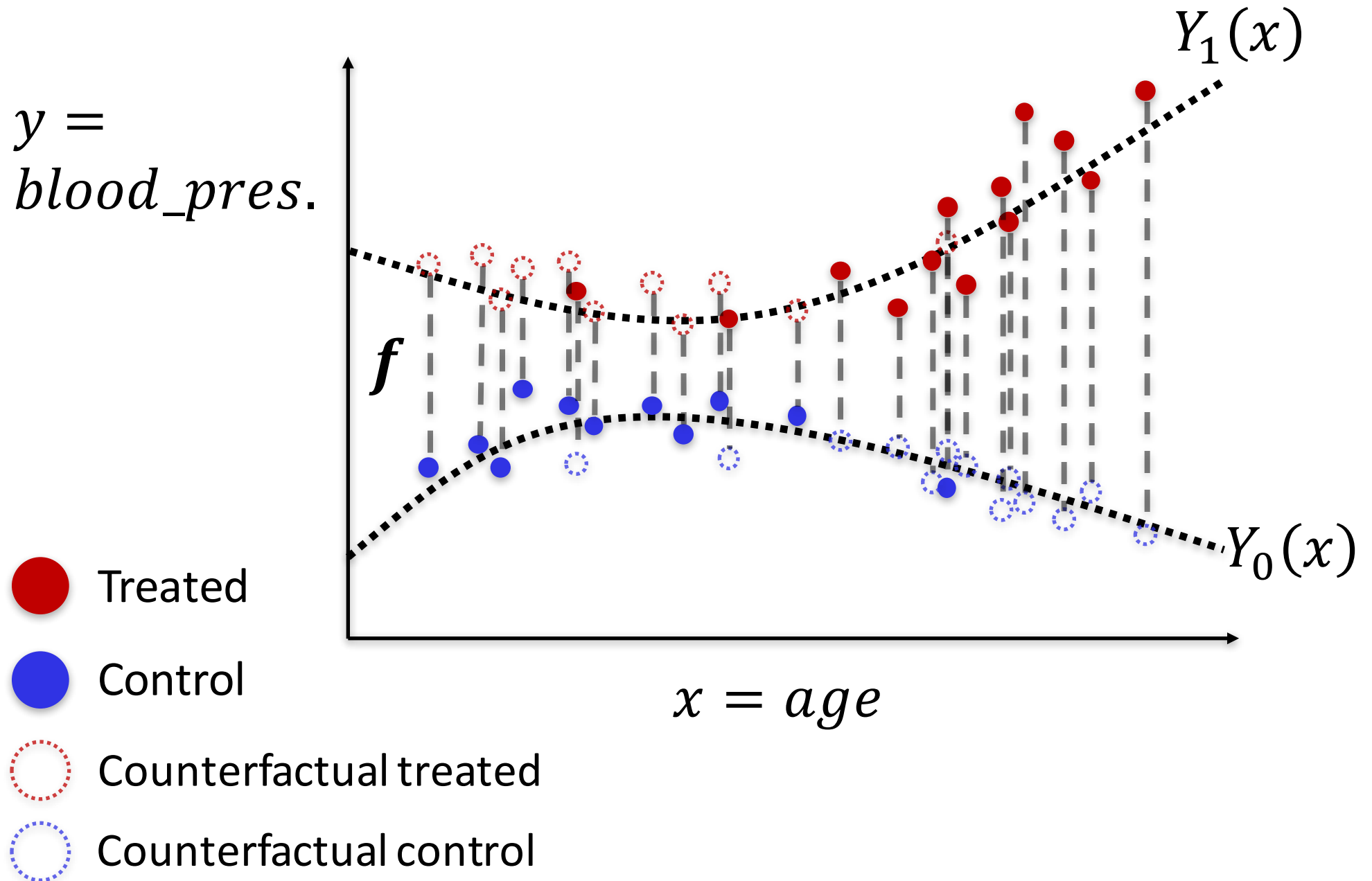
- Fit a model  $f(x, t) \approx \mathbb{E}[Y_t | T = t, x]$

$$\widehat{CATE}(x_i) = f(x_i, 1) - f(x_i, 0)$$

# Covariate adjustment



# Covariate adjustment



# Example of how covariate adjustment fails when there is no overlap

