



Munich Center for Machine Learning

Causal ML for predicting treatment outcomes

**Prof. Stefan Feuerriegel &
Valentyn Melnychuk**

Institute of AI in Management
LMU Munich
<https://www.ai.bwl.lmu.de>



VISION

Promises of Causal ML

Estimating treatment effects for vulnerable groups



Augmenting evidence from RCTs



Finding optimal dosages



ML for treatment effect estimation

Estimating post-approval efficacy, including side effects



Guiding treatment choice when a standard of care is absent



Estimating treatment effects for long-term outcomes



Designing treatment recommendations for rare diseases





Munich Center for Machine Learning

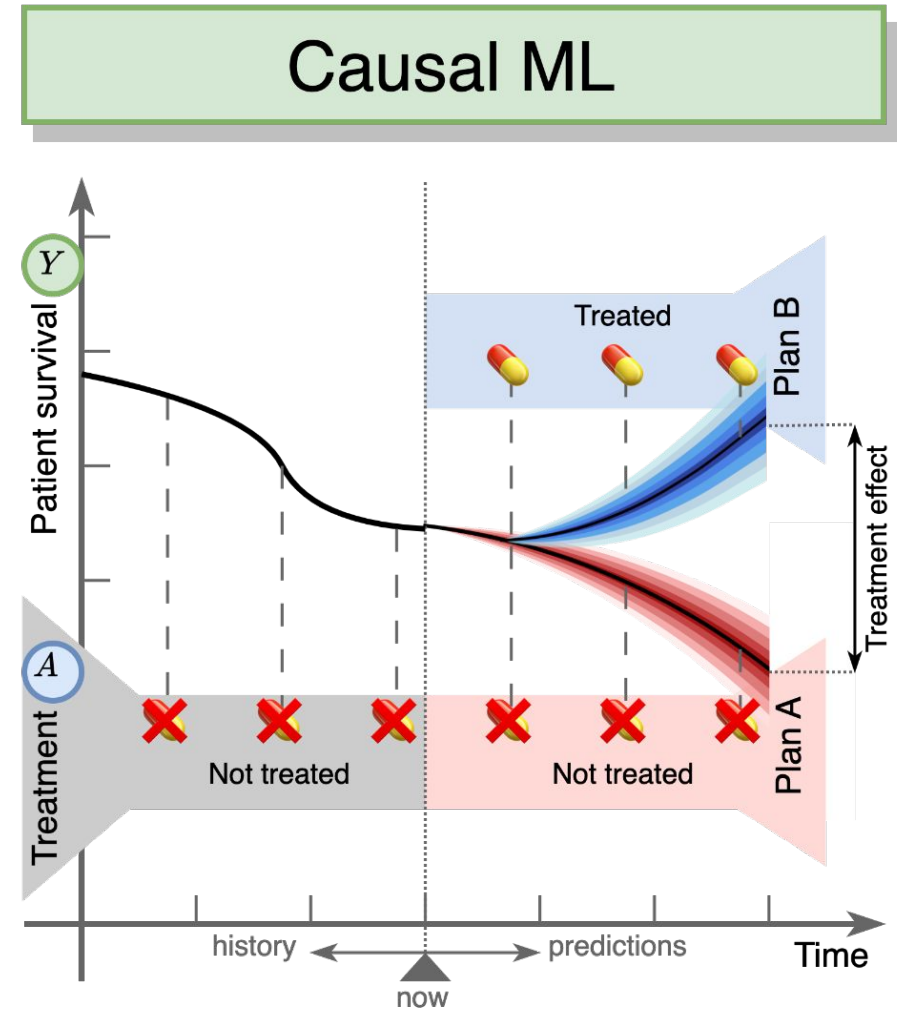
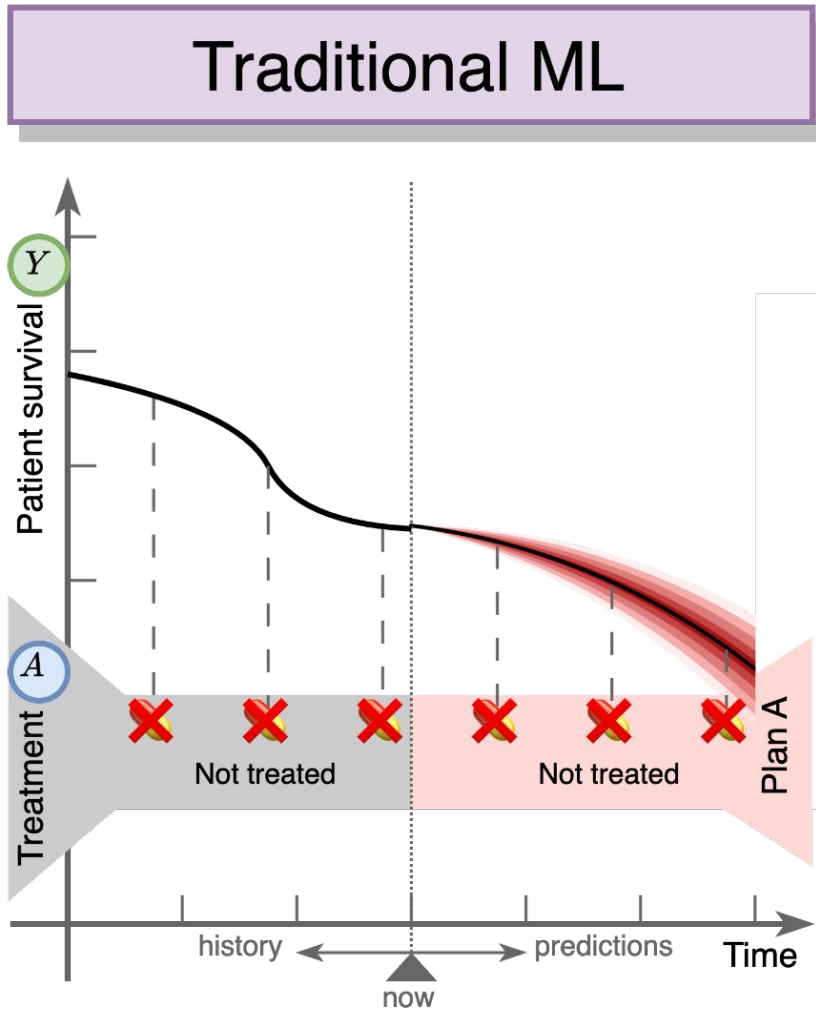
Why do we need Causal ML in medicine?

Reference:

Feuerriegel, S., Frauen, D., Melnychuk, V., Schweisthal, J., Hess, K., Curth, A., Bauer, S., Kilbertus, N., Kohane, I.S. and van der Schaar, M., 2024. **Causal machine learning for predicting treatment outcomes.** *Nature Medicine*, 30(4), pp.958-968.

TERMINOLOGY

Moving from diagnostics to therapeutics: Estimating treatment effects with ML



TERMINOLOGY

Real-world data (RWD) vs. real-world evidence (RWE) to support medicine

The US Food and Drug Administration (FDA) defines ^{1,2,3}:



Real-world data (RWD)

- Data relating to patient health status and the delivery of healthcare
- **Examples:** electronic health records (EHRs), claims and billing activities, disease registries, ...
- Naming: observational data (≠ experimental data)



Real-world evidence (RWE)

- Analysis of RWD regarding usage and effectiveness
- Vision: greater personalization of care
- Disclaimer: should not replace but augment RCTs

1) Real-World Evidence — Where Are We Now? <https://www.nejm.org/doi/full/10.1056/NEJMp2200089>
2) Real-World Evidence — What Is It and What Can It Tell Us? <https://www.nejm.org/doi/full/10.1056/nejmsb1609216>
3) Real-World Evidence and Real-World Data for Evaluating Drug Safety and Effectiveness <https://jamanetwork.com/journals/jama/fullarticle/2697359>

TERMINOLOGY

Real-world data (RWD) vs. real-world evidence (RWE) to support medicine

The US Food and Drug Administration (FDA) defines ^{1,2,3}:



Real-world data (RWD)

- Data relating to patient health status and the delivery of healthcare
- **Examples:** electronic health records (EHRs), claims and billing activities, disease registries, ...
- Naming: observational data (≠ experimental data)



- **Aim:** estimate treatment effectiveness
- **Challenges:** representativeness (selection bias), no proper randomization, ...
- **Custom methodologies:** target trial emulation, **causal machine learning**, ...



Real-world evidence (RWE)

- Analysis of RWD regarding usage and effectiveness
- Vision: greater personalization of care
- Disclaimer: should not replace but augment RCTs

1) Real-World Evidence — Where Are We Now? <https://www.nejm.org/doi/full/10.1056/NEJMp2200089>
2) Real-World Evidence — What Is It and What Can It Tell Us? <https://www.nejm.org/doi/full/10.1056/nejmsb1609216>
3) Real-World Evidence and Real-World Data for Evaluating Drug Safety and Effectiveness <https://jamanetwork.com/journals/jama/fullarticle/2697359>

Application scenarios of RWD

RWD helps to guide decision-making (beyond RCTs):

- 1 ... in the absence of a standard of care**
 - Specific subtypes of diseases with no standard of care yet (e.g., oncology)
 - New or experimental drugs (e.g., orphan drugs, is Biontech vs. Moderna vaccine more effective for subcohort X?)

- 2 ... in complex, high-dimensional decision problems**
 - Complex dosaging problems (e.g., chemotherapy, combi-treatments)

- 3 ... when RCTs are unethical**
 - Vulnerable populations (pregnant women, children, severely ill, etc.) ¹

- 4 ... when a greater personalization is desired**
 - Highly granular subpopulations that cannot be really placed in RCTs (e.g., women, above 60, with comorbidity X, Y & Z or generally specific patient trajectories)
→ maybe a subpopulations responds different for a specific drug, or a second line of treatment is more effective than the first line?
 - Personalization based on genome data (e.g., precision medicine)

EXAMPLE

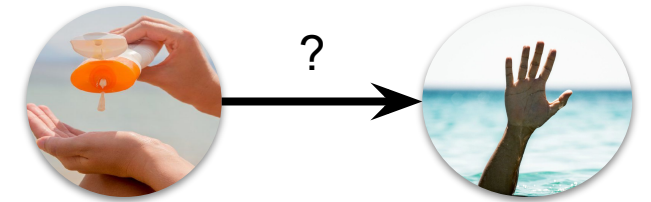
Real-world data (RWD) vs. real-world evidence (RWE) to support medicine

Why is getting a **meaningful** RWE challenging?



Real-world data
(RWD)

- Observational data of
 - sunscreen usage (binary treatment)
 - number of drowning-related deaths (outcome)



-
- **Aim:** effect of sunscreen on the chance of drowning



Real-world evidence
(RWE)

-
- Evidence: The higher the usage of sunscreen -> the more likely is the chance of drowning
 - This is counterintuitive: Is there something we didn't account for?

EXAMPLE

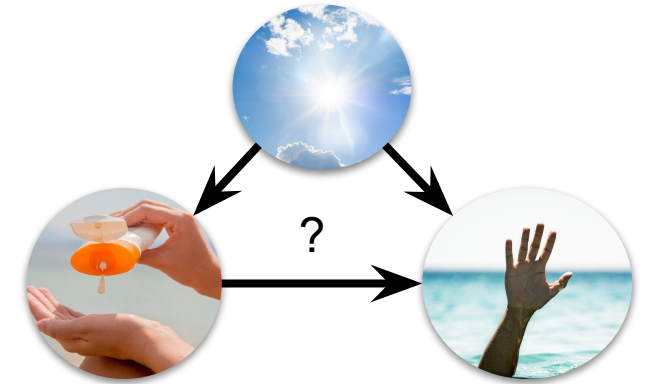
Real-world data (RWD) vs. real-world evidence (RWE) to support medicine

Why is getting a **meaningful** RWE challenging?



Real-world data
(RWD)

- Observational data of
 - sunscreen usage (binary treatment)
 - number of drowning-related deaths (outcome)
 - **intensity of sunlight (covariates)**



- **Aim:** effect of sunscreen on the chance of drowning for **different intensities of sunlight**



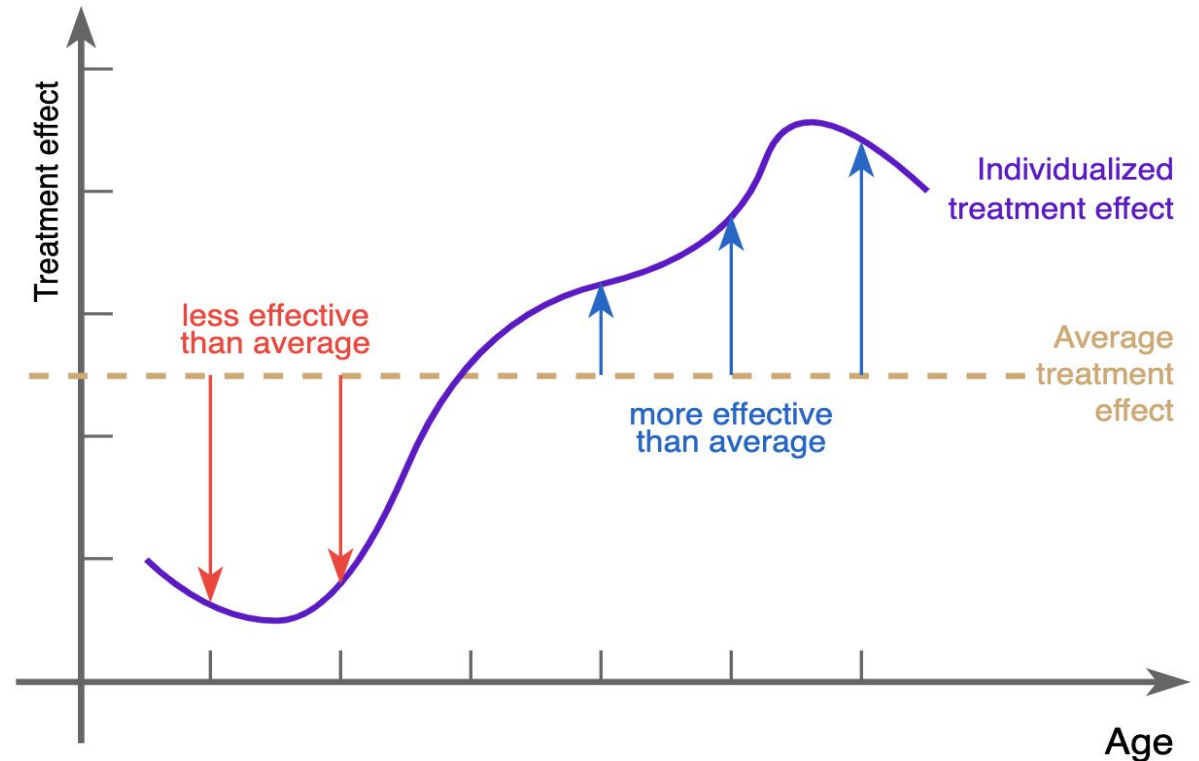
Real-world evidence
(RWE)

- Evidence: no association between sunscreen usage and chance of drowning in each group of sunlight
- Comparing with the previous slide: Intensity of sunlight is a **confounder**

AIM

Understanding heterogeneity in the treatment effect

- Focus is often on **average** treatment effect (ATE)
- ATE is aggregated across the population
- ATE **cannot** tell whether a treatment works for some or not
→ e.g., medication works only for women but not for men, but RCT was done with all patients
- NB: both RCTs and target trial emulation focus on ATEs



To personalize treatment recommendations, we need to understand the **individualized** treatment effect (ITE)



Munich Center for Machine Learning

Short introduction to causal machine learning

Reference:

Feuerriegel, S., Frauen, D., Melnychuk, V., Schweisthal, J., Hess, K., Curth, A., Bauer, S., Kilbertus, N., Kohane, I.S. and van der Schaar, M., 2024. Causal machine learning for predicting treatment outcomes. *Nature Medicine*, 30(4), pp.958-968.

Ladder of causation

Pearl's layers of causation

Level (Symbol)	Typical Activity	Typical Questions	Examples
1. Association $P(y x)$	Seeing	What is? How would seeing X change my belief in Y ?	What does a symptom tell me about a disease? What does a survey tell us about the election results?
2. Intervention $P(y do(x), z)$	Doing Intervening	What if? What if I do X ?	What if I take aspirin, will my headache be cured? What if we ban cigarettes?
3. Counterfactuals $P(y_x x', y')$	Imagining, Retrospection	Why? Was it X that caused Y ? What if I had acted differently?	Was it the aspirin that stopped my headache? Would Kennedy be alive had Oswald not shot him? What if I had not been smoking the past 2 years?



Causal Hierarchy Theorem: statistical inference for a layer requires the information from the same or higher layer. For the inference from lower layer data, we need to make **additional assumptions**.

Ladder of causation

Pearl's layers of causation	Level (Symbol)	Typical Activity	Typical Questions	Examples	Traditional ML
	1. Association $P(y x)$	Seeing	What is? How would seeing X change my belief in Y ?	What does a symptom tell me about a disease? What does a survey tell us about the election results?	
	2. Intervention $P(y do(x), z)$	Doing Intervening	What if? What if I do X ?	What if I take aspirin, will my headache be cured? What if we ban cigarettes?	
	3. Counterfactuals $P(y_x x', y')$	Imagining, Retrospection	Why? Was it X that caused Y ? What if I had acted differently?	Was it the aspirin that stopped my headache? Would Kennedy be alive had Oswald not shot him? What if I had not been smoking the past 2 years?	



Causal Hierarchy Theorem: statistical inference for a layer requires the information from the same or higher layer. For the inference from lower layer data, we need to make **additional assumptions**.

Ladder of causation

Level (Symbol)	Typical Activity	Typical Questions	Examples
1. Association $P(y x)$	Seeing	What is? How would seeing X change my belief in Y ?	What does a symptom tell me about a disease? What does a signal about the election tell me?
2. Intervention $P(y do(x), z)$	Doing Intervening	What if? What if I do X ?	What if I take aspirin, will my headache be cured? What if we ban cigarettes?
3. Counterfactuals $P(y_x x', y')$	Imagining, Retrospection	Why? Was it X that caused Y ? What if I had acted differently?	Was it the aspirin that stopped my headache? Would Kennedy be alive had Oswald not shot him? What if I had not been smok- ing the past 2 years?

Pearl's
layers of
causation

Causal ML



Causal Hierarchy Theorem: statistical inference for a layer requires the information from the same or higher layer. For the inference from lower layer data, we need to make **additional assumptions**.

Estimating the potential outcomes of treatments

Problem
formulation

- Given i.i.d. observational dataset

$$\mathcal{D} = \{x_i, a_i, y_i\}_{i=1}^n \sim \mathbb{P}(X, A, Y)$$

- X covariates
- A (binary) treatments
- Y continuous (factual) outcomes

- We want to identify & estimate treatment outcomes:

- **treatment effects**

$$Y[1] - Y[0]$$

- **potential outcomes**

(separately) $Y[0]$ $Y[1]$

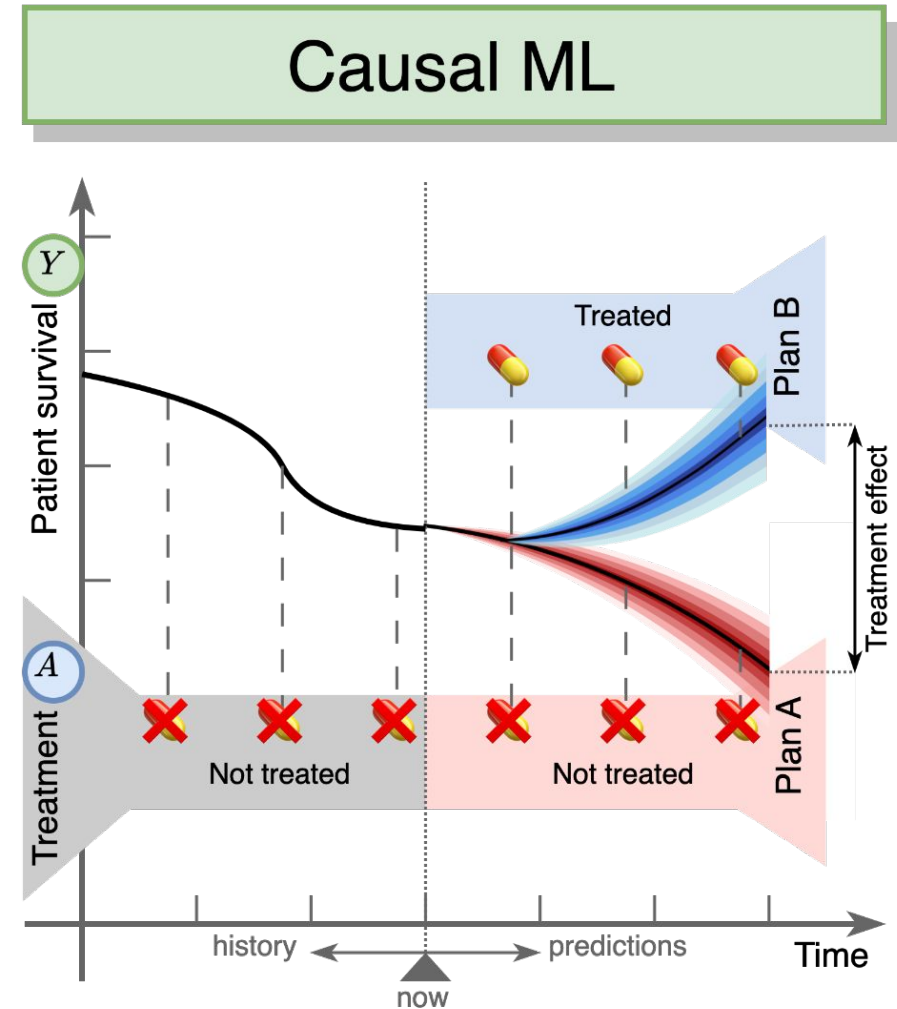
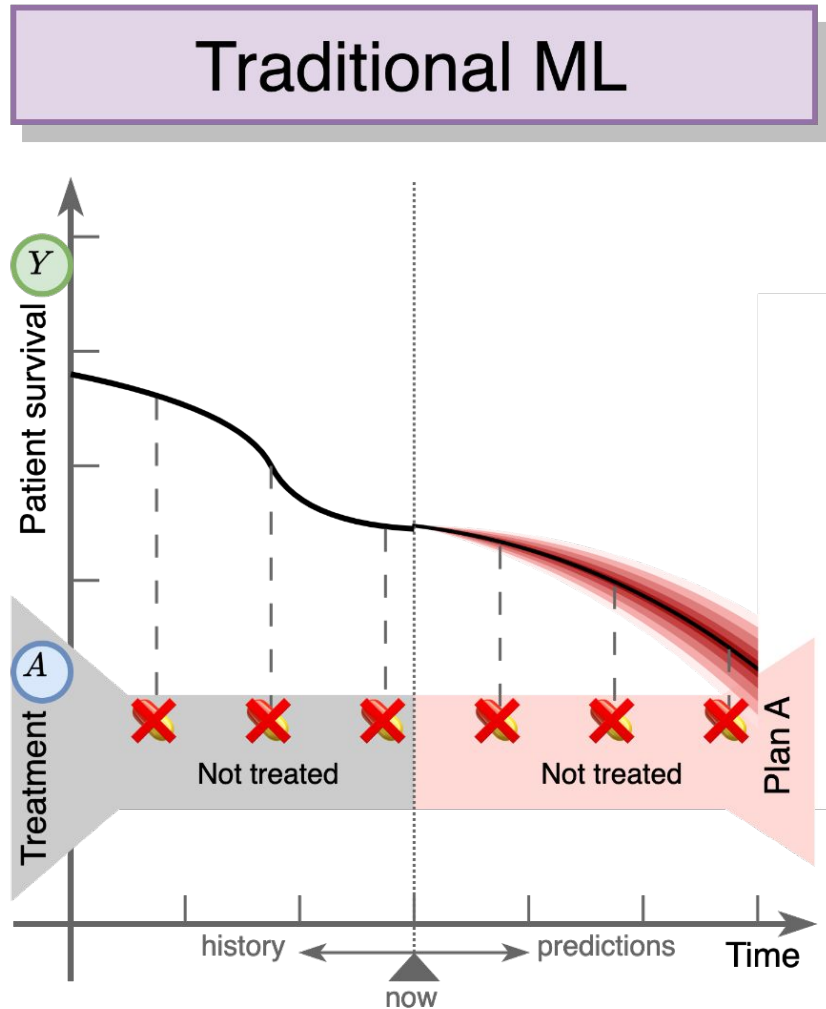
- Fundamental problem:** never observing both potential outcomes!

Patient	Covariates X	Treatment A	Y = Y(0)	Y = Y(1)
		0	-1.0	/
		1	/	2.3
		1	/	0.3
...

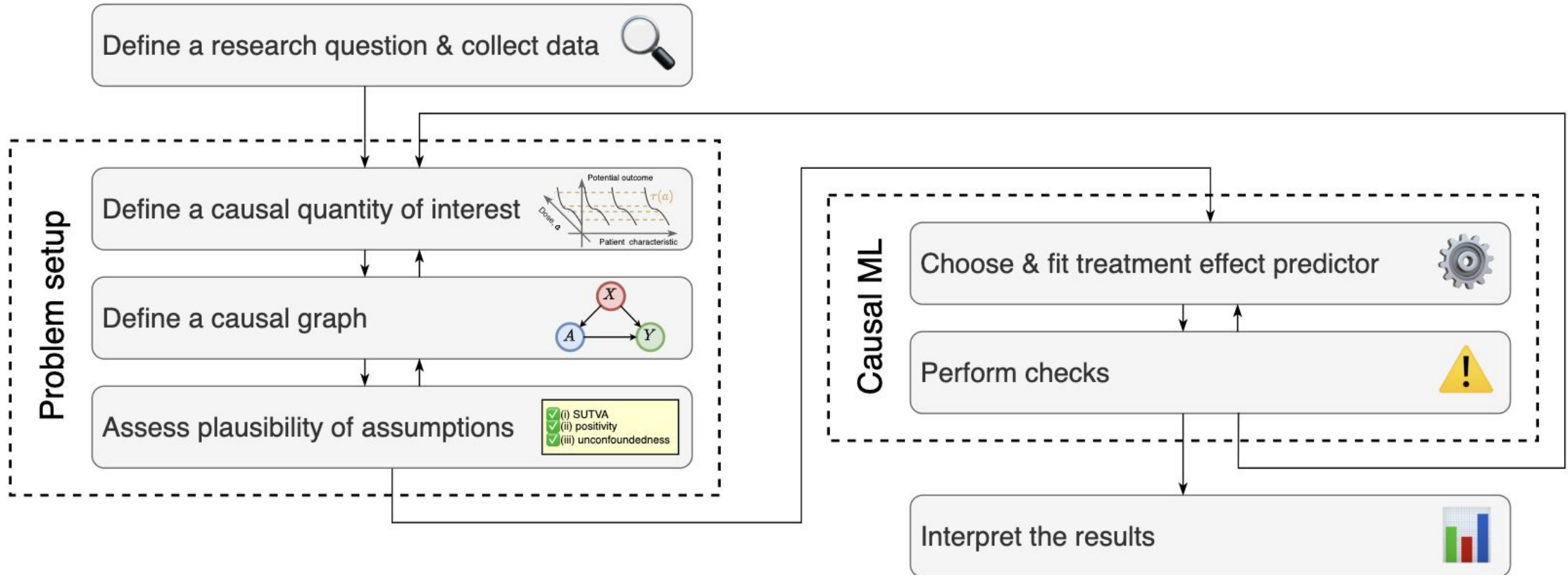
Patient	Covariates X	Potential outcomes Y(0)	Y(1)	Treatment effect Y(1) - Y(0)
		?	?	?
		?	?	?
...

Traditional ML vs. Causal ML

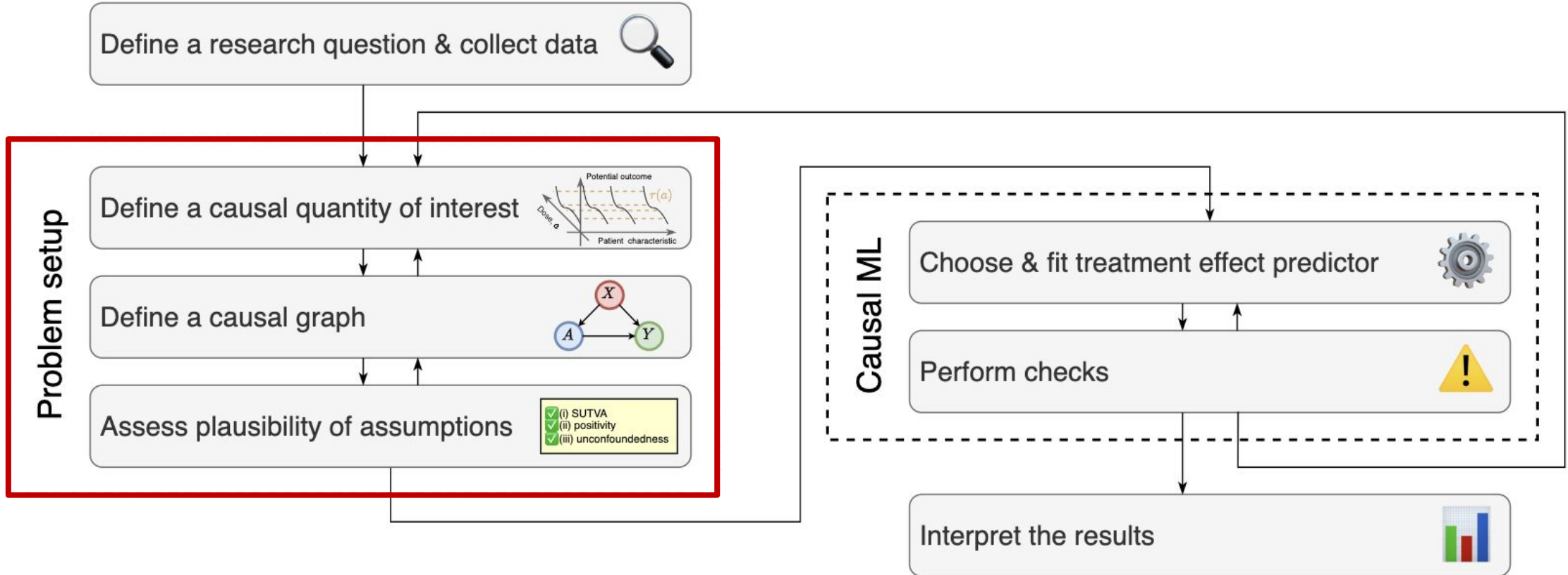
Traditional ML vs. Causal ML



Causal ML Workflow

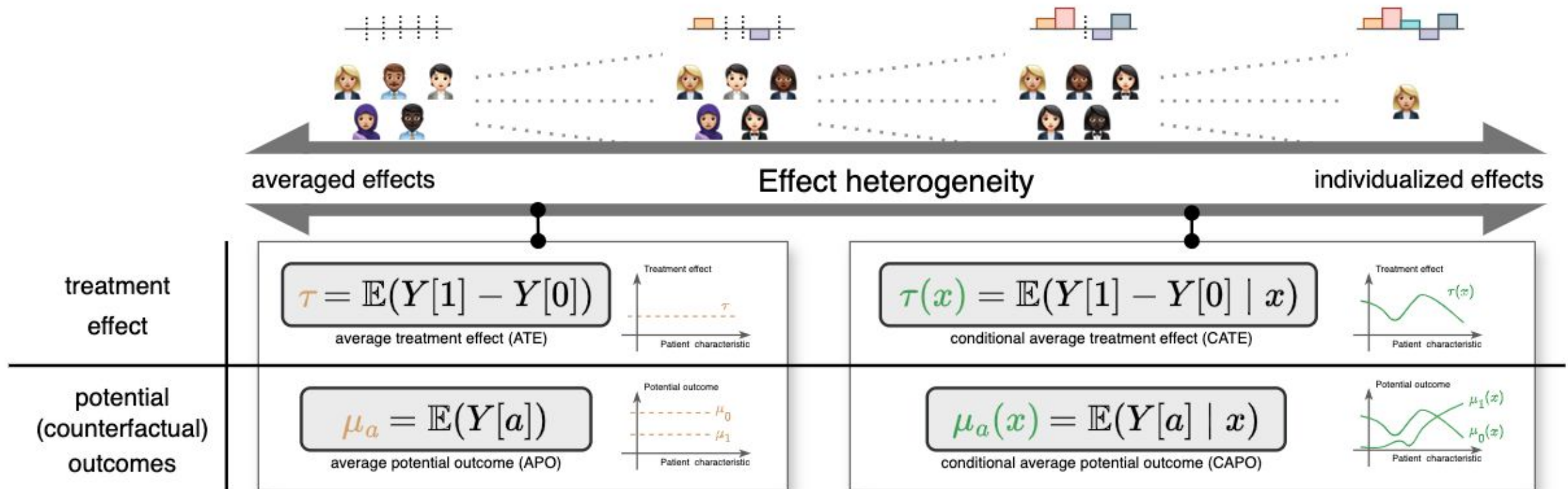


Causal ML Workflow



PROBLEM SETUP

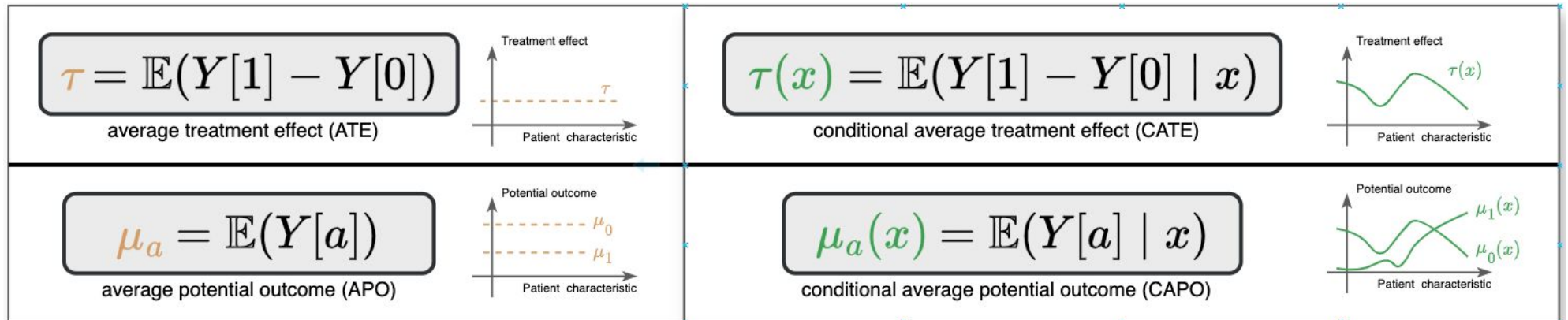
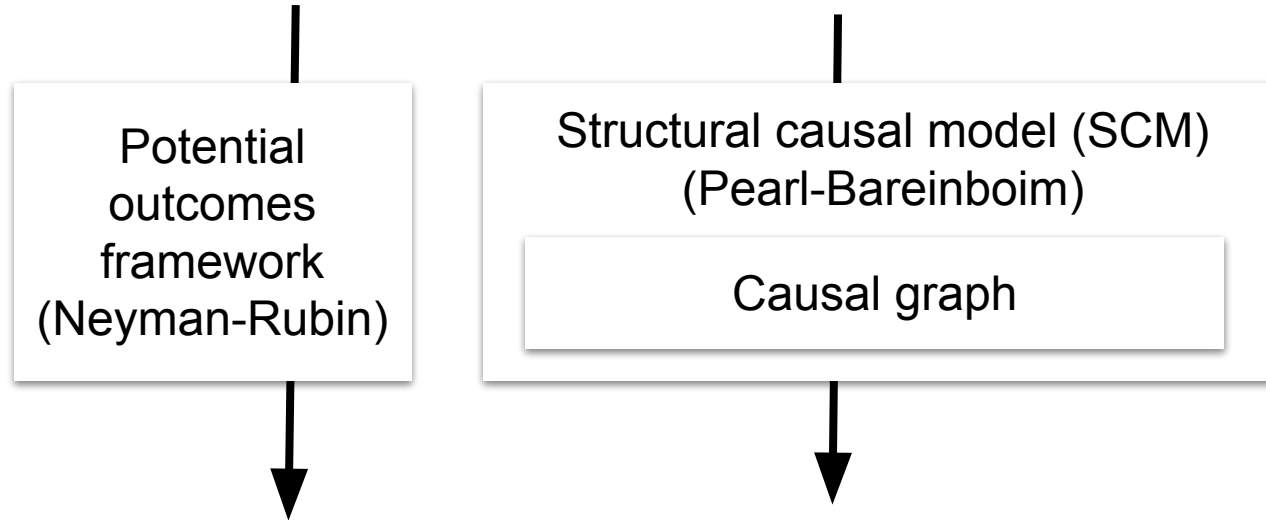
Causal quantities of interest



PROBLEM SETUP

Assumption frameworks

$$\mathcal{D} = \{x_i, a_i, y_i\}_{i=1}^n \sim \mathbb{P}(X, A, Y)$$



PROBLEM SETUP

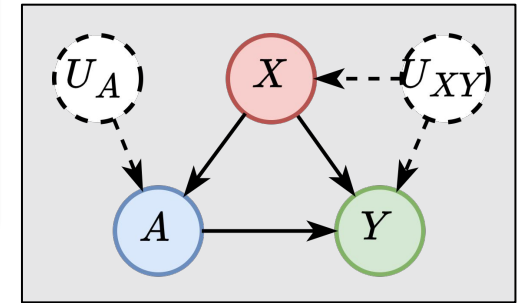
Assumption frameworks: SCMs and causal graphs

$$\mathcal{D} = \{x_i, a_i, y_i\}_{i=1}^n \sim \mathbb{P}(X, A, Y)$$

Potential outcomes framework (Neyman-Rubin)

Structural causal model (SCM) (Pearl-Bareinboim)
Causal graph

Assumptions stem from structural knowledge



$\tau = \mathbb{E}(Y[1] - Y[0])$ <p>average treatment effect (ATE)</p>	$\tau(x) = \mathbb{E}(Y[1] - Y[0] x)$ <p>conditional average treatment effect (CATE)</p>
$\mu_a = \mathbb{E}(Y[a])$ <p>average potential outcome (APO)</p>	$\mu_a(x) = \mathbb{E}(Y[a] x)$ <p>conditional average potential outcome (CAPO)</p>

PROBLEM SETUP

Assumption frameworks: Potential outcomes framework

$$\mathcal{D} = \{x_i, a_i, y_i\}_{i=1}^n \sim \mathbb{P}(X, A, Y)$$

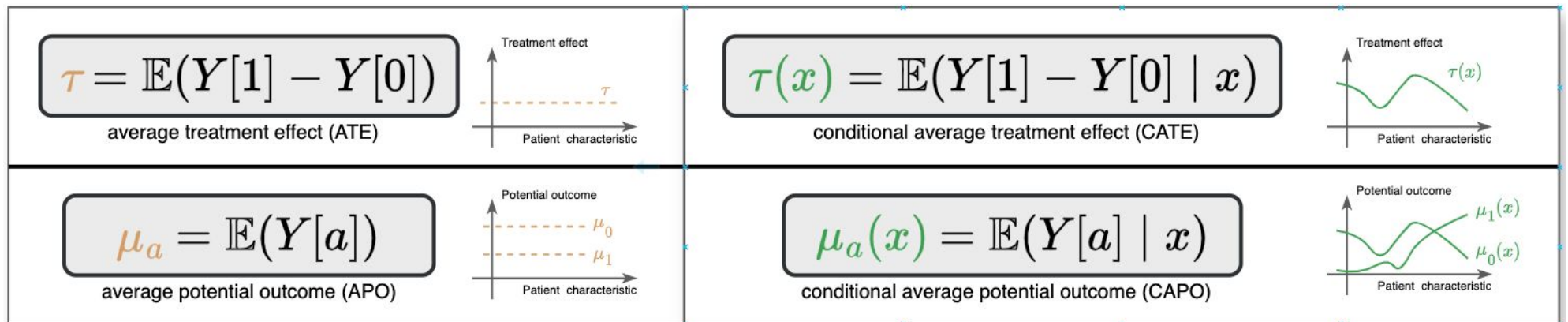
More general

- (i) Consistency
- (ii) Positivity (Overlap)
- (iii) Exchangeability (Ignorability)

Potential outcomes framework (Neyman-Rubin)

Structural causal model (SCM) (Pearl-Bareinboim)

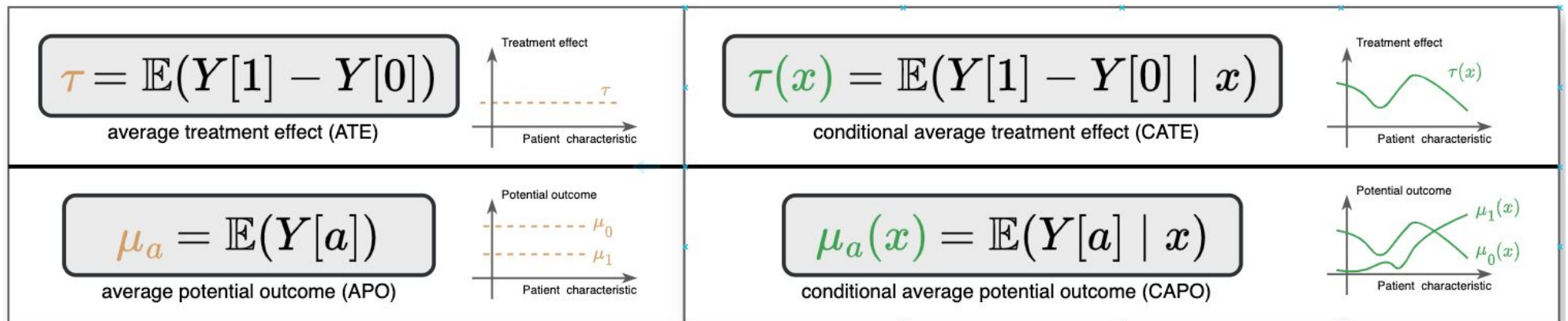
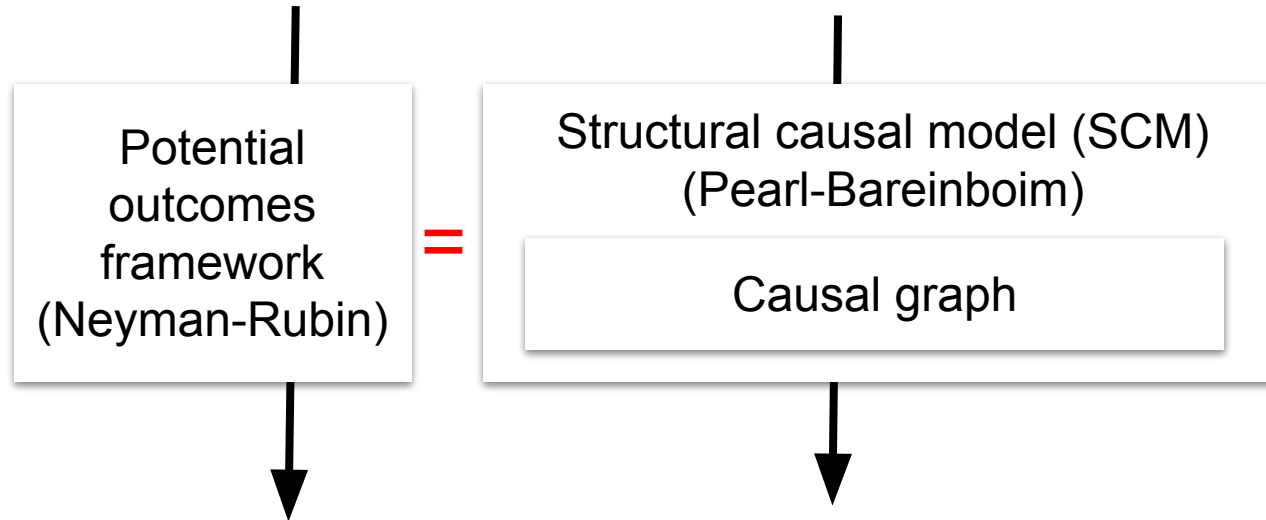
Causal graph



PROBLEM SETUP

Assumption frameworks

$$\mathcal{D} = \{x_i, a_i, y_i\}_{i=1}^n \sim \mathbb{P}(X, A, Y)$$



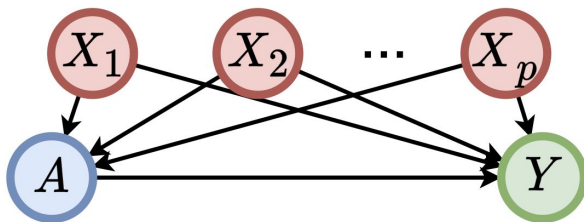
PROBLEM SETUP

Example of a case study

Aim: estimate heterogeneous treatment effect of development aid on SDG outcomes

- **Treatment A :** development aid earmarked to end the HIV/AIDS epidemic
- **Outcome Y :** relative reduction in HIV infection rate
- **Covariates X :** control for differences in country characteristics

Causal graph



Causal quantity of interest

$$\mu_a(x) = \mathbb{E}(Y[a] \mid x)$$

conditional average potential outcome (CAPO)

Assumptions

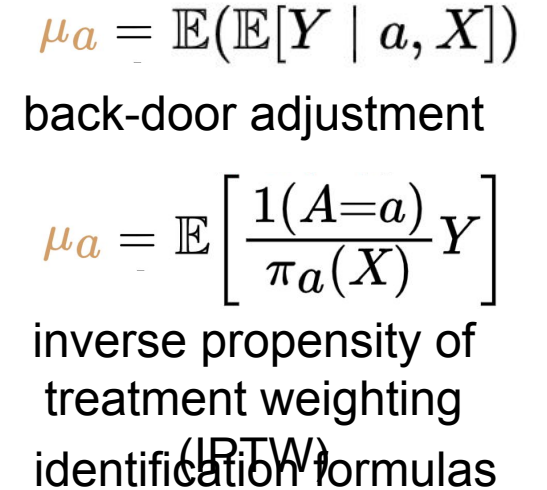
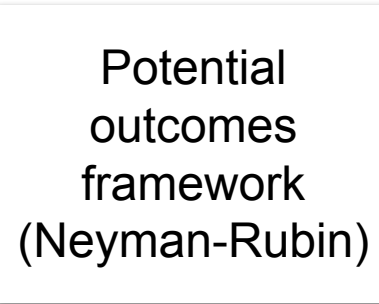
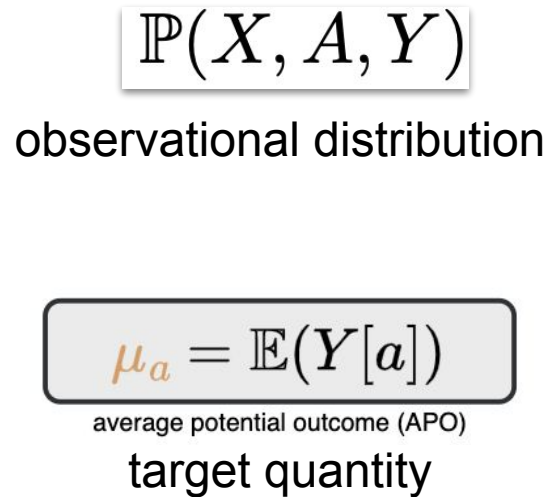
Consistency: $Y = Y(a)$ if $A = a$

Positivity: $0 < p(A = a \mid X = x) < 1, \forall a \in \mathcal{A}$

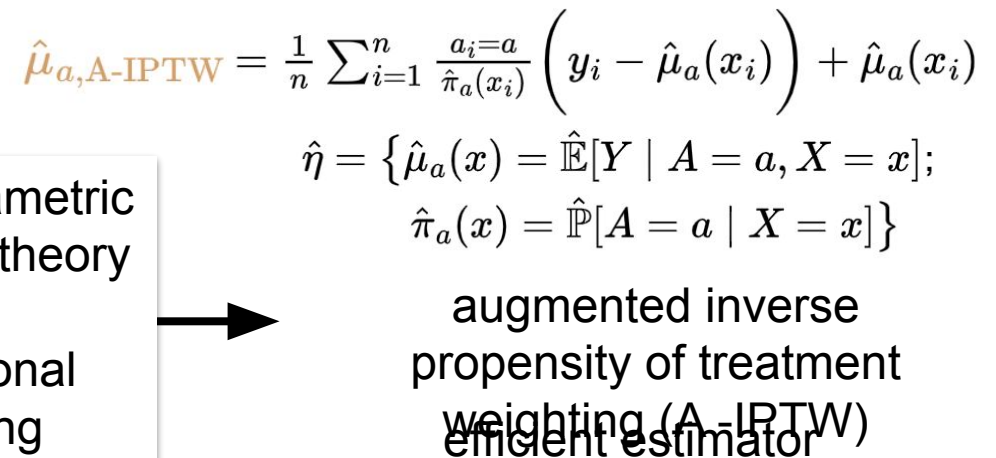
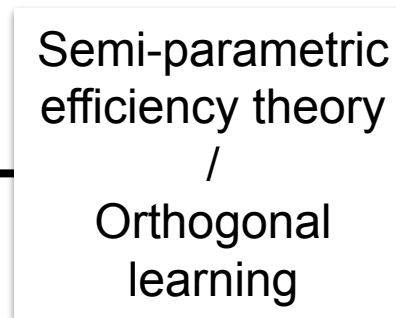
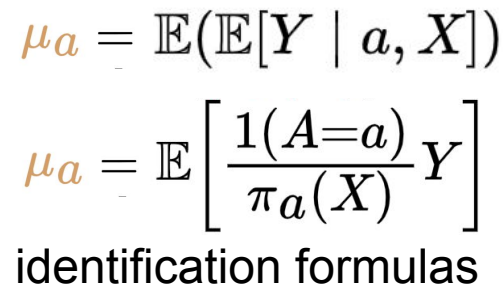
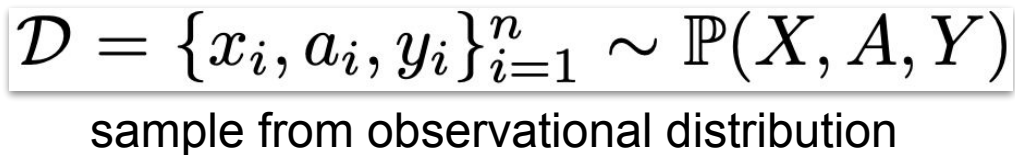
Ignorability: $Y(a) \perp\!\!\!\perp A \mid X = x, \forall a \in \mathcal{A}$

Primer: Identification vs. Estimation

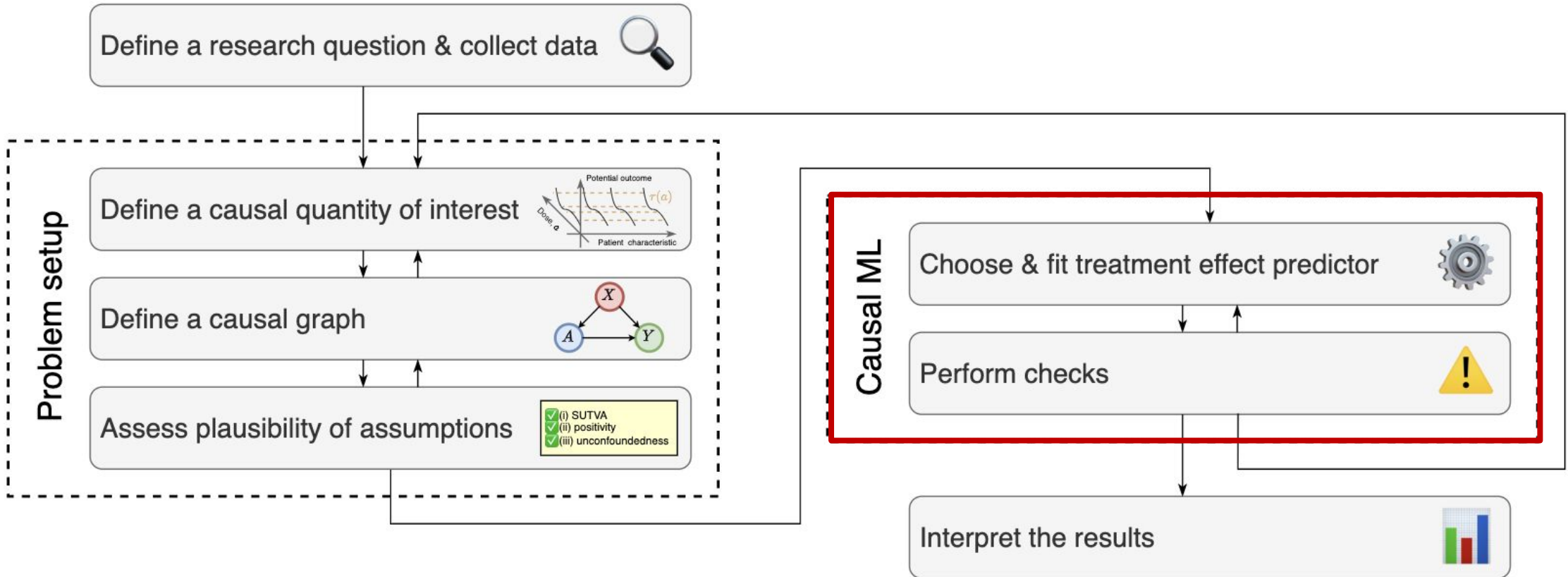
**Identification
(infinite data)**



**Estimation
(finite data)**



Causal ML Workflow



Challenges and open questions fitting an ML model

Challenges

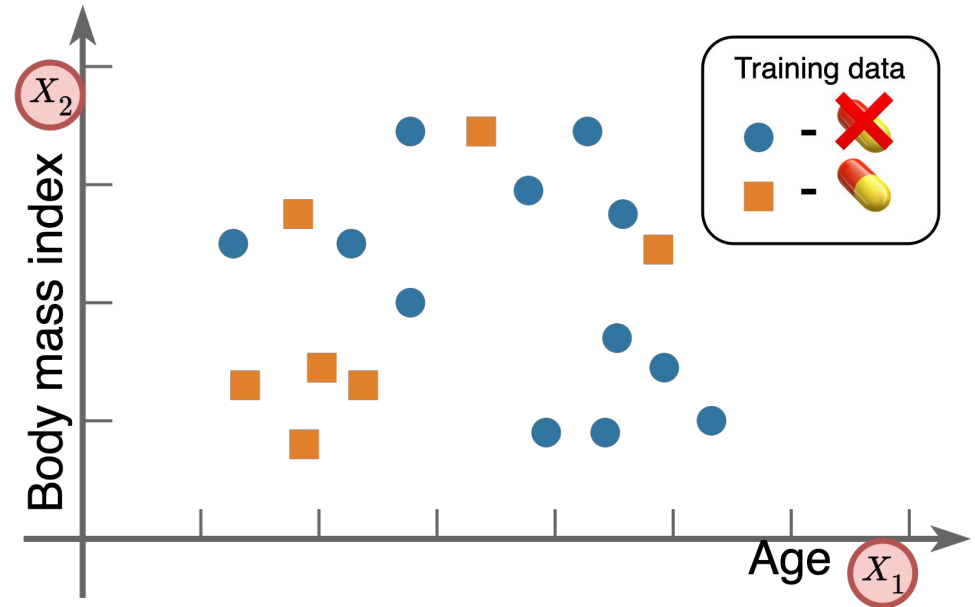
$$\mu_a(x) = \mathbb{E}(Y[a] | x)$$

conditional average potential outcome (CAPO)

$$\tau(x) = \mathbb{E}(Y[1] - Y[0] | x)$$

conditional average treatment effect (CATE)

Open problems



Challenges and open questions fitting an ML model

Challenges

$$\mu_a(x) = \mathbb{E}(Y[a] | x)$$

conditional average potential outcome (CAPO)

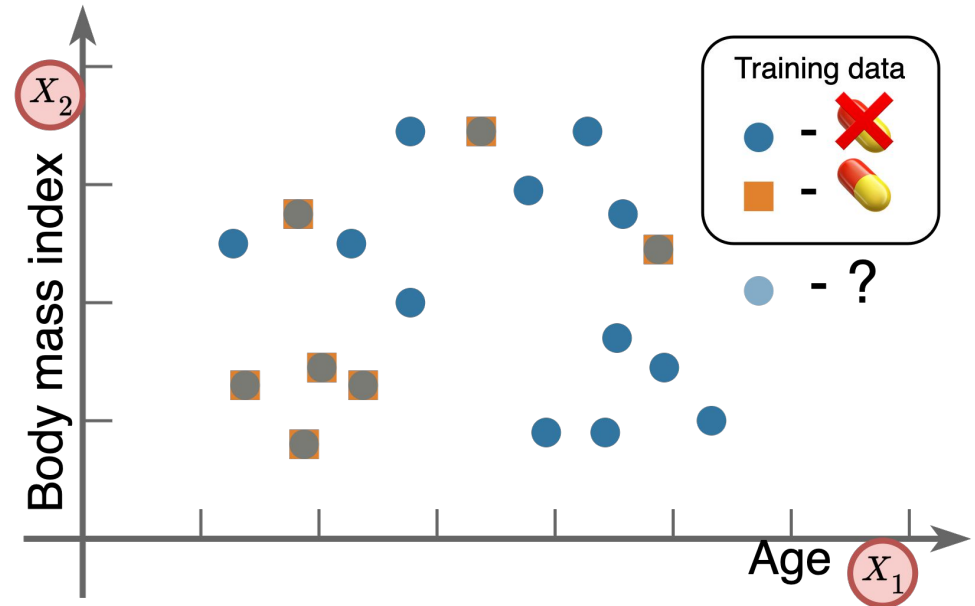
- **Selection bias:** parts of the population rarely gets treated

$$\tau(x) = \mathbb{E}(Y[1] - Y[0] | x)$$

conditional average treatment effect (CATE)

- **Selection bias:** parts of the population rarely gets treated

Open problems



Challenges and open questions fitting an ML model

Challenges

$$\mu_a(x) = \mathbb{E}(Y[a] | x)$$

conditional average potential outcome (CAPO)

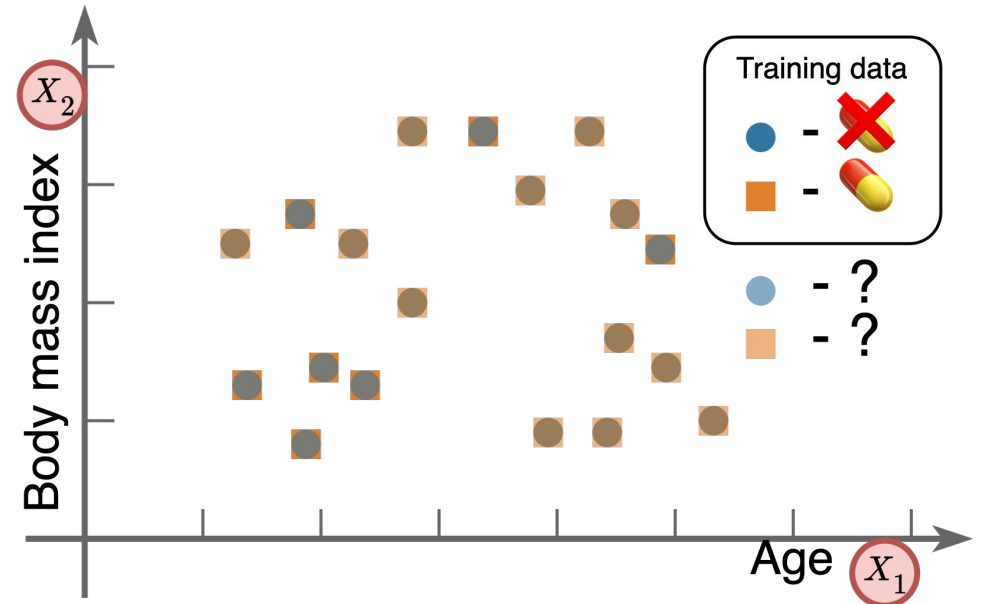
- **Selection bias:** parts of the population rarely gets treated

$$\tau(x) = \mathbb{E}(Y[1] - Y[0] | x)$$

conditional average treatment effect (CATE)

- **Selection bias:** parts of the population rarely gets treated
- **Fundamental problem:** never observing a difference of potential outcomes

Open problems



Challenges and open questions fitting an ML model

Challenges

$$\mu_a(x) = \mathbb{E}(Y[a] | x)$$

conditional average potential outcome (CAPO)

- **Selection bias:** parts of the population rarely gets treated

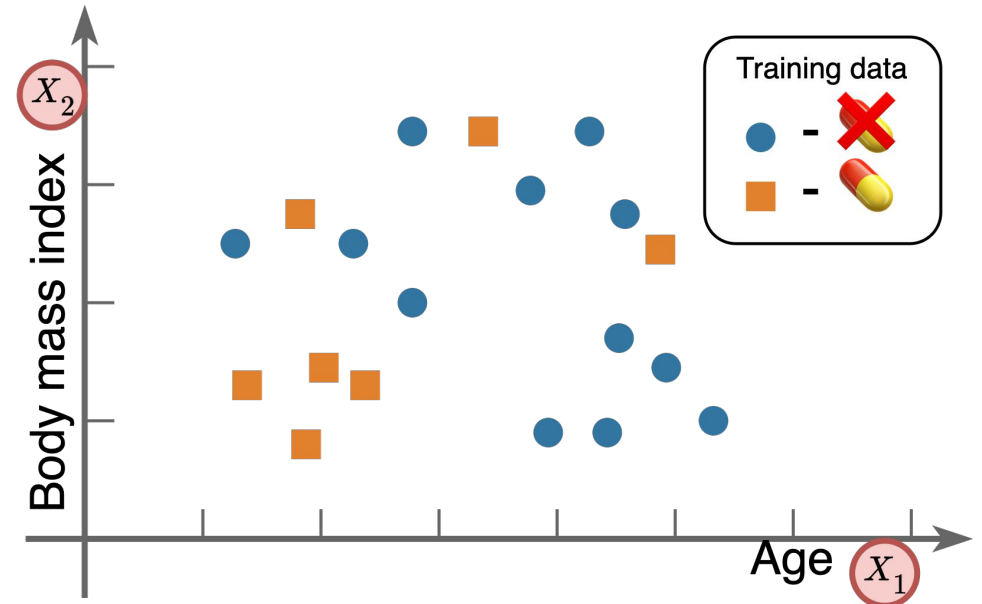
$$\tau(x) = \mathbb{E}(Y[1] - Y[0] | x)$$

conditional average treatment effect (CATE)

- **Selection bias:** parts of the population rarely gets treated
- **Fundamental problem:** never observing a difference of potential outcomes

Open problems

- How to effectively address selection bias?
- How to incorporate inductive biases, e.g., regularize CAPO / CATE models?



CAUSAL ML Methods

Meta-learners

- Meta-learners (Kunzel 2019) are model-agnostic methods for CATE estimation
- Can be used for treatment effect estimation in combination with an arbitrary ML model of choice (e.g., a decision tree, a neural network)

Model-based learners

- Model-specific methods make adjustments to existing ML models to address statistical challenges arising in treatment effect estimation
- Prominent **examples** are the causal tree (Athey 2016) and the causal forest (Wager 2018, Athey 2019)
- Others adapt representation learning to leverage neural networks (Shalit 2017, Shi 2019)

1. Kunzel, Sören R., et al. "Metalearners for estimating heterogeneous treatment effects using machine learning." Proceedings of the national academy of sciences 116.10 (2019): 4156-4165.
2. Athey, Susan, and Guido Imbens. "Recursive partitioning for heterogeneous causal effects." Proceedings of the National Academy of Sciences 113.27 (2016): 7353-7360.
3. Athey, Susan, and Stefan Wager. "Estimating treatment effects with causal forests: An application." Observational studies 5.2 (2019): 37-51.
4. Shalit, Uri, Fredrik D. Johansson, and David Sontag. "Estimating individual treatment effect: generalization bounds and algorithms." International conference on machine learning. PMLR, 2017.
5. Shi, Claudia, David Blei, and Victor Veitch. "Adapting neural networks for the estimation of treatment effects." Advances in neural information processing systems 32 (2019).

CAUSAL ML Methods

Meta-learners

One-stage learners

- “Plug-in learners”: fit a **single** regression model with a treatment as an input or **two** regression models for each treated and control sub-groups
- Examples: S-learner and T-learner

Two-stage learners

- Two-stages of learning: derive and estimate pseudo-outcomes as surrogates, which has the same expected value as the CATE
- Examples: DR-learner and R-learner

Model-based learners

- Model-specific methods make adjustments to existing ML models to address statistical challenges arising in treatment effect estimation
- Prominent **examples** are the causal tree (Athey 2016) and the causal forest (Wager 2018, Athey 2019)
- Others adapt representation learning to leverage neural networks (Shalit 2017, Shi 2019)

1. Künzel, Sören R., et al. "Metalearners for estimating heterogeneous treatment effects using machine learning." Proceedings of the national academy of sciences 116.10 (2019): 4156-4165.
2. Athey, Susan, and Guido Imbens. "Recursive partitioning for heterogeneous causal effects." Proceedings of the National Academy of Sciences 113.27 (2016): 7353-7360.
3. Athey, Susan, and Stefan Wager. "Estimating treatment effects with causal forests: An application." Observational studies 5.2 (2019): 37-51.
4. Shalit, Uri, Fredrik D. Johansson, and David Sontag. "Estimating individual treatment effect: generalization bounds and algorithms." International conference on machine learning. PMLR, 2017.
5. Shi, Claudia, David Blei, and Victor Veitch. "Adapting neural networks for the estimation of treatment effects." Advances in neural information processing systems 32 (2019).

CAUSAL ML Methods

Meta-learners

One-stage learners

- “Plug-in learners”: fit a **single** regression model with a treatment as an input or **two** regression models for each treated and control sub-groups
- Examples: S-learner and T-learner

Two-stage learners

- Two-stages of learning: derive and estimate pseudo-outcomes as surrogates, which has the same expected value as the CATE
- Examples: DR-learner and R-learner

Model-based learners

- Model-specific methods make adjustments to existing ML models to address statistical challenges arising in treatment effect estimation
- Prominent **examples** are the causal tree (Athey 2016) and the causal forest (Wager 2018, Athey 2019)
- Others adapt representation learning to leverage neural networks (Shalit 2017, Shi 2019)

1. Künzel, Sören R., et al. "Metalearners for estimating heterogeneous treatment effects using machine learning." Proceedings of the national academy of sciences 116.10 (2019): 4156-4165.
2. Athey, Susan, and Guido Imbens. "Recursive partitioning for heterogeneous causal effects." Proceedings of the National Academy of Sciences 113.27 (2016): 7353-7360.
3. Athey, Susan, and Stefan Wager. "Estimating treatment effects with causal forests: An application." Observational studies 5.2 (2019): 37-51.
4. Shalit, Uri, Fredrik D. Johansson, and David Sontag. "Estimating individual treatment effect: generalization bounds and algorithms." International conference on machine learning. PMLR, 2017.
5. Shi, Claudia, David Blei, and Victor Veitch. "Adapting neural networks for the estimation of treatment effects." Advances in neural information processing systems 32 (2019).

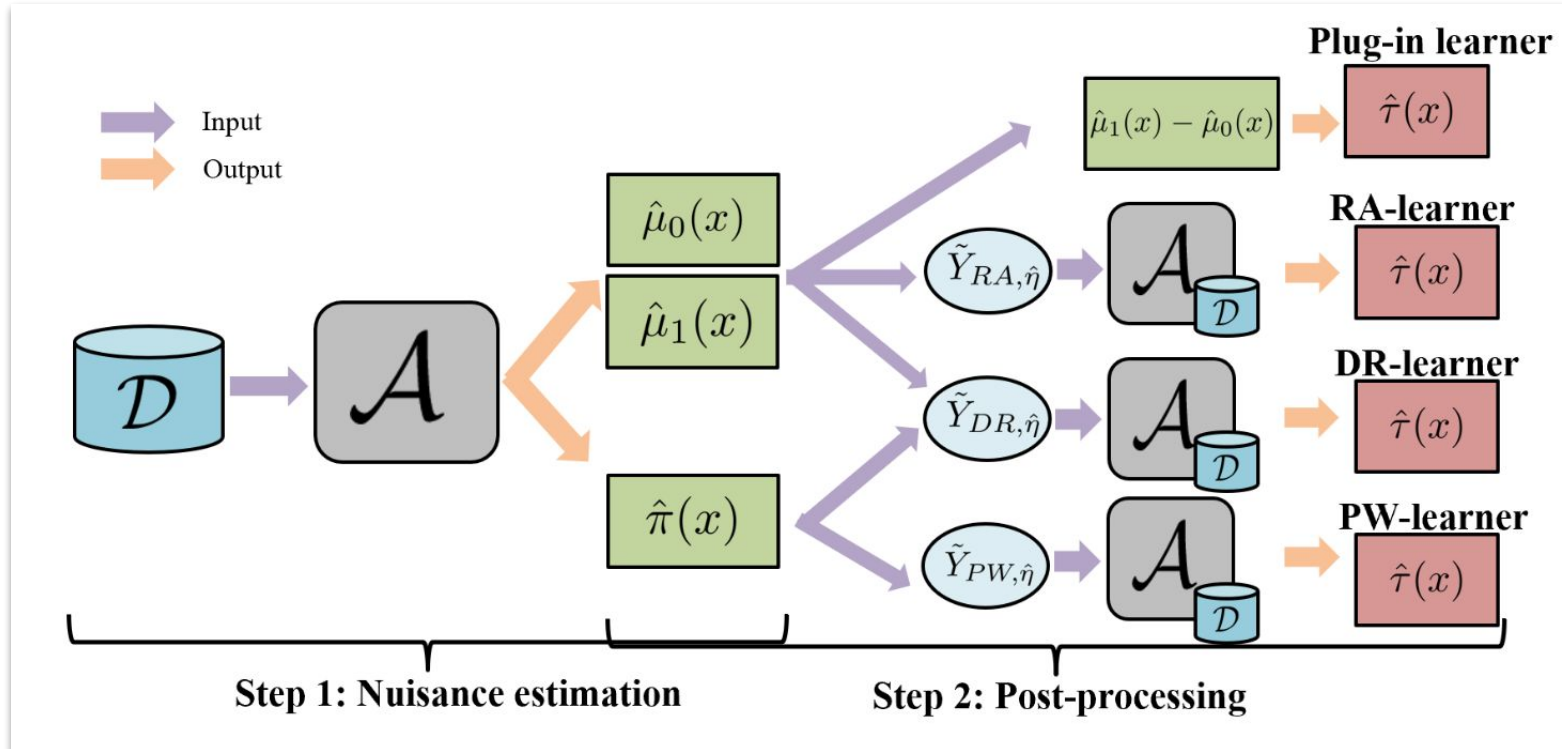
One-stage and two-stage meta-learners

Example: meta-learners for CATE

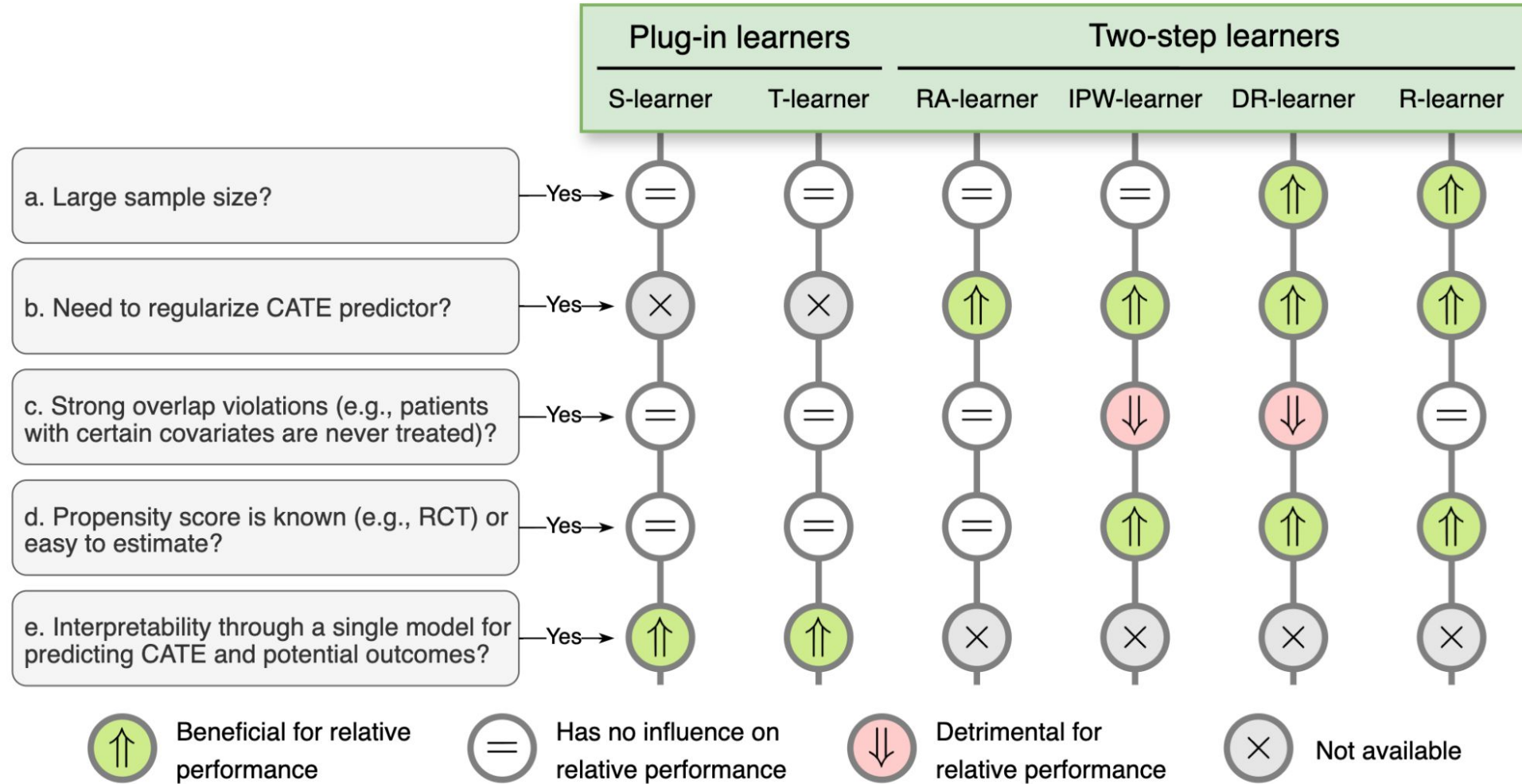
$$\tau(x) = \mathbb{E}(Y[1] - Y[0] \mid x)$$

conditional average treatment effect (CATE)

Method: Using any ML model to fit relevant parts of the observed distribution, namely, **nuisance functions**. Then, we can use the nuisance functions estimators for the final CATE model.



Comparison of meta-learners



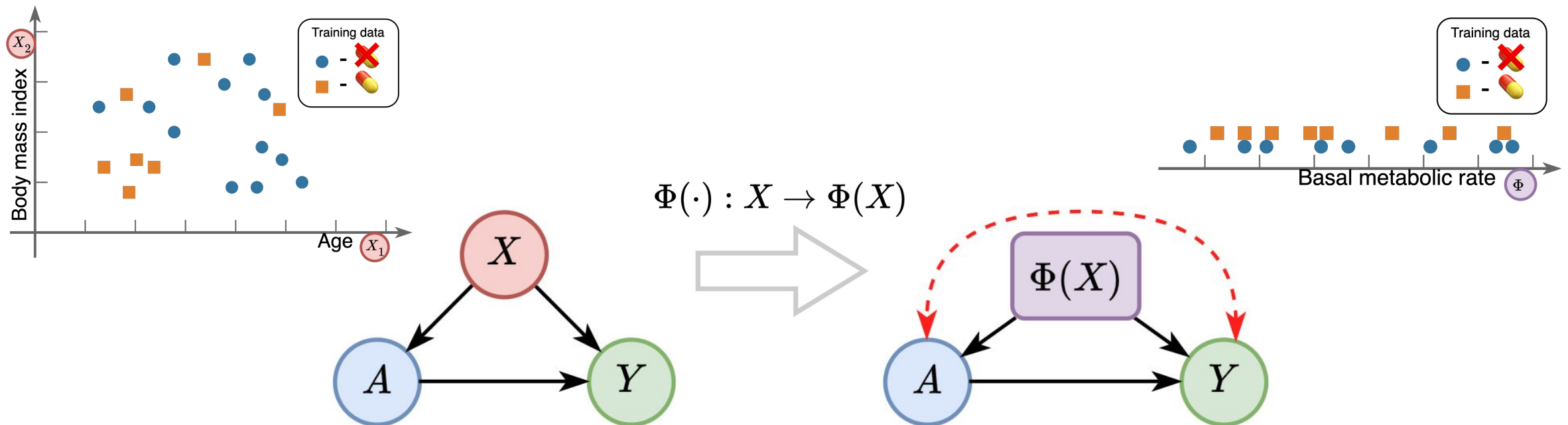
Model-based learners: Representation learning

Example: TarNET / CFRNet for CATE

$$\tau(x) = \mathbb{E}(Y[1] - Y[0] \mid x)$$

conditional average treatment effect (CATE)

Method: Learning a low-dimensional (balanced) representation $\Phi(\cdot)$ of high-dimensional covariates. Then, we can fit a CATE model based on the representations.



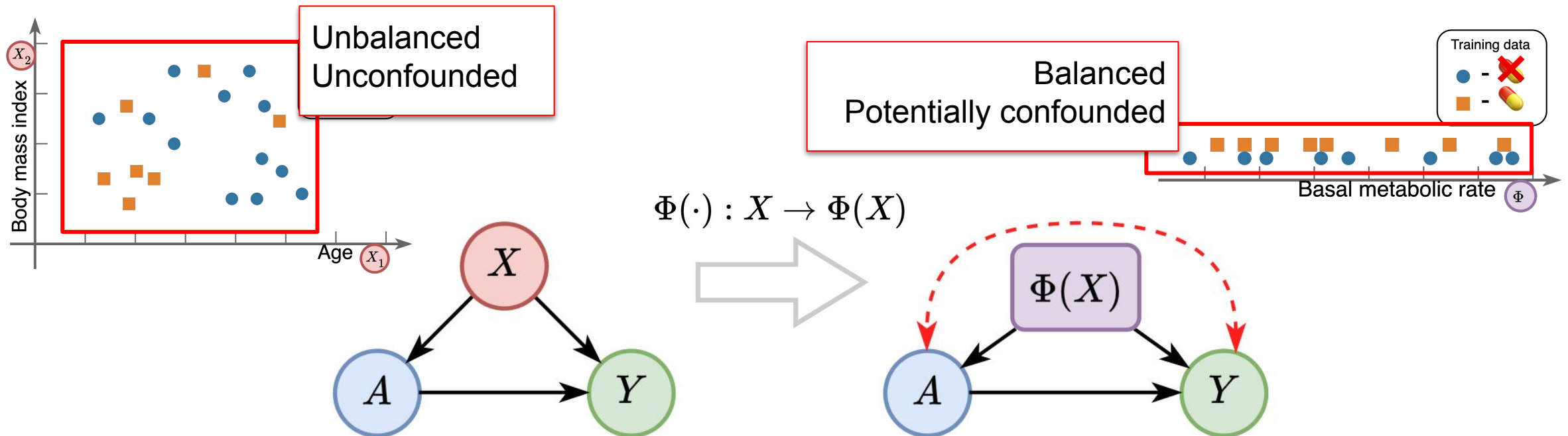
Model-based learners: Representation learning

Example: TarNET / CFRNet for CATE

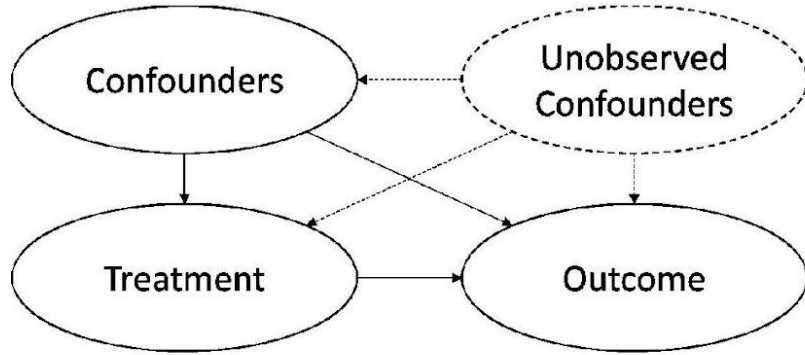
$$\tau(x) = \mathbb{E}(Y[1] - Y[0] \mid x)$$

conditional average treatment effect (CATE)

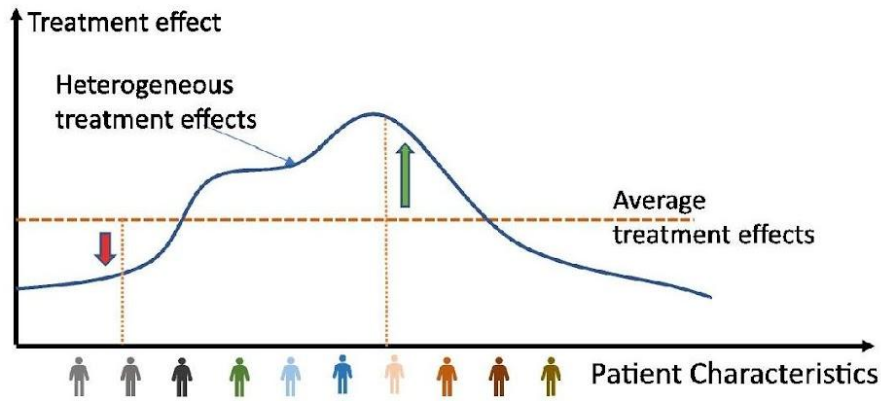
Method: Learning a low-dimensional (balanced) representation $\Phi(\cdot)$ of high-dimensional covariates. Then, we can fit a CATE model based on the representations.



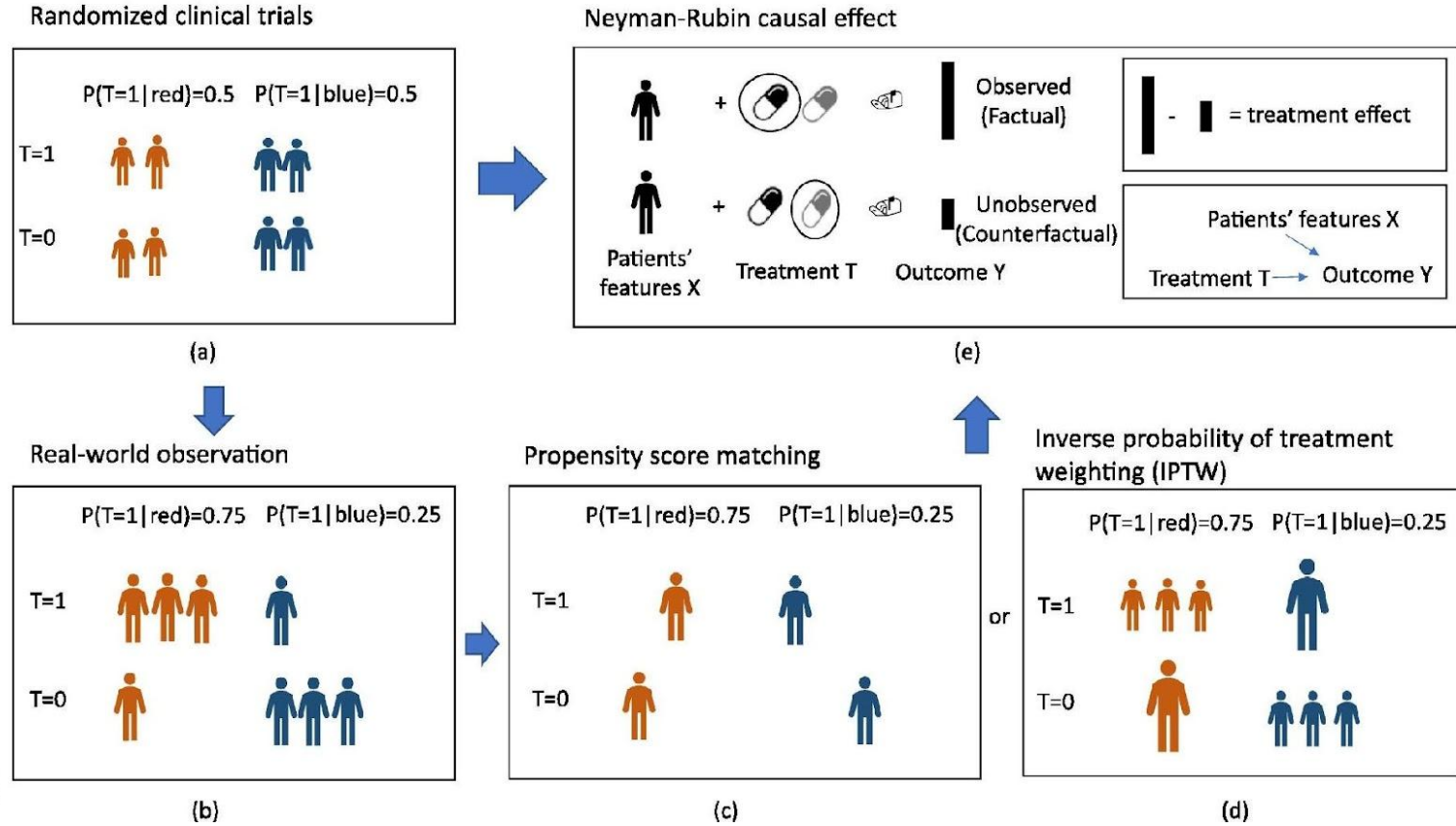
A. An example of treatment effects



C. An illustration of heterogeneous treatment effects



B. Comparison between estimating treatment effects from RCTs and from observational data



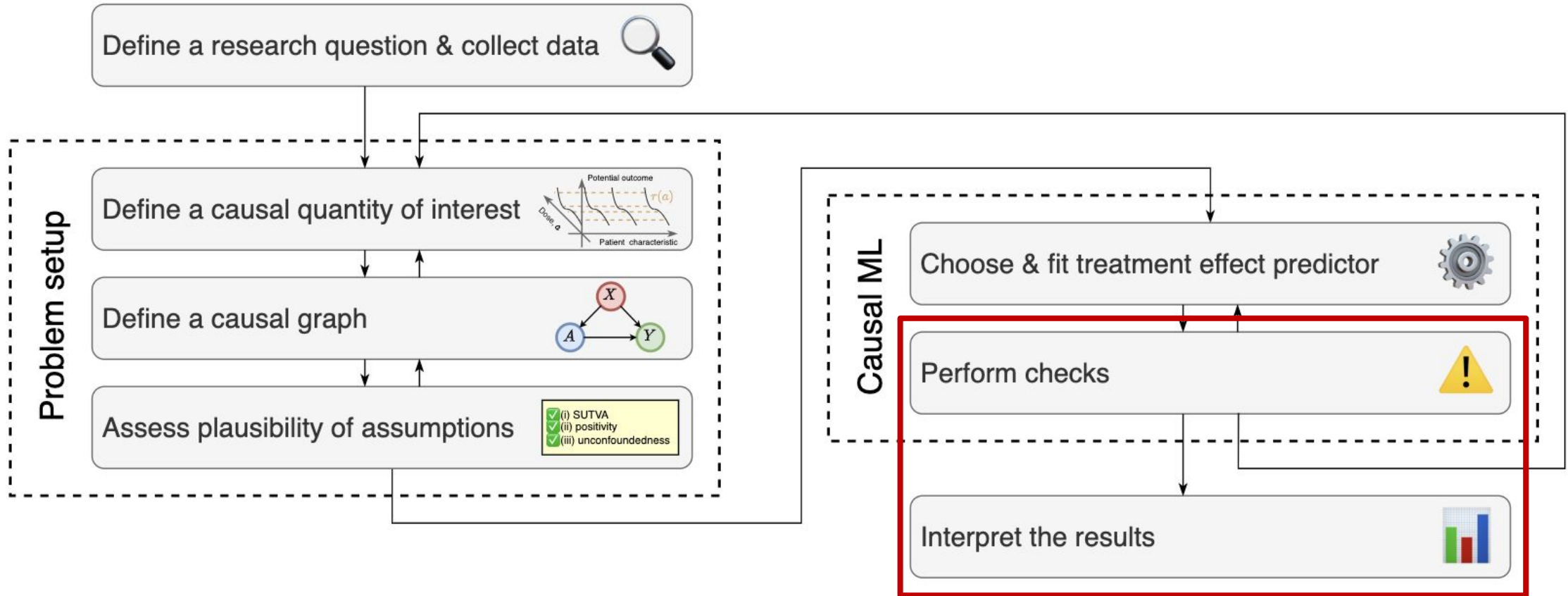


Munich Center for Machine Learning

Where we are (and what is still needed): Current state of causal ML research



Causal ML Workflow



Extensions & Open research problems

1 Model validity

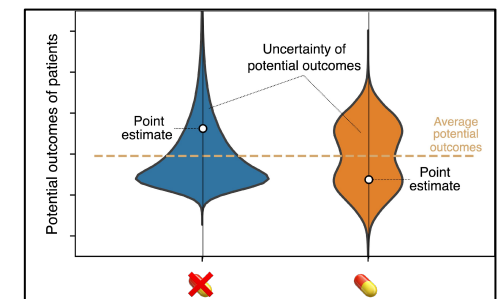
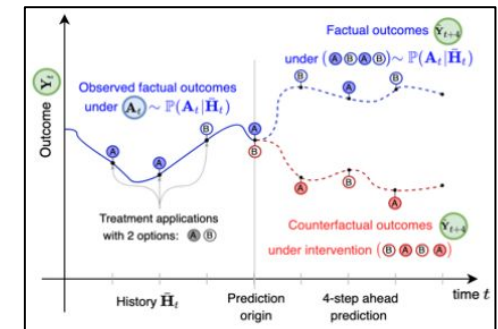
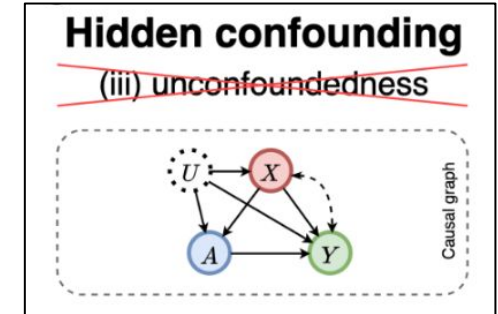
- Selection and validation of CATE models
 - Unlike traditional ML, we do not have a ground truth validation subset
- Robustness checks wrt. violation of assumptions
 - Sensitivity models
 - Spillover effects

2 Flexibility

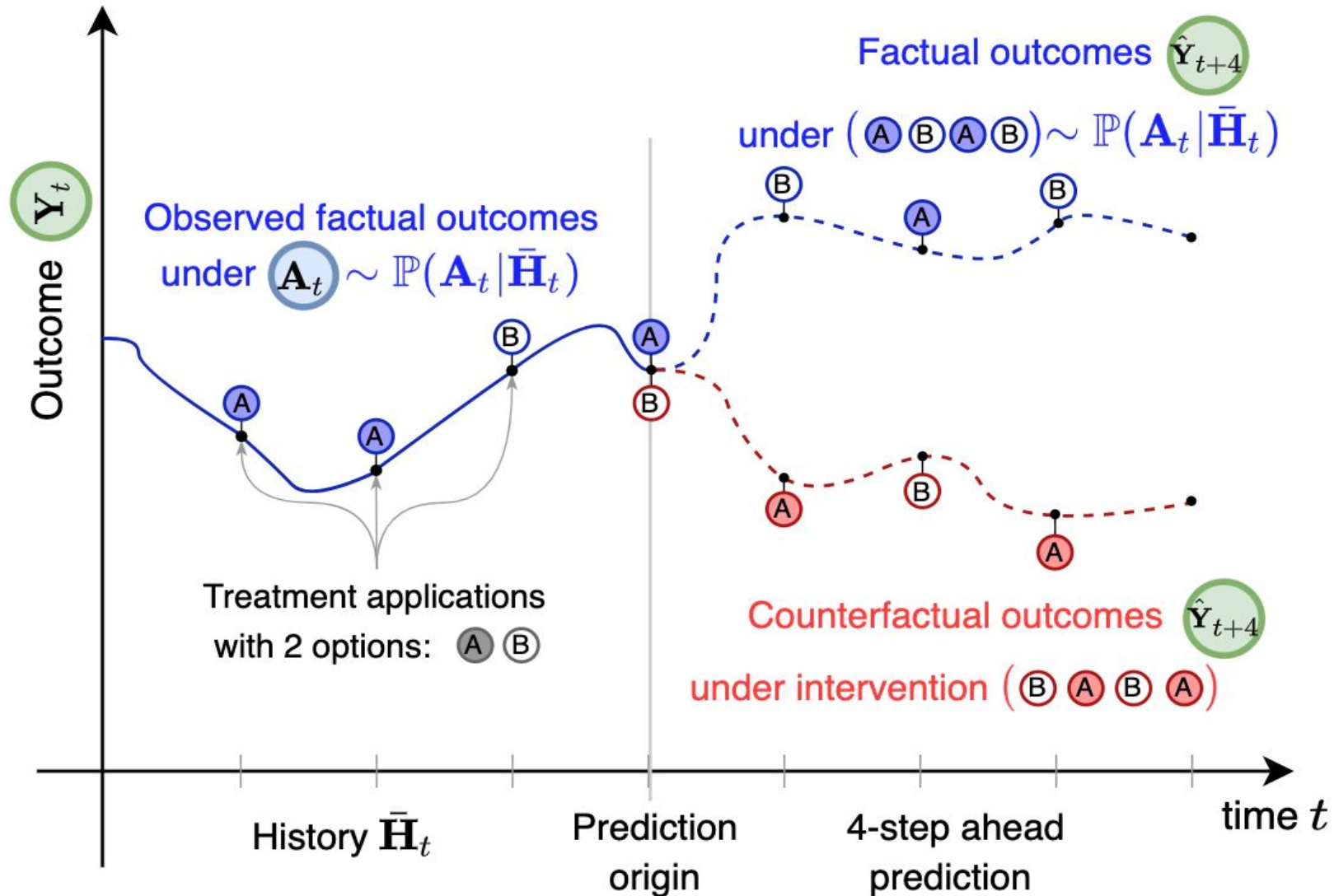
- Extensions to more complicated settings
 - continuous / high-dimensional treatments
 - time-varying potential outcomes and treatment effects
- Data fusion from multiple environments
- Constrained ML: interpretability, privacy enforcement

3 Uncertainty quantification

- Uncertainty quantification
 - uncertainty of estimation (aka confidence intervals)
 - predictive uncertainty (aka predictive intervals)

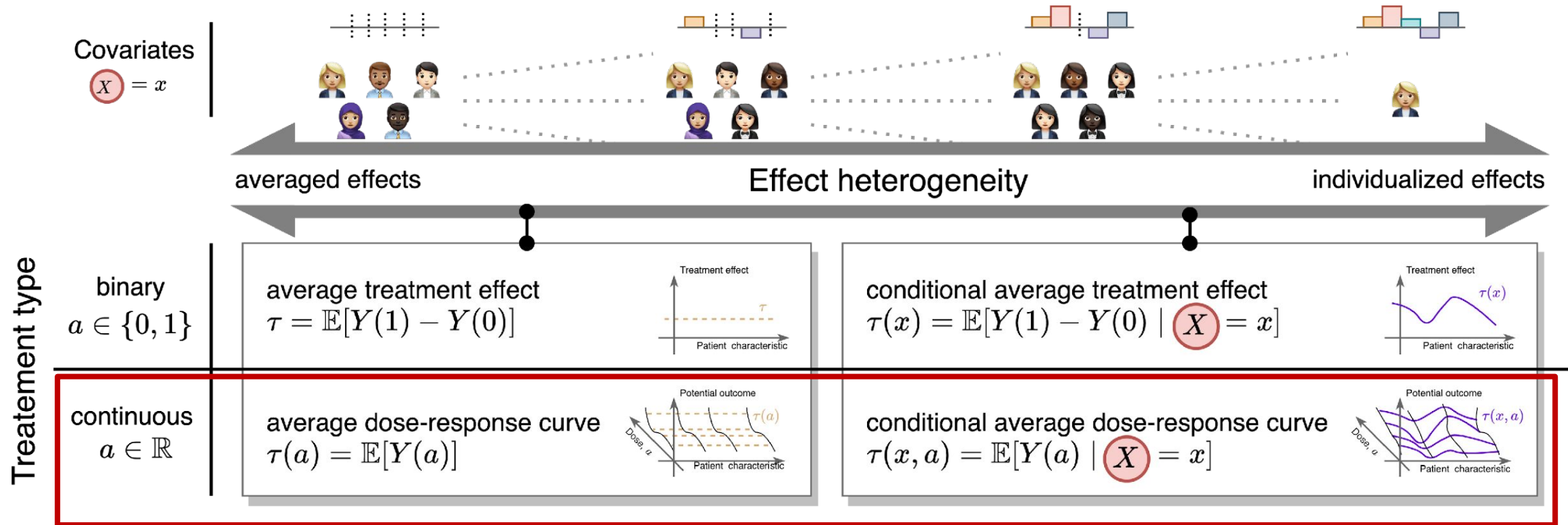


Flexibility: Causal ML for predicting outcomes over time

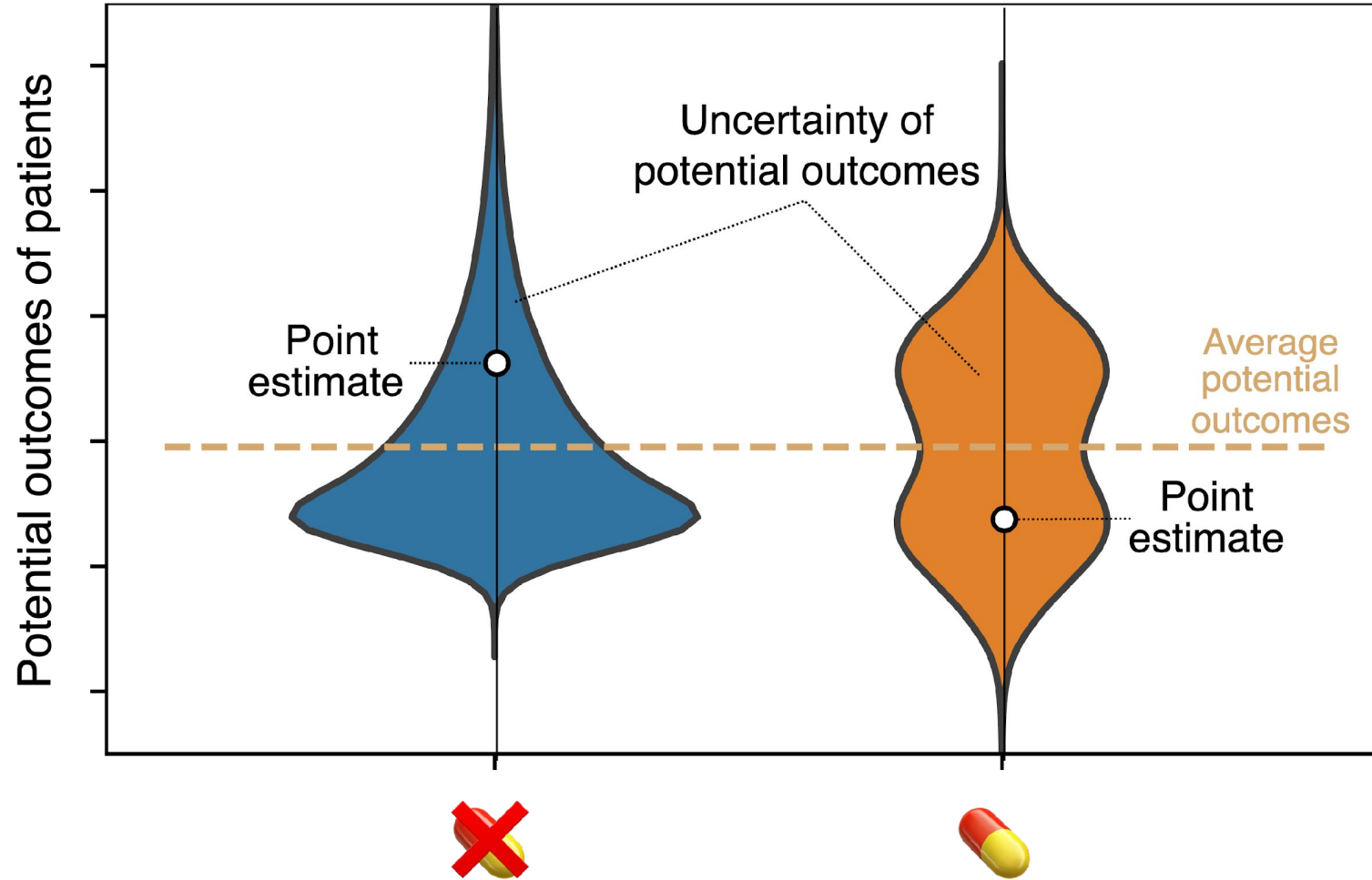


EXTENSIONS & OPEN RESEARCH QUESTIONS

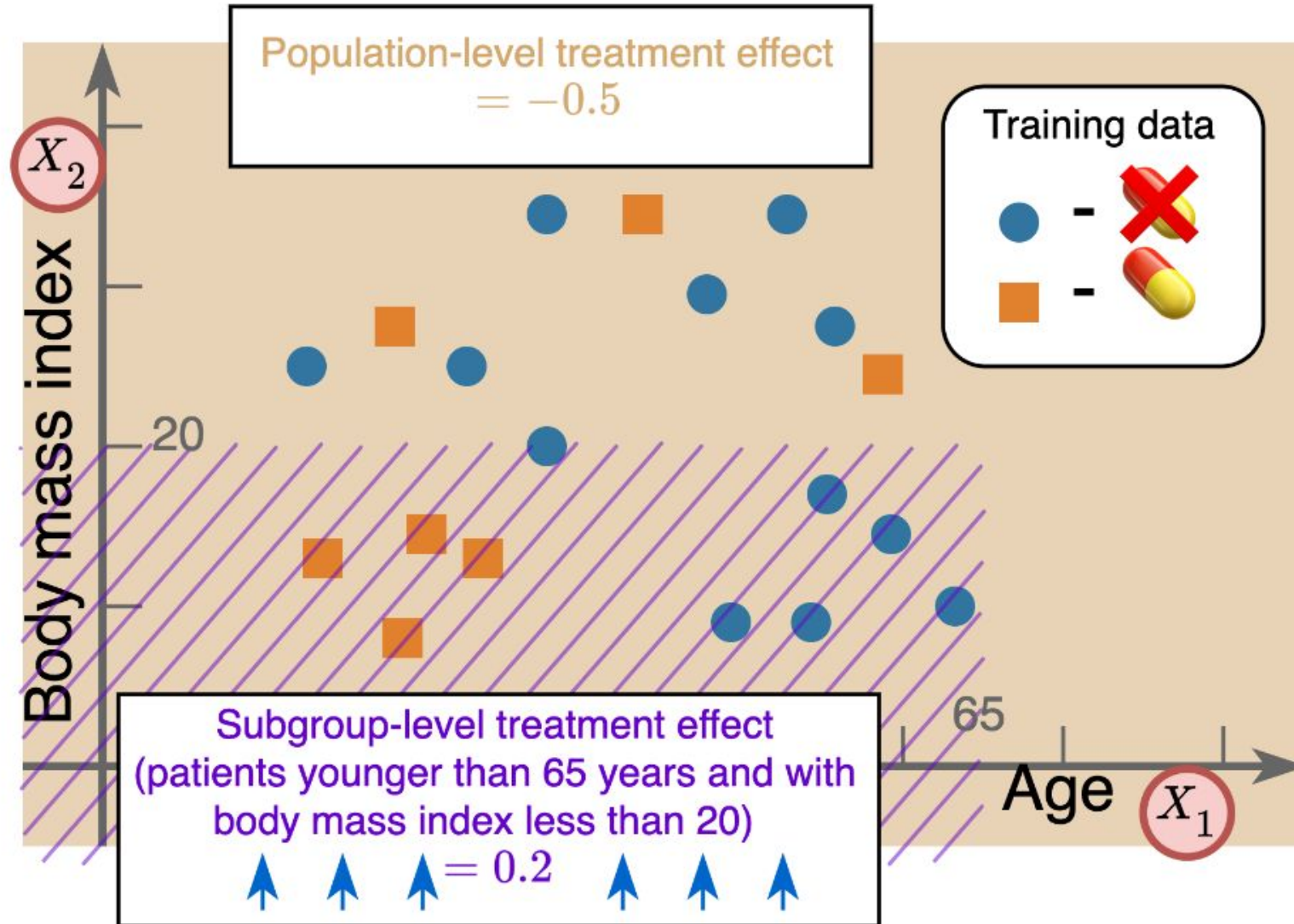
Flexibility: Continuous / high-dimensional treatments



Uncertainty quantification



Identifying predictive biomarkers (=treatment responders)



VISION

Promises of Causal ML

Estimating treatment effects for vulnerable groups



Augmenting evidence from RCTs



Finding optimal dosages



ML for treatment effect estimation

Estimating post-approval efficacy, including side effects



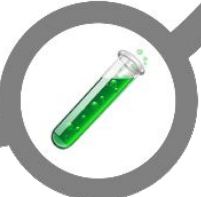
Guiding treatment choice when a standard of care is absent



Estimating treatment effects for long-term outcomes



Designing treatment recommendations for rare diseases





LUDWIG-
MAXIMILIANS-
UNIVERSITÄT
MÜNCHEN

LMU MUNICH
SCHOOL OF
MANAGEMENT

INSTITUTE OF AI IN MANAGEMENT



Institute of AI in Management
Prof. Dr. Stefan Feuerriegel

<http://www.ai.bwl.lmu.de>

 @stfeuerriegel  stefan-feuerriegel

Artificial intelligence | Impact