

Generalizing a causal effect from a trial to a target population

Bénédicte Colnet — Wednesday, 28 June 2023 — Ph.D. Defense

Jury members

Ph.D. advisors

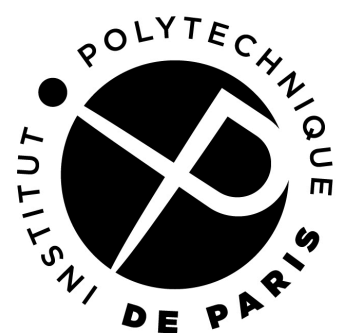
- Julie Josse (Inria)
- Erwan Scornet (École polytechnique)
- Gaël Varoquaux (Inria)

Reviewers

- Nicolai Meinshausen (ETH Zürich)
- Stinj Vansteelandt (Ghent University)

Examiners

- Trevor Hastie (Stanford)
- Erwan Le Pennec (École polytechnique)
- Elizabeth Ogburn (John Hopkins)
- Philippe Ravaud (Paris' hospitals)



ECOLE
DOCTORALE
DE MATHÉMATIQUES
HADAMARD

Inria



Outline

1. Introduction

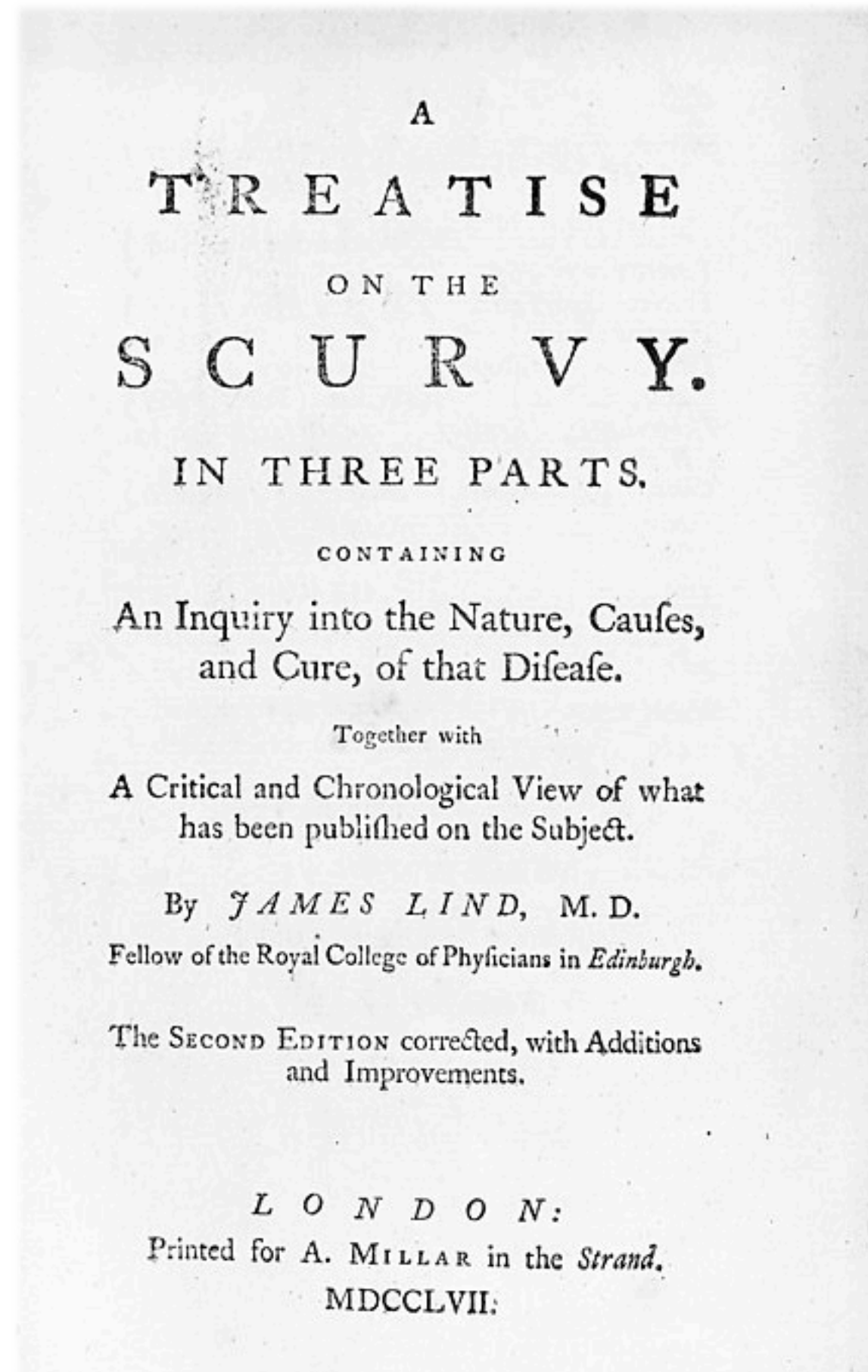
- A. Motivating example from critical care medicine
- B. State-of-the-art

— Focus on two contributions —

2. Finite and large sample analysis of the IPSW estimator

3. Extension to different causal measures

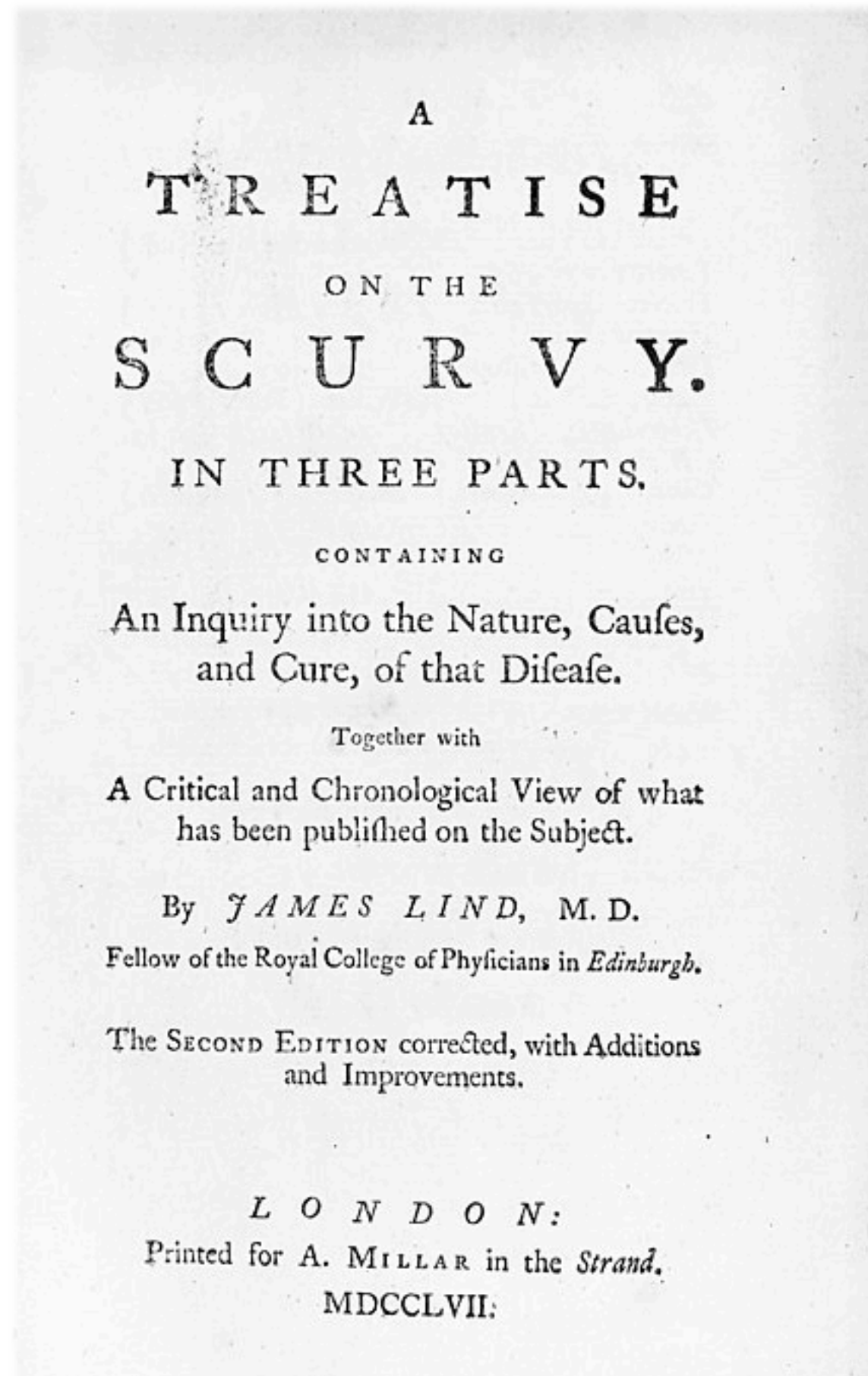
A longstanding presence of Randomized Controlled Trials (RCTs)



James Lind experiment on scorbout in **1757**

Source: Wikipedia

A longstanding presence of Randomized Controlled Trials (RCTs) ... now being the gold-standard



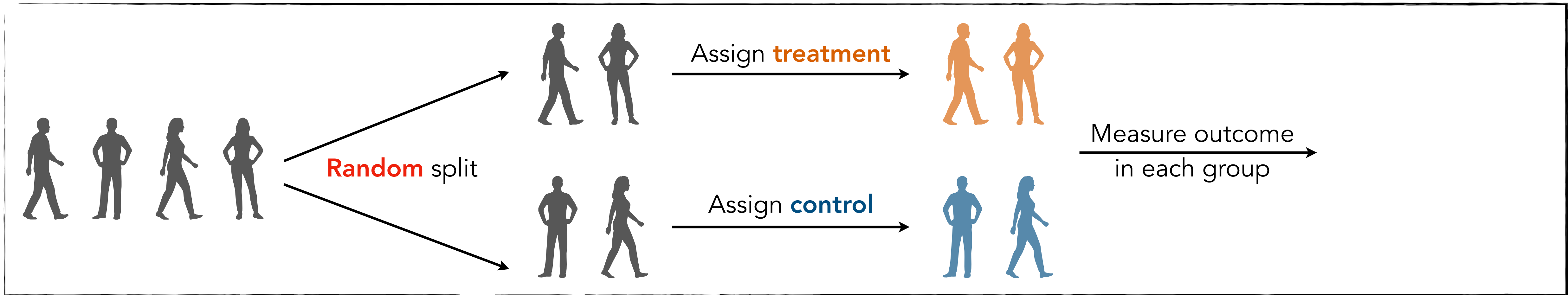
James Lind experiment on scorbout in **1757**
Source: Wikipedia

Drug Trials Snapshot	Active Ingredient	Date of FDA Approval	What is it Approved For
CABENUVA	cabotegravir and rilpivirine	January 20, 2021	Treatment of HIV-1 infection.
LUPKYNIS	voclosporin	January 22, 2021	Treatment of lupus nephritis
VERQUVO	vericiguat	January 19, 2021	Treatment of chronic heart failure
GEMTESA	vibegron	December 23, 2020	Treatment of symptoms of overactive bladder
EBANGA	ansuvimab-zykl	December 21, 2020	Treatment of Zaire ebolavirus infection
ORGOVYX	relugolix	December 18, 2020	Treatment of advanced prostate cancer

Recently approved drugs by the Food and Drug Administration (FDA), all with their corresponding RCT snapshot and information.
Source: www.fda.gov - **2022**

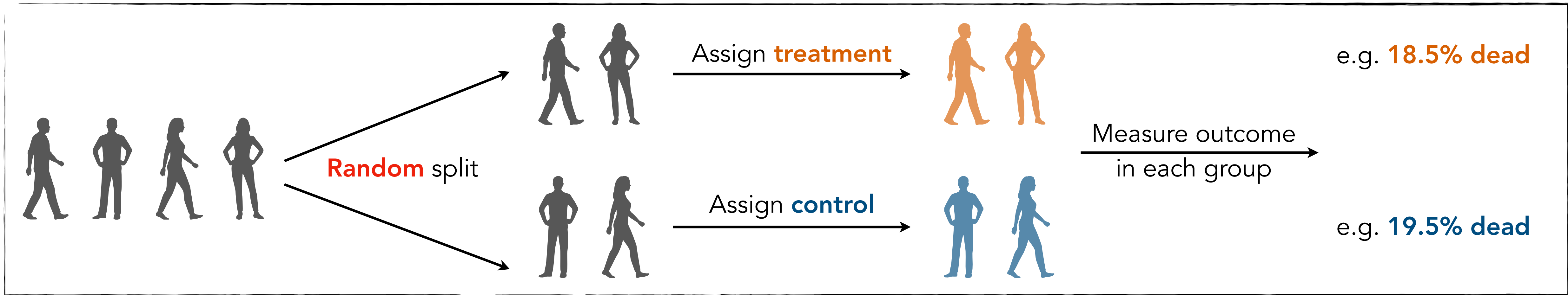
RCTs' principle : estimating a causal effect

Principle



RCTs' principle : estimating a causal effect

Principle

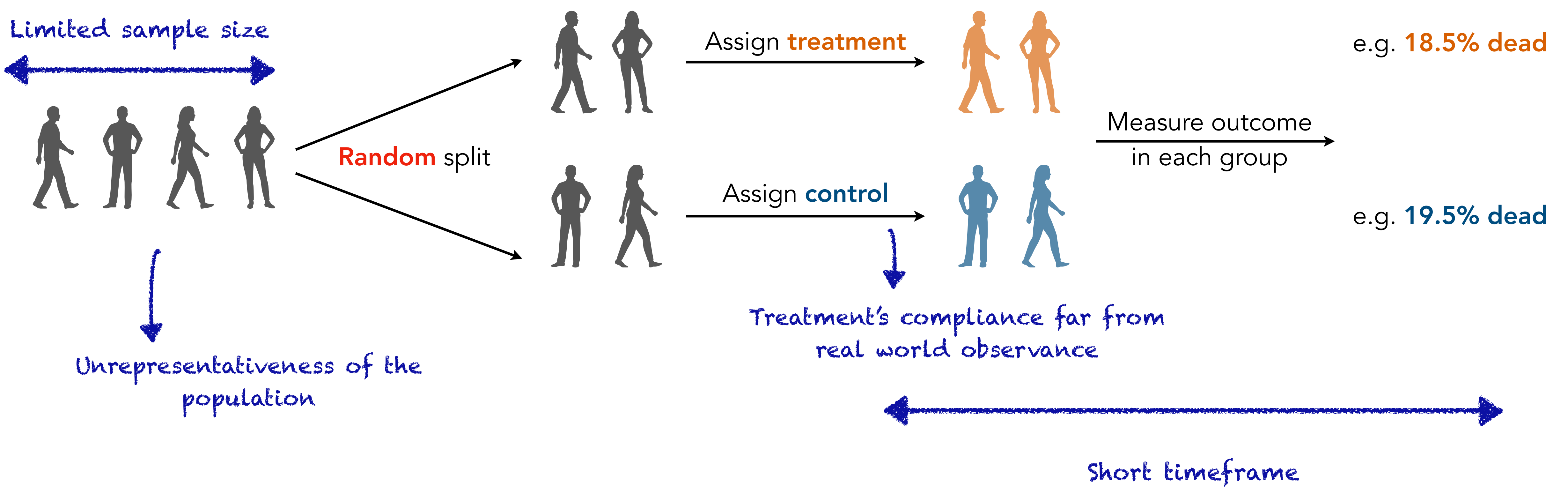


In practice : the CRASH-3 trial investigating Tranexamic Acid effect on brain injured (TBI) related death

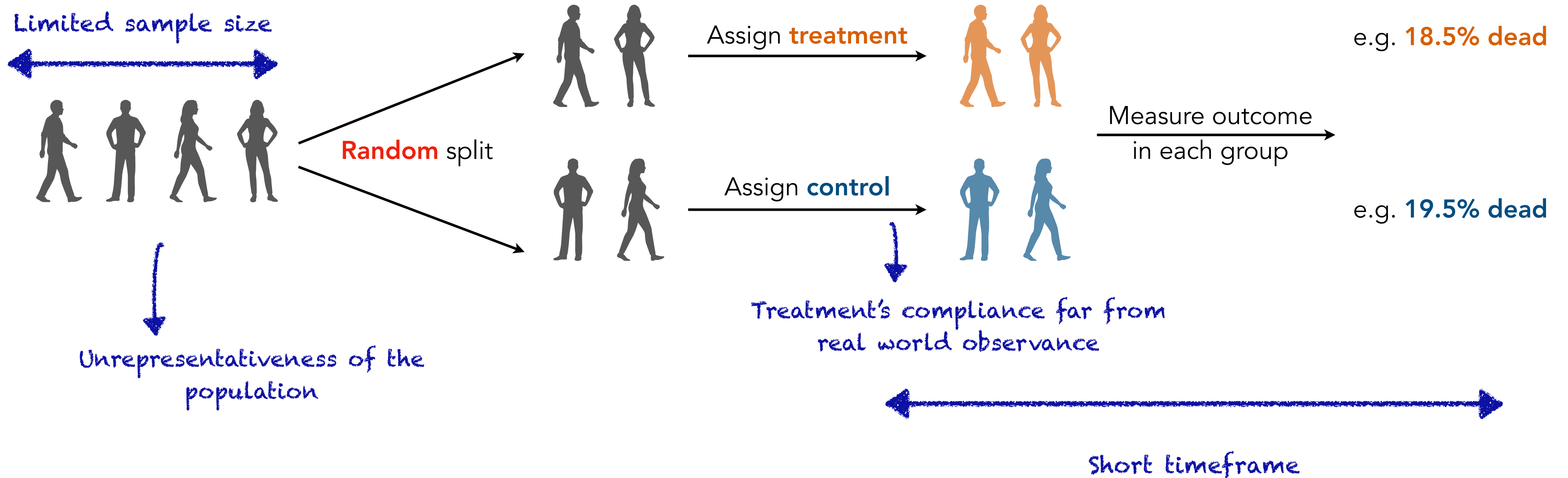
Results Between July 20, 2012, and Jan 31, 2019, we **randomly** allocated 12 737 patients with TBI to receive **tranexamic acid** (6406 [50·3%] or **placebo** [6331 [49·7%], of whom 9202 (72·2%) patients were treated within 3 h of injury. Among patients treated within 3 h of injury, the risk of head injury-related death was **18·5%** in the tranexamic acid group versus **19·8%** in the placebo group (855 vs 892 events; risk ratio [RR] 0·94 [95% CI 0·86–1·02]).

Source: Screenshot from the Lancet (CRASH-3 main report)

The scope of RCTs is increasingly under scrutiny



The scope of RCTs is increasingly under **scrutiny**



“‘External validity’ asks the question of generalizability: to what populations, settings, treatment variables, and measurement variables can this effect be generalized?” — Campbell and Stanley (1963), p. 5

The **promise** of detailed and larger observational or *real world* data sets

Estimate the efficacy in real-world conditions

- Using large cohorts like hospital data bases
 - To **emulate a target trial**⁽¹⁾ leveraging observed confounding variables
 - Solving both representativity and effective treatment given
- 📦 *Large sample enabling more personalization (i.e stratified effects)*

(1) Hernán and Robins, Using Big Data to Emulate a Target Trial When a Randomized Trial Is Not Available, *Am J Epidemiol*, 2016



Source: FDA's website

The example of a large French national cohort — The Traumabase

- 30,000 patients of unique size and granularity in Europe (~9,000 suffering from TBI)
- But randomisation does not hold, e.g. severe trauma are more likely to be treated

Among control
16% dead

Among treated
38% dead



Confusion problem

The example of a large French national cohort — The Traumabase

- 30,000 patients of unique size and granularity in Europe (~9,000 suffering from TBI)
- But randomisation does not hold, e.g. severe trauma are more likely to be treated



After adjustment on confounding covariates (Glasgow score, age, blood pressure, ...), the null hypothesis of no effect can not be rejected⁽²⁾.

CRASH-3 key results

Is there a paradox?

The risk of head injury-related death reduced with tranexamic acid in patients with mild-to-moderate head injury (RR 0.78 [95% CI 0.64–0.95]) but not in patients with severe head injury (0.99 [95% CI 0.91–1.07])

(2) Mayer et al., Doubly robust treatment effect estimation with missing attributes, *Annals of Applied Statistics* 2019

Idea — Using both types of data : experimental and observational

Fear of **unobserved confounding** in the observational sample.

Idea — Using both types of data : experimental and observational

Fear of **unobserved confounding** in the observational sample.

Both **Randomized Controlled Trial (RCT)** data and **observational** data have limitations and advantages.

The idea is to **combine** them to get the **best of both worlds**.

Causal inference methods for combining randomized trials and observational studies: a review

Bénédicte Colnet¹, Imke Mayer¹, Guanhua Chen, Awa Dieng, Ruohong Li, Gaël Varoquaux, Jean-Philippe Vert, Julie Josse², Shu Yang²

Accepted for publication in Statistical Science

Idea — Using both types of data

Fear of **unobserved confounding** in the observational sample.

Both **Randomized Controlled Trial (RCT)** data and **observational** data have limitations and advantages.

The idea is to **combine** them to get the **best of both worlds**.

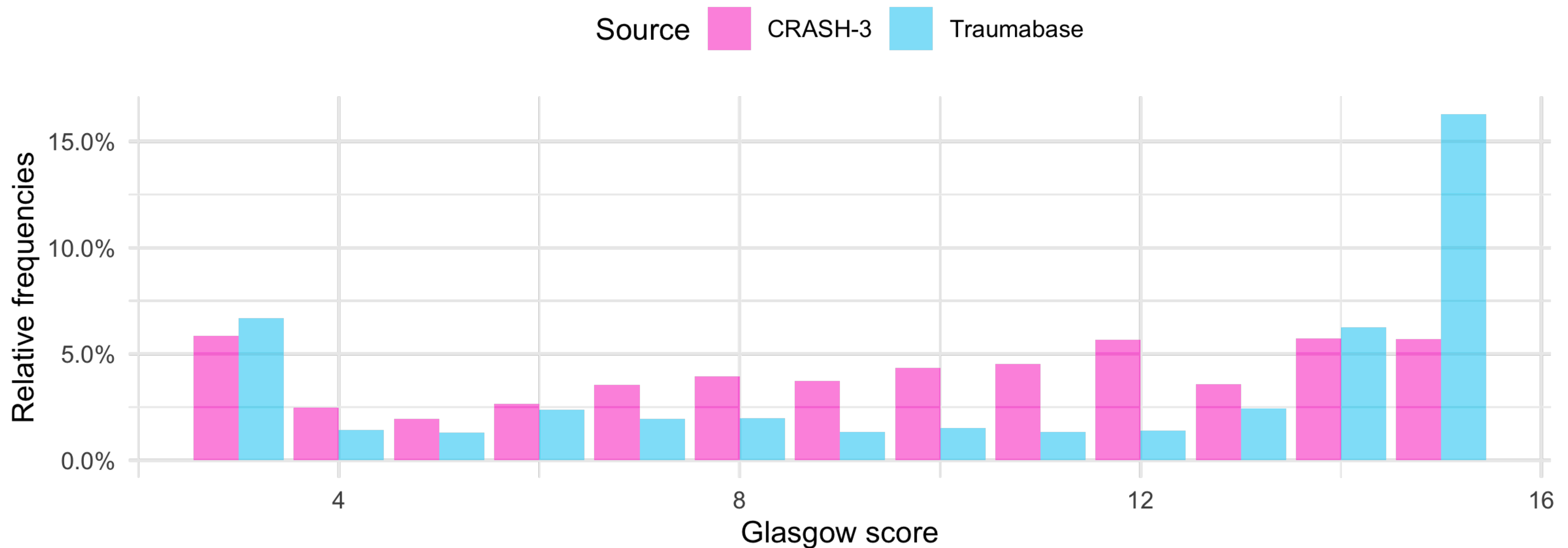
Causal inference methods for combining randomized trials and observational studies: a review

Bénédicte Colnet¹, Imke Mayer¹, Guanhua Chen, Awa Dieng, Ruohong Li, Gaël Varoquaux, Jean-Philippe Vert, Julie Josse², Shu Yang²

— Using observational data to improve trial's representativity

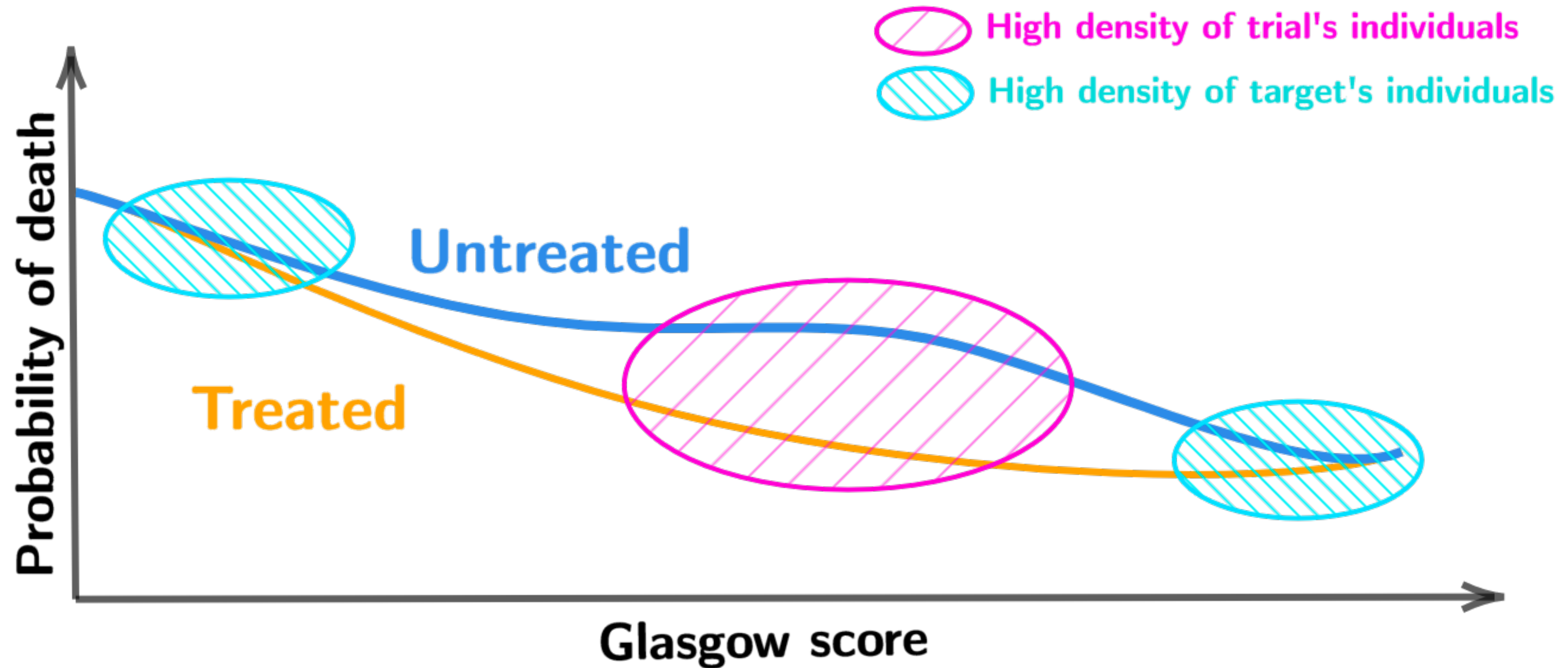
Accepted for publication in Statistical Science

Generalizing or transporting CRASH-3 findings to the Traumabase population



“What would have been measured as an effect in CRASH-3 if the trial was sampled in the Traumabase?”

Generalizing or transporting CRASH-3 findings to the Traumabase population



Hypothetical drawing of how the Glasgow score could modulate treatment effect

State-of-the art in a Nutshell

- Foundational work in epidemiological books (Rothman & Greenland, 2000)
- Idea of using two data sets (Stuart et al. 2010 and Pearl & Barenboim 2011)
- Flourishing field in statistics!
- Usually clinical papers focus on characterising the lack of representativeness
 - Comparison of Table 1
 - % of patients actually treated that would have been eligible

Notations

For each individual i , consider each of the possible outcomes for **treated** $Y^{(1)}$, and **control** $Y^{(0)}$.

characteristics binary treatment

	X	A	$Y^{(1)}$	$Y^{(0)}$	Y
F	1	0	NA	0	0
M	2	0	NA	1	1
M	1	1	0	NA	0
F	3	0	NA	1	1
F	2	1	1	NA	1

Comparison of two potential outcomes

Individual effect $\Delta_i := Y_i^{(1)} - Y_i^{(0)}$

Notations

For each individual i , consider each of the possible outcomes for **treated** $Y^{(1)}$, and **control** $Y^{(0)}$.

characteristics binary treatment

	X	A	$Y^{(1)}$	$Y^{(0)}$	Y
F	1	0	NA	0	0
M	2	0	NA	1	1
M	1	1	0	NA	0
F	3	0	NA	1	1
F	2	1	1	NA	1

Individual effect $\Delta_i := Y_i^{(1)} - Y_i^{(0)}$

Can not be observed!

Average effect $ATE \equiv \tau := \mathbb{E} [\Delta_i]$

The potential outcomes framework for generalization

Denoting,

- **A** the binary treatment
- **X** the covariates
- **Y** the observed outcome

Two samples,

- A **trial** of size **n** sampled from a population **$p_R(X)$** ,
- A data set of size **m** sampled from **$p_T(X)$** the **target** population of interest.

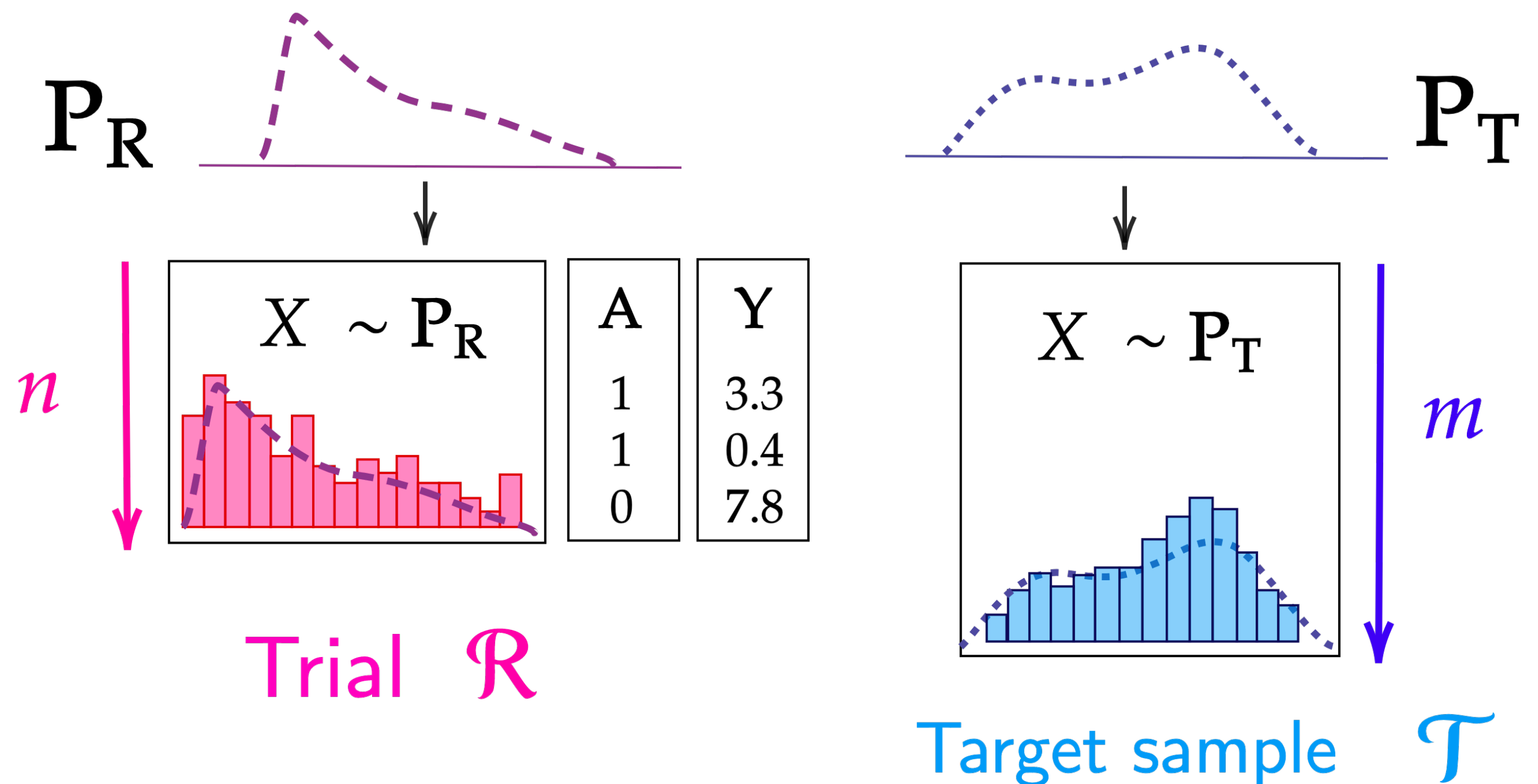
The potential outcomes framework for generalization

Denoting,

- \mathbf{A} the binary treatment
- \mathbf{X} the covariates
- \mathbf{Y} the observed outcome

Two samples,

- A **trial** of size n sampled from a population $p_{\mathbf{R}}(\mathbf{X})$,
- A data set of size m sampled from $p_{\mathbf{T}}(\mathbf{X})$ the **target** population of interest.



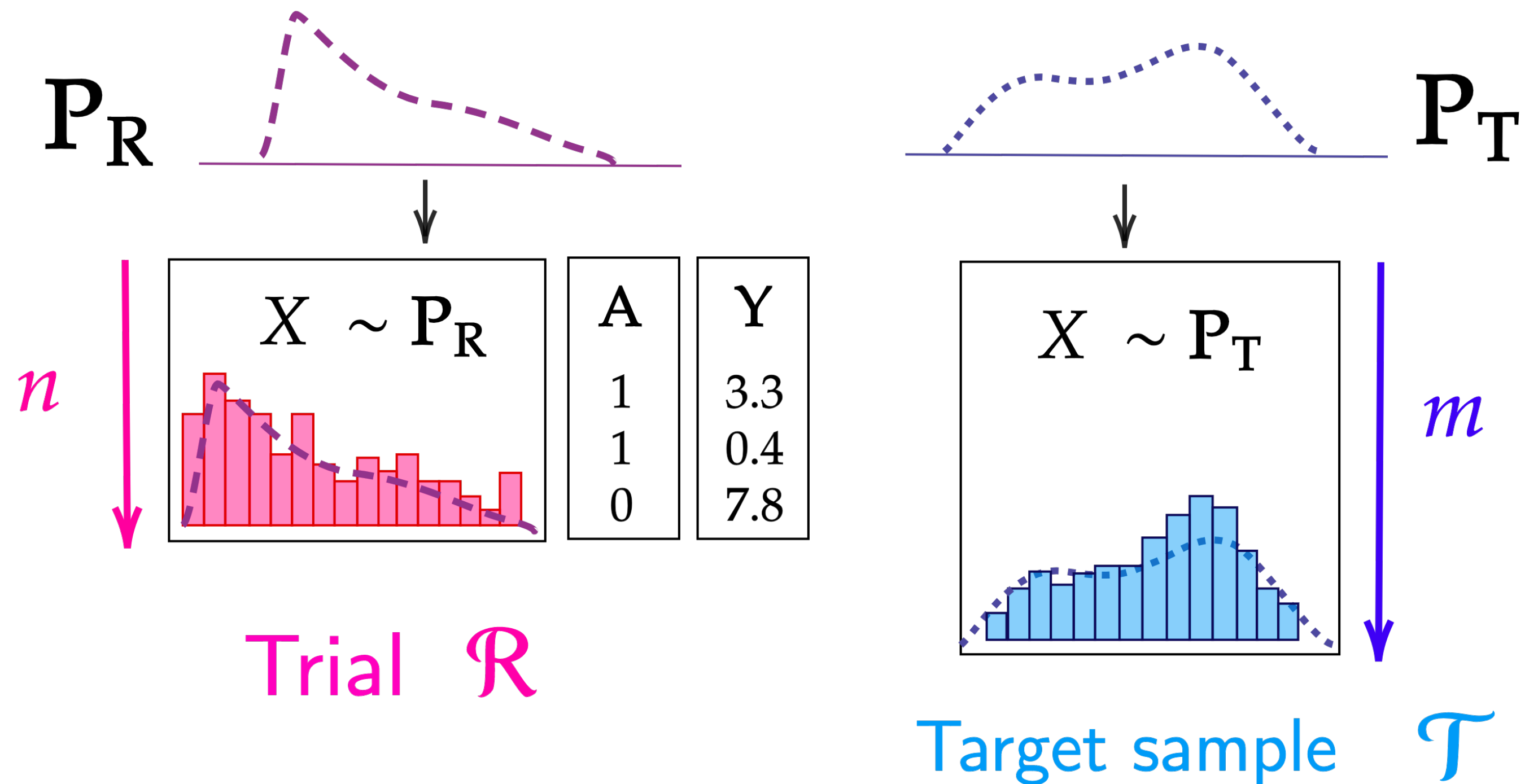
The potential outcomes framework for generalization

Denoting,

- \mathbf{A} the binary treatment
- \mathbf{X} the covariates
- \mathbf{Y} the observed outcome

Two samples,

- A **trial** of size n sampled from a population $p_R(\mathbf{X})$,
- A data set of size m sampled from $p_T(\mathbf{X})$ the **target** population of interest.



$$p_R(x) \neq p_T(x) \Rightarrow \underbrace{\tau_R := \mathbb{E}_R[Y(1) - Y(0)]}_{\text{ATE in the RCT}} \neq \underbrace{\mathbb{E}_T[Y(1) - Y(0)] := \tau}_{\text{Target ATE}}$$

Generalization's *causal* assumptions

Transportability assumption

$$\forall x \in \mathbb{X}, \quad \mathbb{P}_R(Y^{(1)} - Y^{(0)} \mid X = x) = \mathbb{P}_T(Y^{(1)} - Y^{(0)} \mid X = x)$$

— Needed covariates are **shifted** treatment effect **modifiers**.

↑ spirit of ignobility assumption for a single observational data set

Generalization's *causal* assumptions

Transportability assumption

$$\forall x \in \mathbb{X}, \quad \mathbb{P}_R(Y^{(1)} - Y^{(0)} \mid X = x) = \mathbb{P}_T(Y^{(1)} - Y^{(0)} \mid X = x)$$

— Needed covariates are **shifted** treatment effect **modifiers**.

Several versions in practice

e.g. of a lighter version

$$\forall x \in \mathbb{X}, \quad \mathbb{E}_R [Y^{(1)} - Y^{(0)} \mid X = x] = \mathbb{E}_T [Y^{(1)} - Y^{(0)} \mid X = x]$$

Dahabreh et al. 2020

Most common notation where S denotes the sample's indicator

$$Y^{(1)} - Y^{(0)} \perp\!\!\!\perp S \mid X$$

Nguyen et al. 2017

$$\{Y^{(1)}, Y^{(0)}\} \perp\!\!\!\perp S \mid X$$

Stuart et al. 2011

 Stronger assumption

Generalization's *causal* assumptions

Transportability assumption

$$\forall x \in \mathbb{X}, \quad \mathbb{P}_R(Y^{(1)} - Y^{(0)} \mid X = x) = \mathbb{P}_T(Y^{(1)} - Y^{(0)} \mid X = x)$$

— Needed covariates are **shifted** treatment effect **modifiers**.

Positivity assumption

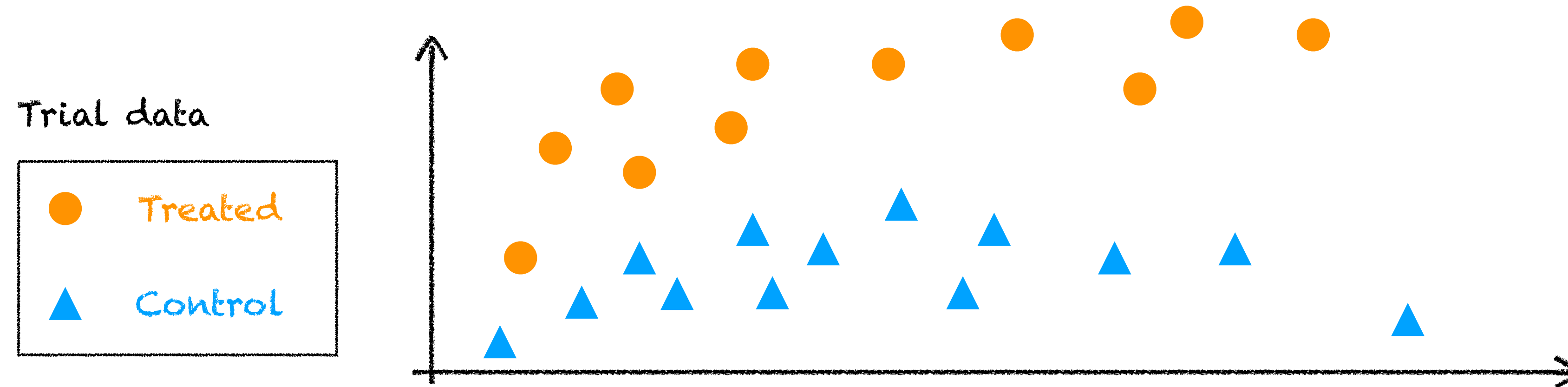
$$\text{supp}(P_T(X)) \subset \text{supp}(P_R(X))$$

— Each individuals in the target population has to be represented in the trial.

Also found as
 $P(S=1|X) > 1$

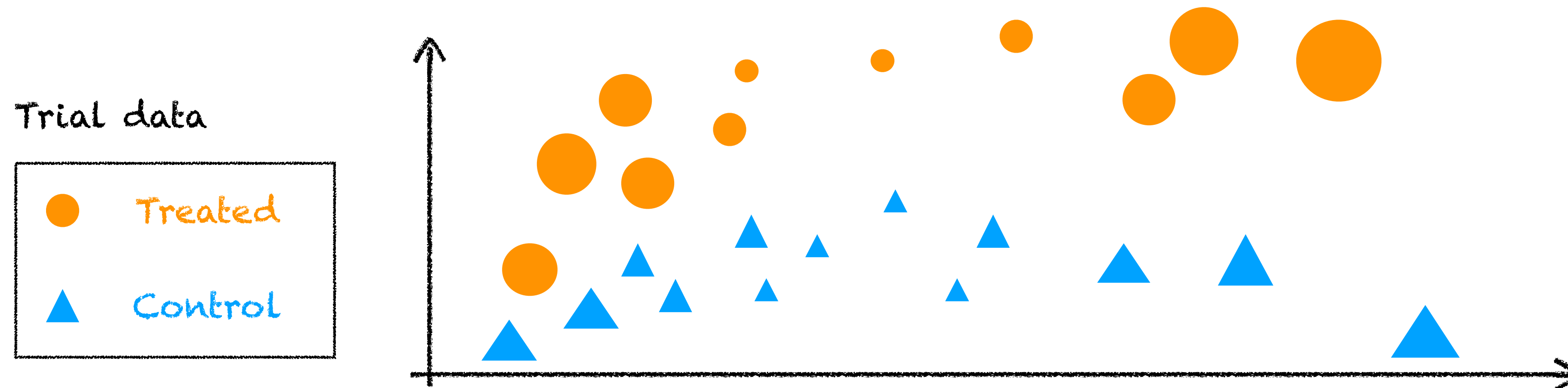
2 main approaches to generalize

1. **Re-weight** the trial individuals — *Inverse Propensity Sampling Weighting*



2 main approaches to generalize

1. **Re-weight** the trial individuals — *Inverse Propensity Sampling Weighting*



Definition

Spirit of IPW $\hat{\tau}_{IPSW,n,m} = \frac{1}{n} \sum_{i \in \text{Trial}} \hat{w}_{n,m}(X_i) \left(\frac{Y_i A_i}{\pi} - \frac{Y_i (1 - A_i)}{1 - \pi} \right)$

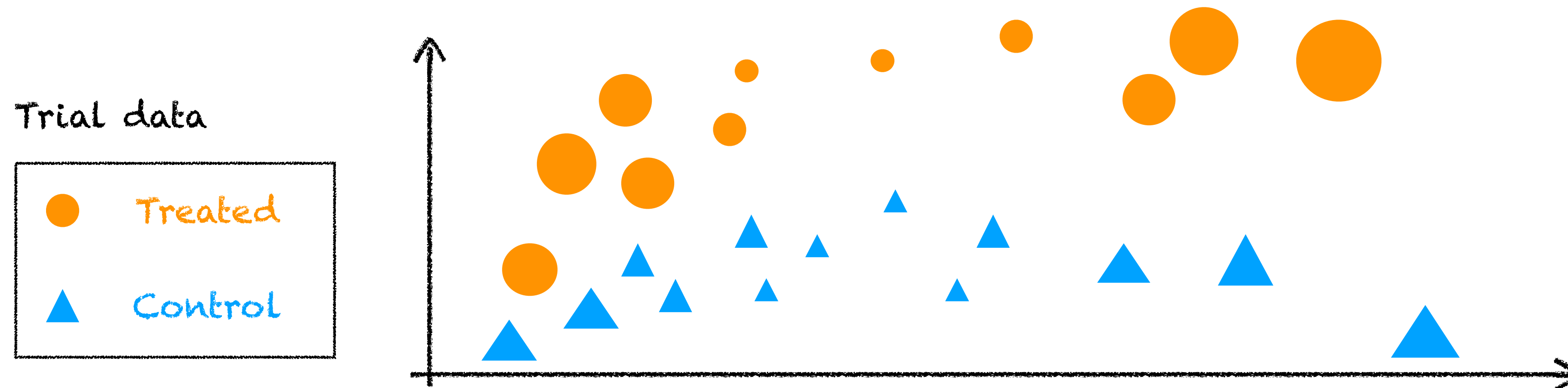
Weights (pointing to $\hat{w}_{n,m}(X_i)$)

Trial only (pointing to the fraction)

$\pi := P_{\text{RCT}}(A=1)$
Typically $\pi = 0.5$

2 main approaches to generalize

1. **Re-weight** the trial individuals — *Inverse Propensity Sampling Weighting*



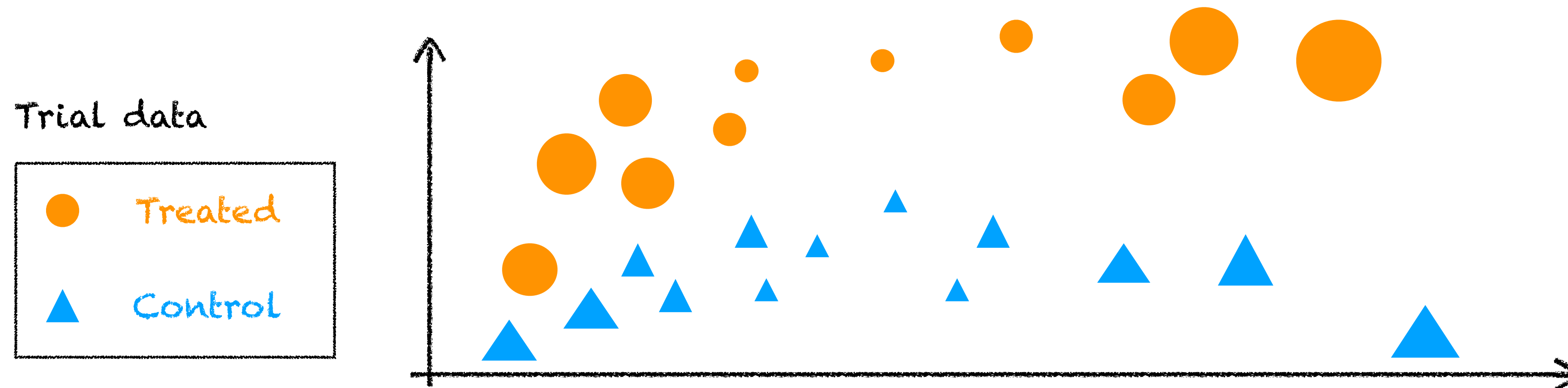
Definition

Spirit of IPW $\hat{\tau}_{IPSW,n,m} = \frac{1}{n} \sum_{i \in \text{Trial}} \frac{\hat{p}_{T,m}(X_i)}{\hat{p}_{R,n}(X_i)} \left(\frac{Y_i A_i}{\pi} - \frac{Y_i(1 - A_i)}{1 - \pi} \right)$

$\pi = P_{RCT}(A=1)$
Typically $\pi = 0.5$

2 main approaches to generalize

1. Re-weight the trial individuals — *Inverse Propensity Sampling Weighting*



Consistency Assuming that Y is square integrable, and that

$$(H1) \quad \sup_{x \in \mathcal{X}} \left| \hat{w}_{n,m}(x) - \frac{p_T(x)}{p_R} \right| = \epsilon_{n,m} \xrightarrow[n,m \rightarrow \infty]{a.s.} 0$$

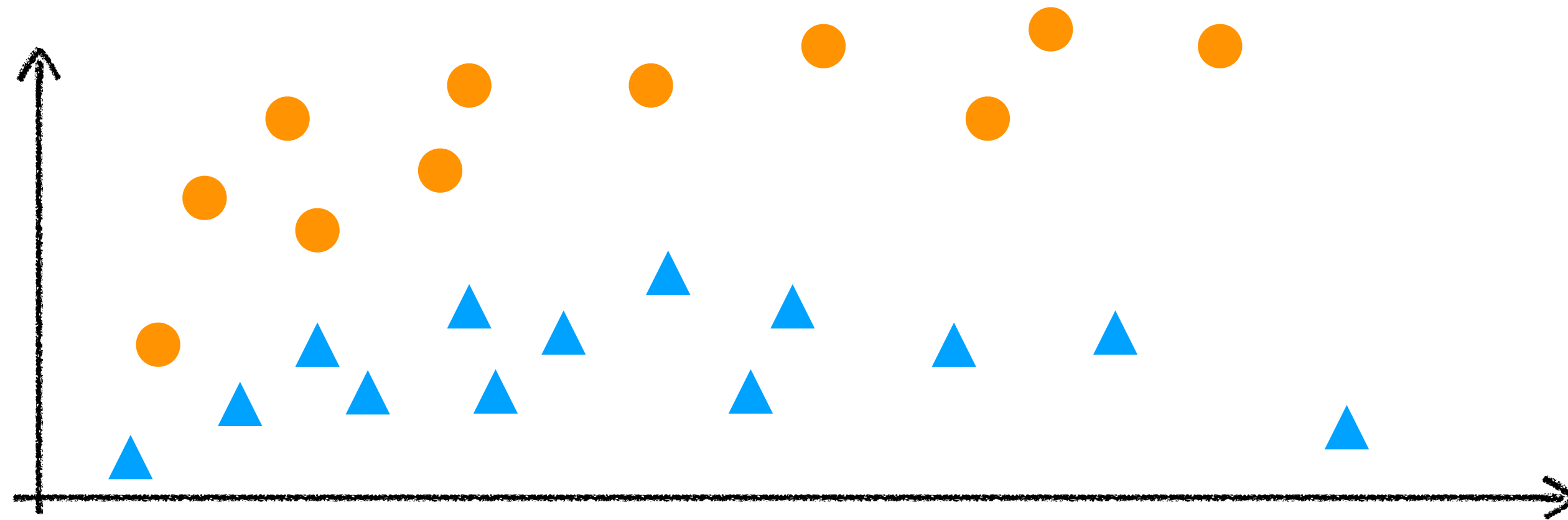
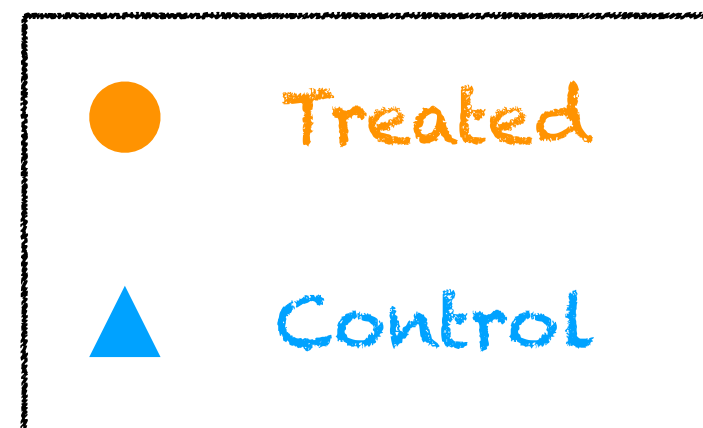
then $\hat{\tau}_{IPSW,n,m} \xrightarrow[n,m \rightarrow \infty]{L^1} \tau_T$

$$(H2) \quad \mathbb{E}[\epsilon_{n,m}^2] \xrightarrow[n,m \rightarrow \infty]{a.s.} 0$$

2 main approaches to generalize

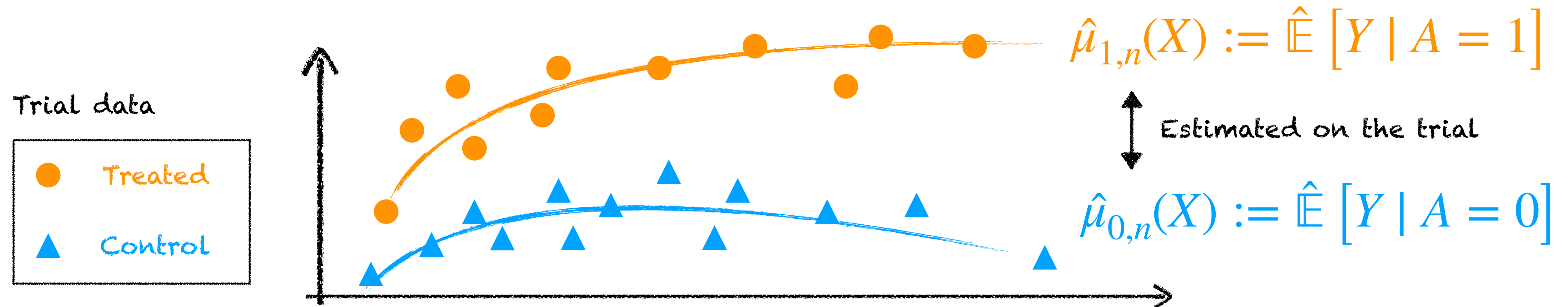
1. **Re-weight** the trial individuals — *Inverse Propensity Sampling Weighting*
2. **Model the response** on the trial and impute the target sample — *plug-in G-formula*

Trial data



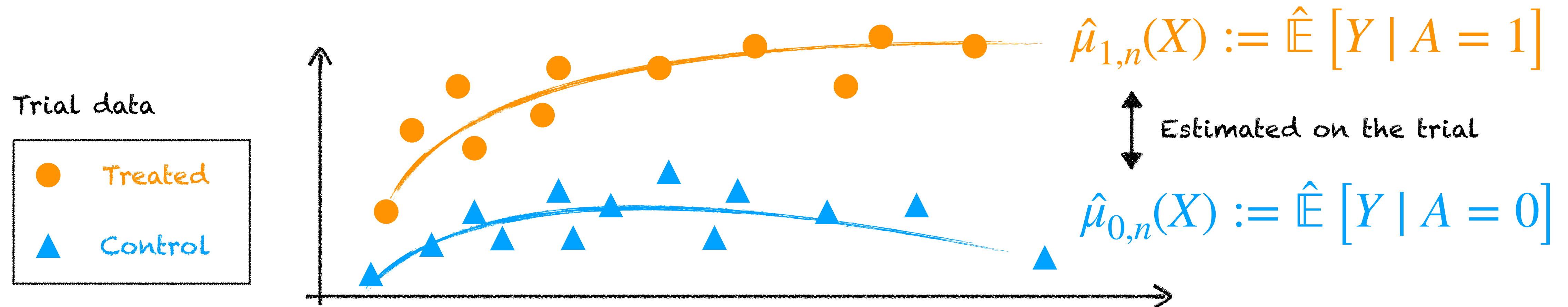
2 main approaches to generalize

1. **Re-weight** the trial individuals — *Inverse Propensity Sampling Weighting*
2. **Model the response** on the trial and impute the target sample — *plug-in G-formula*



2 main approaches to generalize

1. **Re-weight** the trial individuals — *Inverse Propensity Sampling Weighting*
2. **Model the response** on the trial and impute the target sample — *plug-in G-formula*



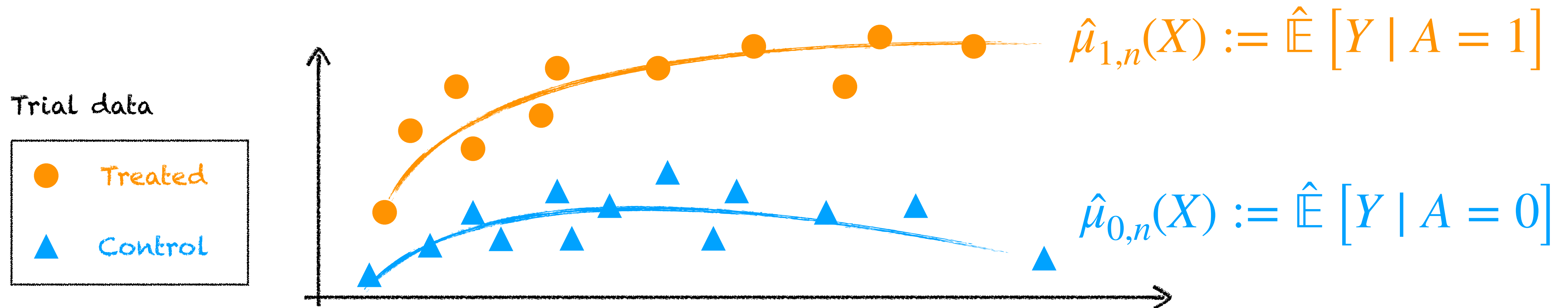
Definition

$$\hat{\tau}_{G,n,m} := \frac{1}{m} \sum_{i \in \text{Target}} \hat{\mu}_{1,n}(X_i) - \hat{\mu}_{0,n}(X_i)$$

Marginalised
on the target
sample

2 main approaches to generalize

1. **Re-weight** the trial individuals — *Inverse Propensity Sampling Weighting*
2. **Model the response** on the trial and impute the target sample — *plug-in G-formula*



Consistency

$$(H1) \quad \mathbb{E} \left[|\hat{\mu}_{a,n}(X) - \mu_a(X)| \mid T \right] \xrightarrow[n \rightarrow \infty]{p} 0$$

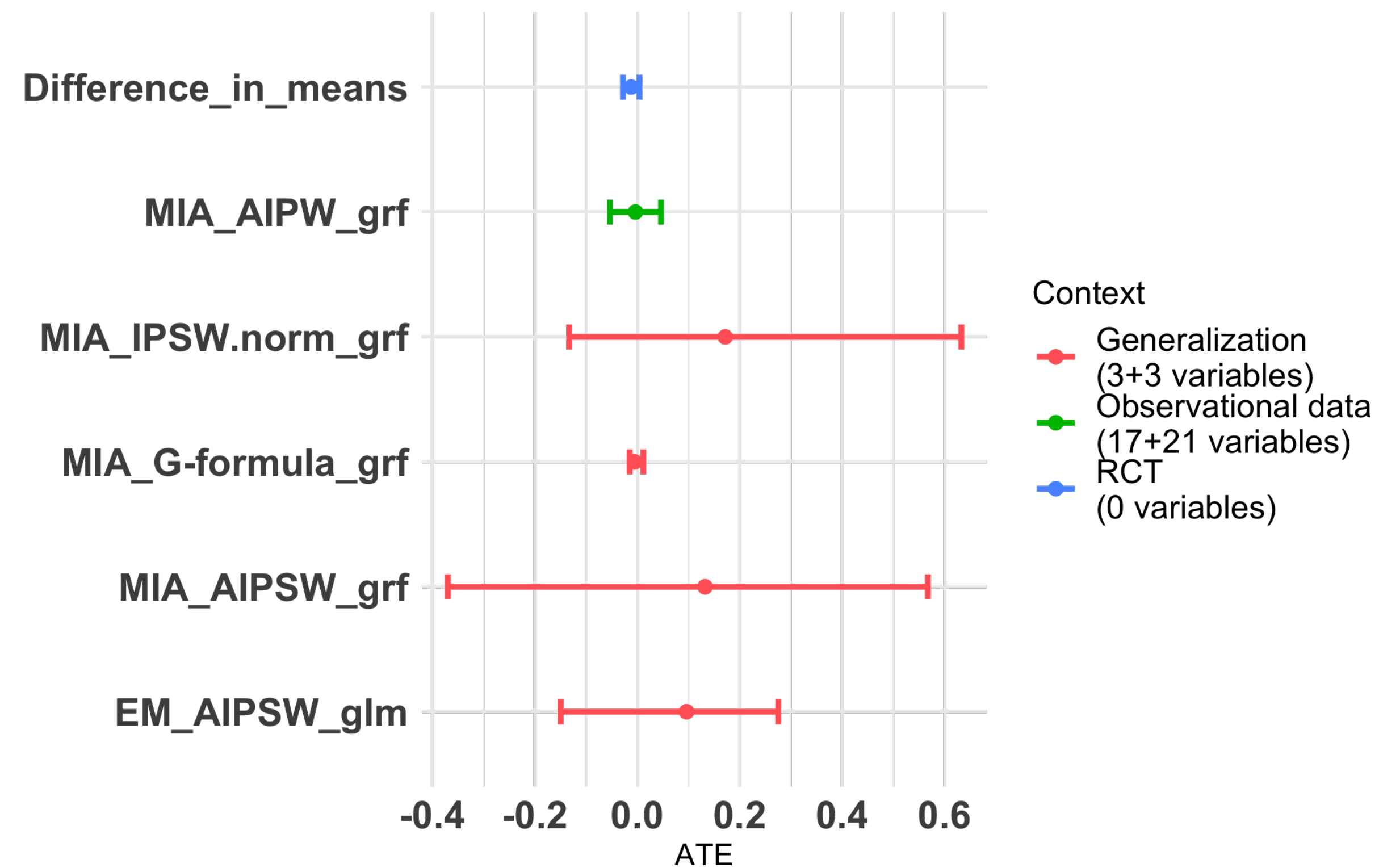
$$(H2) \quad \exists C_1, N_1 \quad \forall n \geq N_1, \quad \mathbb{E}[\hat{\mu}_{a,n}^2(X) \mid \mathcal{D}_n] \leq C_1$$

then

$$\hat{\tau}_{G,n,m} \xrightarrow[n,m \rightarrow \infty]{L^1} \tau_T$$

Application on the CRASH-3 & Traumabase example

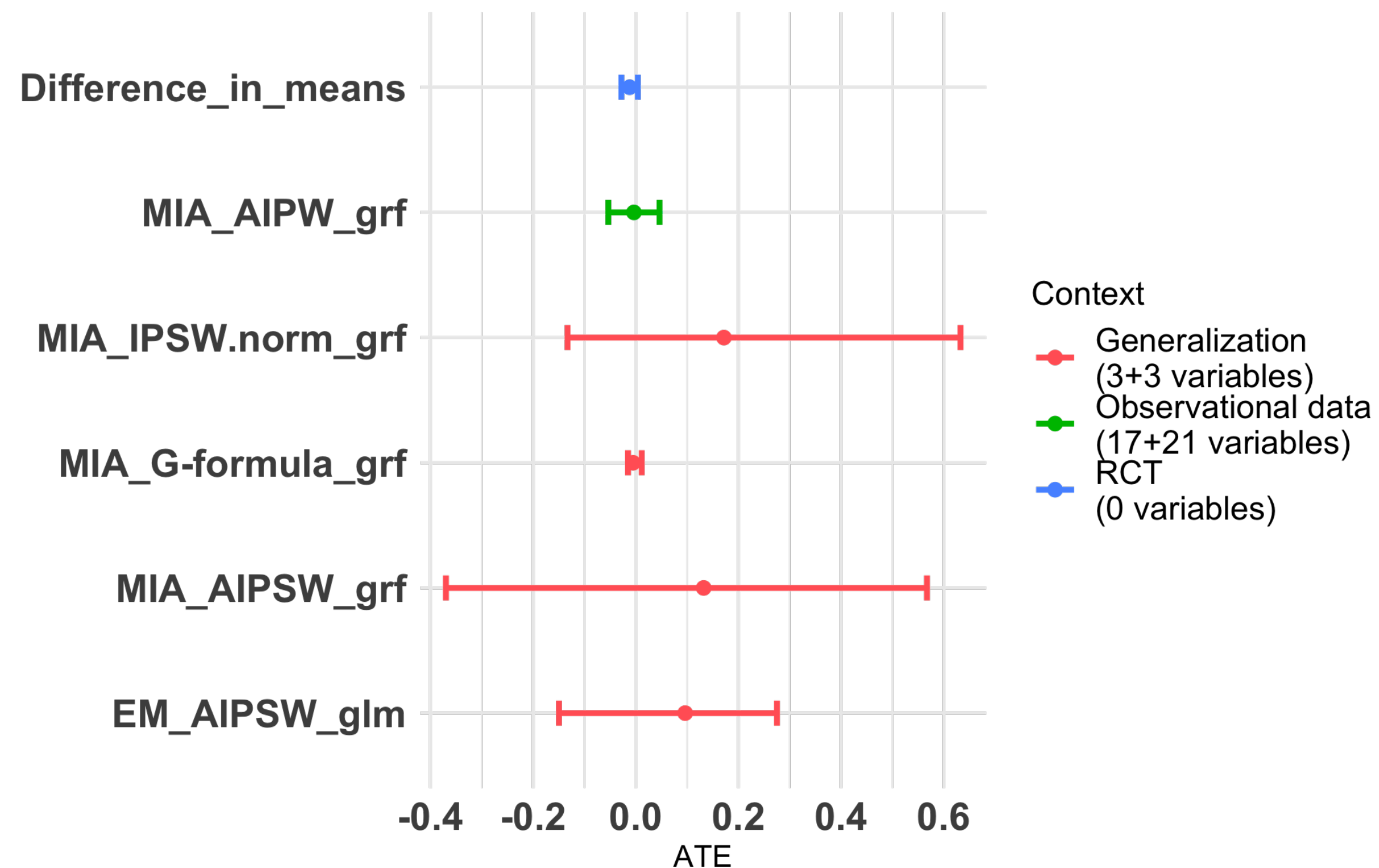
Widely varying results!



Extract of the applied results published in Statistical sciences.

Application on the CRASH-3 & Traumabase example

Widely varying results!



List of open questions

- Effect of finite sample?
- Which covariate to include? — would adding prognostic variables reduce the variance as in the classical case?
- Clinicians collaborators where rather interested in the ratio, rather than the difference

Extract of the applied results published in Statistical sciences.

Contributions

1. A review of methods to combine experimental and observational data

— *Causal inference methods for combining randomized trials and observational studies: a review*, co-authored with Imke Mayer, Statistical Science, 2022

2. Consistency proofs and sensitivity analysis for generalisation

— *Causal effect on a target population: A sensitivity analysis to handle missing covariates*, Journal of Causal Inference, 2022

3. Properties of IPWS and discussion on covariates selection

— *Reweighting the RCT for generalization: finite sample error and variable selection*, in revision in JRRS-A

4. Extension of generalization to other causal measures than the difference

— *Risk ratio, odds ratio, risk difference... Which causal measure is easier to generalize?*, submitted to Stat. In Med.

Contributions

1. A review of methods to combine experimental and observational data

— *Causal inference methods for combining randomized trials and observational studies: a review*, co-authored with Imke Mayer, Statistical Science, 2022

2. Consistency proofs and sensitivity analysis for generalisation

— *Causal effect on a target population: A sensitivity analysis to handle missing covariates*, Journal of Causal Inference, 2022

3. Properties of IPWS and discussion on covariates selection

— *Reweighting the RCT for generalization: finite sample error and variable selection*, in revision in JRRS-A

4. Extension of generalization to other causal measures than the difference

— *Risk ratio, odds ratio, risk difference... Which causal measure is easier to generalize?*, submitted to Stat. In Med.

Recalling what is done on a classical clinical randomized trial

Horvitz-Thomson estimator

$$\hat{\tau}_{HT,n} = \frac{1}{n} \sum_{i \in \text{Trial}} \left(\frac{Y_i A_i}{\pi} - \frac{Y_i (1 - A_i)}{1 - \pi} \right)$$

Probability to receive treatment, usually 0.5

Recalling what is done on a classical clinical randomized trial

Horvitz-Thomson estimator

$$\hat{\tau}_{HT,n} = \frac{1}{n} \sum_{i \in \text{Trial}} \left(\frac{Y_i A_i}{\pi} - \frac{Y_i (1 - A_i)}{1 - \pi} \right)$$

Probability to receive treatment, usually 0.5

Properties

$$\mathbb{E} [\hat{\tau}_{HT,n}] = \tau_R$$

Unbiased

$$n \text{Var} [\hat{\tau}_{HT,n}] = \frac{\mathbb{E} [(Y^{(1)})^2]}{\pi} + \frac{\mathbb{E} [(Y^{(0)})^2]}{1 - \pi} - \tau^2 := V_{HT}$$

Finite sample variance

Enriching the trial data with the target sample data

$$\hat{\tau}_{IPSW,n,m} = \frac{1}{n} \sum_{i \in \text{Trial}} \frac{\hat{p}_{T,m}(X_i)}{\hat{p}_{R,n}(X_i)} \left(\frac{Y_i A_i}{\pi} - \frac{Y_i (1 - A_i)}{1 - \pi} \right)$$

Depends on n and m !

Same as single RCT

Wished properties?

$$\mathbb{E} [\hat{\tau}_{IPSW,n}] = \tau_T$$

Unbiased

$$n \text{ Var} [\hat{\tau}_{IPSW,n,m}] = ?$$

Theoretical guarantees of IPSW with oracle weights

$$\hat{\tau}_{\pi, T, R, n}^* = \frac{1}{n} \sum_{i \in \mathcal{R}} \frac{p_T(X_i)}{p_R(X_i)} Y_i \left(\frac{A_i}{\pi} - \frac{1 - A_i}{1 - \pi} \right)$$

Theoretical guarantees of IPSW with oracle weights

$$\hat{\tau}_{\pi, T, R, n}^* = \frac{1}{n} \sum_{i \in \mathcal{R}} \frac{p_T(X_i)}{p_R(X_i)} Y_i \left(\frac{A_i}{\pi} - \frac{1 - A_i}{1 - \pi} \right)$$

Finite-sample properties — Oracle weights

$$\mathbb{E} \left[\hat{\tau}_{\pi, T, R, n}^* \right] = \tau_T \quad \text{Var} \left[\hat{\tau}_{\pi, T, R, n}^* \right] = \frac{V_o}{n}$$

where

$$V_o = \text{Var}_R \left[\frac{p_T(X)}{p_R(X)} \tau(X) \right] + \mathbb{E}_R \left[\left(\frac{p_T(X)}{p_R(X)} \right)^2 V_{HT}(X) \right]$$

How do we estimate weights in practice?

Assumption: assume \mathbf{X} is composed of categorical covariates — e.g. smoking status, gender, ...

$$\hat{\tau}_{\pi, T, n}^* = \frac{1}{n} \sum_{i \in \mathcal{R}} \frac{p_T(X_i)}{\hat{p}_{R, n}(X_i)} Y_i \left(\frac{A_i}{\pi} - \frac{1 - A_i}{1 - \pi} \right) \quad \text{where} \quad \hat{p}_{R, n}(x) := \frac{1}{n} \sum_{i \in \mathcal{R}} 1_{X_i=x}$$

How do we estimate weights in practice?

Assumption: assume \mathbf{X} is composed of categorical covariates — e.g. smoking status, gender, ...

$$\hat{\tau}_{\pi, T, n}^* = \frac{1}{n} \sum_{i \in \mathcal{R}} \frac{p_T(X_i)}{\hat{p}_{R, n}(X_i)} Y_i \left(\frac{A_i}{\pi} - \frac{1 - A_i}{1 - \pi} \right) \quad \text{where} \quad \hat{p}_{R, n}(x) := \frac{1}{n} \sum_{i \in \mathcal{R}} 1_{X_i=x}$$

Finite-sample properties — Semi oracle weights

$$\mathbb{E} \left[\hat{\tau}_{\pi, T, n}^* \right] - \tau = - \sum_{x \in \mathbb{X}} p_T(x) (1 - p_R(x))^n \tau(x)$$

$$\text{Var} \left[\hat{\tau}_{\pi, T, n}^* \right] \leq \frac{2V_{so}}{n+1} + \left(1 - \min_{x \in \mathbb{X}} p_R(x) \right)^n \mathbb{E}_T \left[\tau(X)^2 \right]$$

where $V_{so} := \mathbb{E}_R \left[\left(\frac{p_T(X)}{p_R(X)} \right)^2 V_{HT}(X) \right] = V_o - \text{Var}_R \left[\frac{p_T(X)}{p_R(X)} \tau(X) \right]$

- Positive **but exponentially small bias** compared to the oracle estimate due to undercoverage of some categories in the trial
- Smaller asymptotic variance than the oracle estimate⁽⁴⁾

(4) Robins et al. (1992). Estimating exposure effects by modelling the expectation of exposure conditional on confounders. *Biometrics*.

Theoretical guarantees of IPSW with completely estimated weights

$$\hat{\tau}_{\pi,n,m} = \frac{1}{n} \sum_{i \in \mathcal{R}} \frac{\hat{p}_{T,m}(X_i)}{\hat{p}_{R,n}(X_i)} Y_i \left(\frac{A_i}{\pi} - \frac{1 - A_i}{1 - \pi} \right)$$

Finite-sample properties — Fully estimated weights

$$\mathbb{E} \left[\hat{\tau}_{\pi,T,n}^* \right] - \tau = - \sum_{x \in \mathcal{X}} p_T(x) (1 - p_R(x))^n \tau(x)$$

$$\begin{aligned} \text{Var} \left[\hat{\tau}_{\pi,n,m} \right] &\leq \frac{2V_{so}}{n+1} + \frac{\text{Var}_T [\tau(X)]}{m} \\ &\quad + \frac{2}{m(n+1)} \mathbb{E}_R \left[\frac{p_T(X)(1-p_T(X))}{p_R(X)^2} V_{HT}(X) \right] \\ &\quad + \left(1 - \min_x p_R(x) \right)^{n/2} \mathbb{E}_T [\tau(X)^2] \left(1 + \frac{4}{m} \right) \end{aligned}$$

- Same bias as the semi oracle: bias can only be explained by a limited RCT
- Two sample size: RCT (n) and observational study (m)
 - Additional term decreasing as $1/m$ compared to the semi oracle estimate
 - Consistent if both n and $m \rightarrow \infty$. In this case, the first two terms dominate.

IPSW Large sample properties

Large sample properties — Fully estimated weights

Letting $\lim_{n,m \rightarrow \infty} m/n = \lambda \in [0, \infty]$,

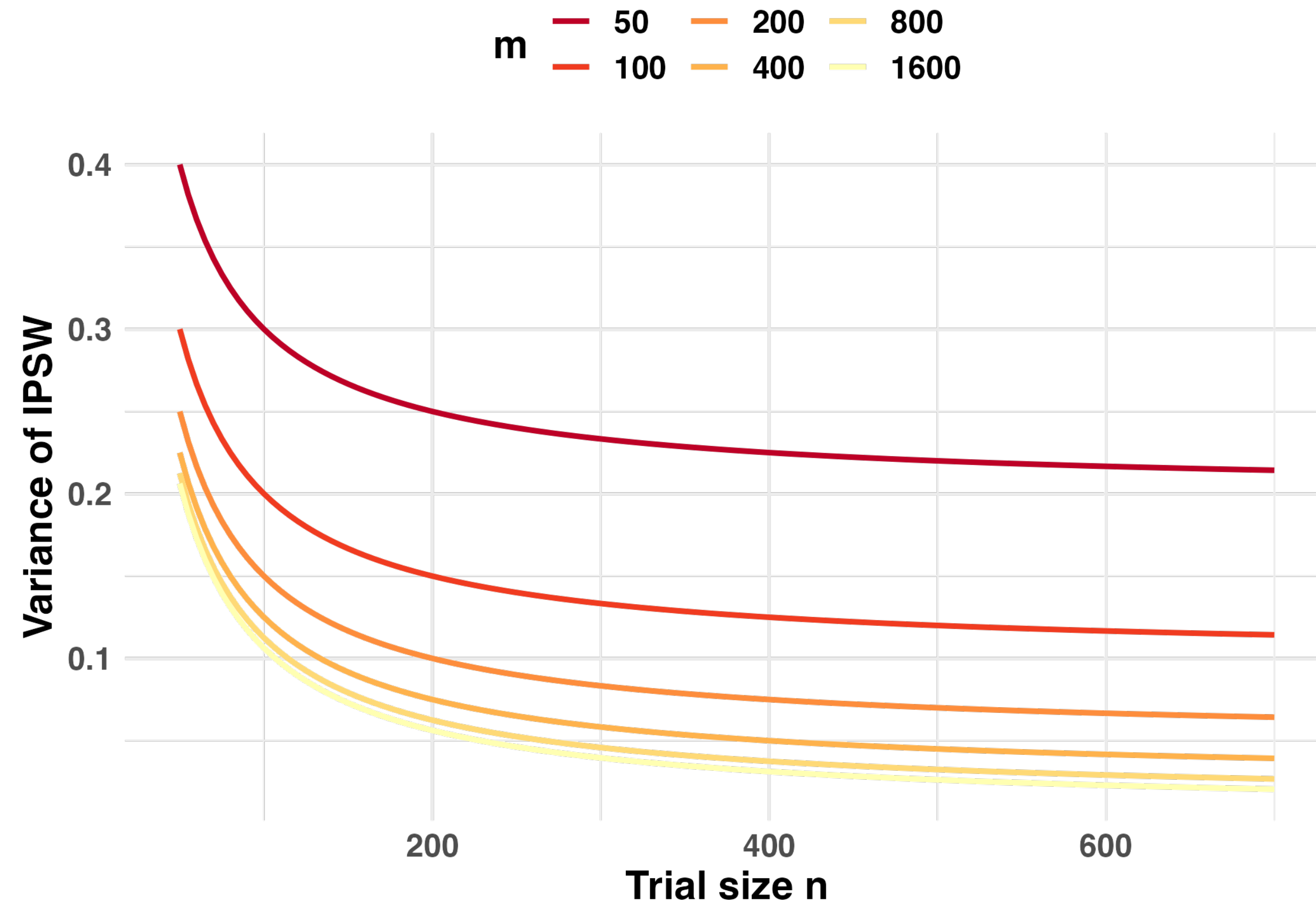
$$\lim_{n,m \rightarrow \infty} \min(n, m) \text{Var} [\hat{\tau}_{\pi, n, m}] = \min(1, \lambda) \left(\frac{\text{Var} [\tau(X)]}{\lambda} + V_{so} \right)$$

Two data samples sizes dictating two asymptotic variance

If **target** \gg **trial** (i.e. $\lambda = \infty$), asymptotic variance = semi-oracle's one and depends on the ratio of probabilities

If **target** \ll **trial** (i.e. $\lambda = 0$), asymptotic variance = conditional treatment effect variance

IPSW Large sample properties - Illustration

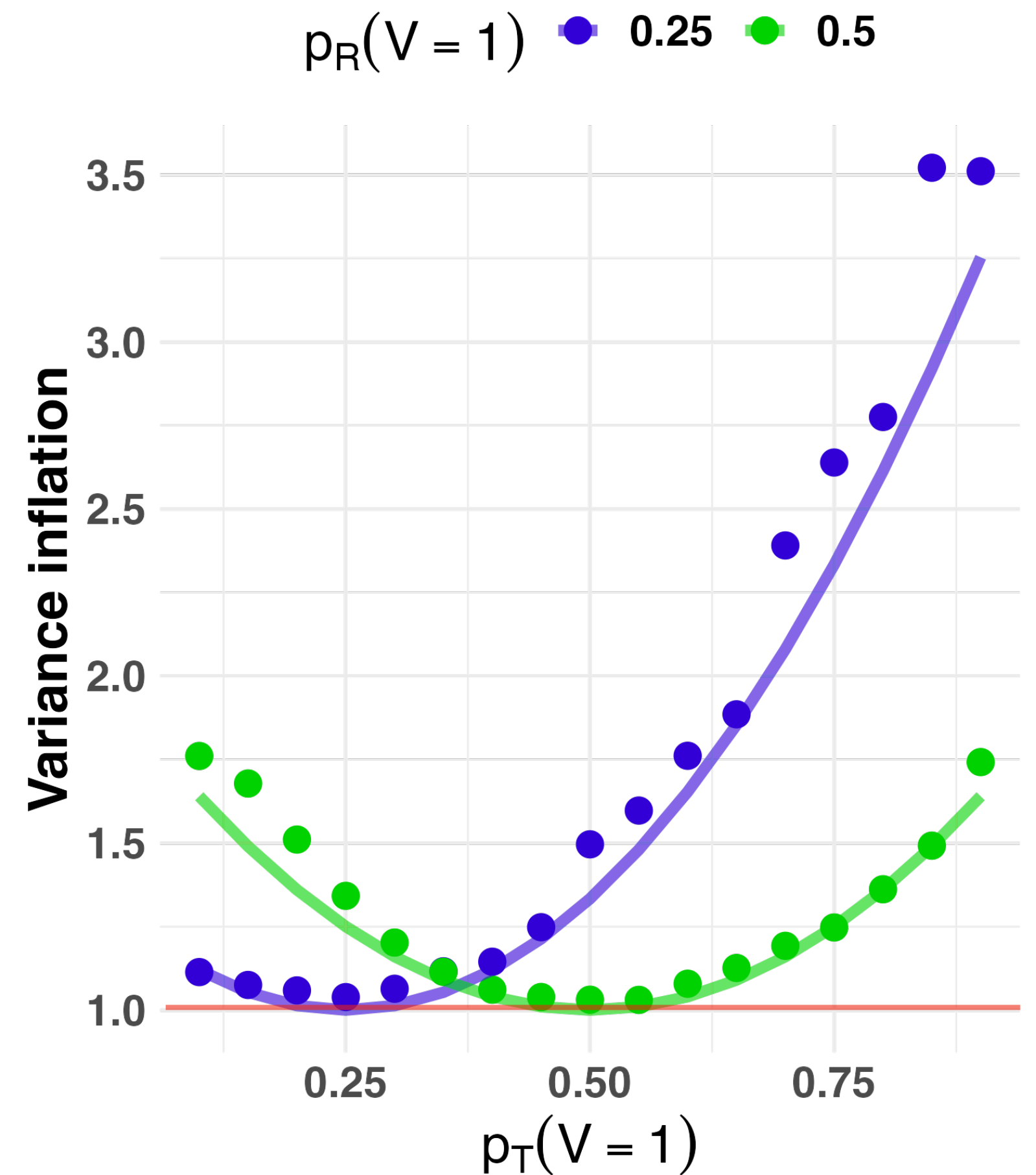


Practical recommendation

e.g. When $n = 200$ and $m = 50$, it is better to double the size of the observational data than that of the RCT.

Impact of additional covariates: for the worse?

- Covariates needed to generalize are,
 - **Treatment effect modifier**
A covariate along which the treatment effect is modulated
 - **Shifted**
Not the same proportion in each population
- In practice, one may be tempted to add many covariates
 - It does prevent to miss important ones
 - But what happen if gender is added but is only shifted?



Dots are simulations, plain lines are the theory introduced on next slide

Impact of adding a shifted covariate which is not treatment effect modifier

Non treatment effect modifier

V does not modulate treatment effect, that is

$$\forall v \in \mathbb{V}, \forall s \in \{T, R\}, \quad \mathbb{P}_s(Y^{(1)} - Y^{(0)} \mid X = x, V = v) = \mathbb{P}_s(Y^{(1)} - Y^{(0)} \mid X = x)$$

Impact of adding a shifted covariate which is not treatment effect modifier

Non treatment effect modifier

V does not modulate treatment effect modifier, that is

$$\forall v \in \mathbb{V}, \forall s \in \{T, R\}, \quad \mathbb{P}_s(Y^{(1)} - Y^{(0)} \mid X = x, V = v) = \mathbb{P}_s(Y^{(1)} - Y^{(0)} \mid X = x)$$

Shifted covariate which is not a treatment effect modifier

Consider the semi-oracle IPSW estimator and a set of additional shifted covariates V, independent of X, which are not treatment effect modifier, then

$$\lim_{n \rightarrow \infty} n \operatorname{Var}_R \left[\hat{\tau}_{T,n,m}^*(X, V) \right] = \left(\sum_{v \in \mathcal{V}} \frac{p_T(v)^2}{p_R(v)} \right) \lim_{n \rightarrow \infty} n \operatorname{Var}_R \left[\hat{\tau}_{T,n,m}^*(X) \right]$$

Including non-necessary covariates can seriously damage precision

Impact of adding a non-shifted covariate which is a treatment effect modifier

Non-shifted covariate

V is not shifted, that is

$$\forall v \in \mathbb{V}, p_T(v) = p_R(v).$$

Non-shifted covariate which is a treatment effect modifier

Consider the semi-oracle IPSW estimator and a set of additional non-shifted treatment effect modifier set V, independent of X. Then,

$$\lim_{n \rightarrow \infty} n \text{Var}_R \left[\hat{\tau}_{T,n}^*(X, V) \right] = \lim_{n \rightarrow \infty} n \text{Var}_R \left[\hat{\tau}_{T,n}^*(X) \right] - \mathbb{E}_R \left[\frac{p_T(X)}{p_R(X)} \text{Var} [\tau(X, V) | X] \right]$$

Including non-necessary covariates can improve precision

Semi-synthetic simulation

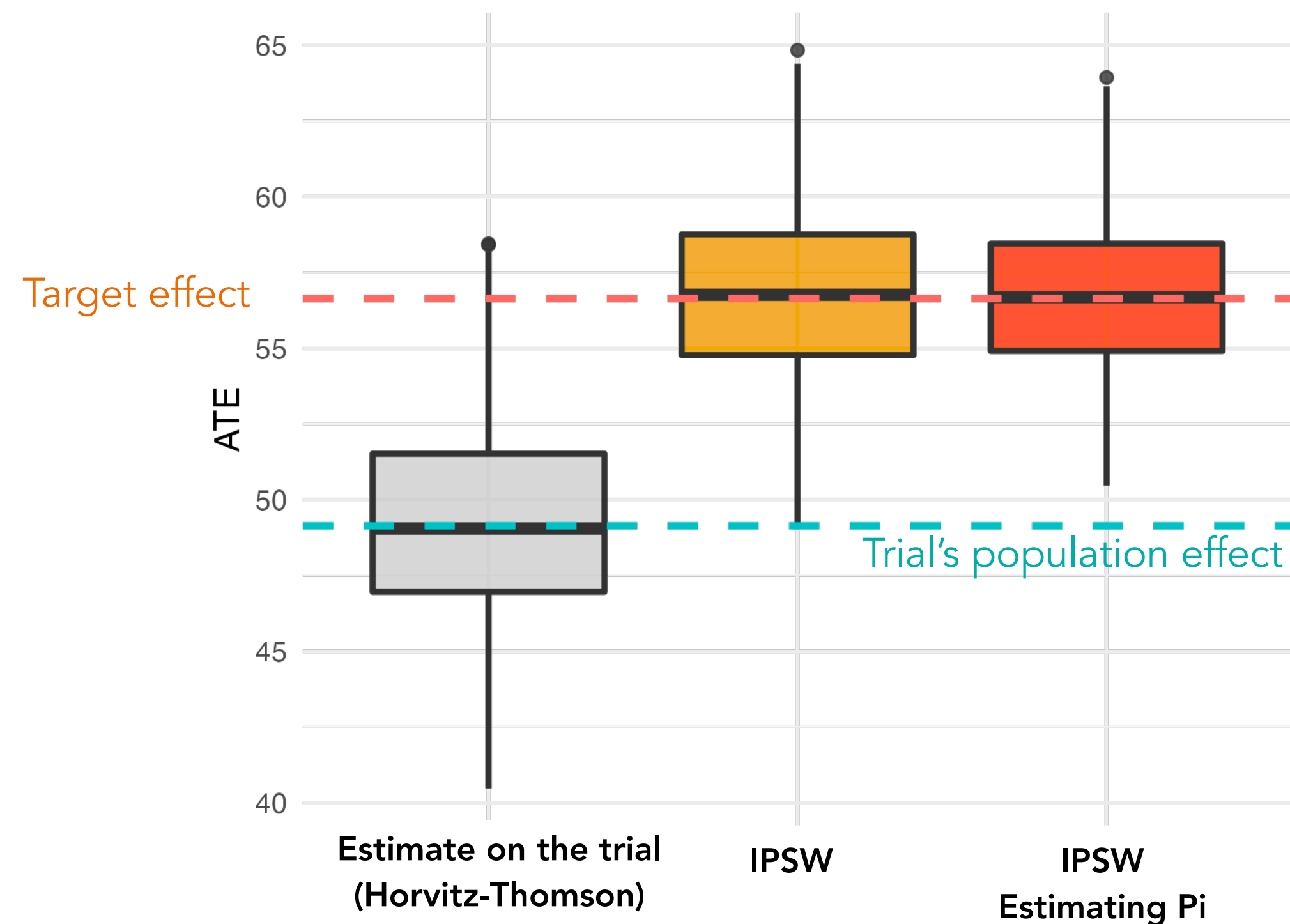
We illustrate the results on semi-synthetic simulations

- Simulations are built from CRASH-3 (~ 9,000 individuals) and Traumabase (~30,000 individuals);
- Doing so, this reflects a real-world shift;
- Covariates are : Glasgow score, gender, time-to-treatment (TTT), blood pressure;
- Time to treatment is simulated as not present in the Traumabase;
- As all covariates are shifted (even a little), a non-shifted treatment effect modifier Z is created
- The outcome is synthetic.

$$Y = 10 - \mathbf{Glasgow} + (\mathbf{if\ Girl:} - 5 \mathbf{else:}0) + A (15(6 - \mathbf{TTT}) + 3 * (\mathbf{Blood.pressure} - 1)^2 + 50Z) + \varepsilon_{TTT}$$

↓
Random gaussian noise whose variance depends on the value of TTT

Results from the semi-synthetic simulations (1)



This simulation does not include Z as the focus is not on adding non-useful covariates

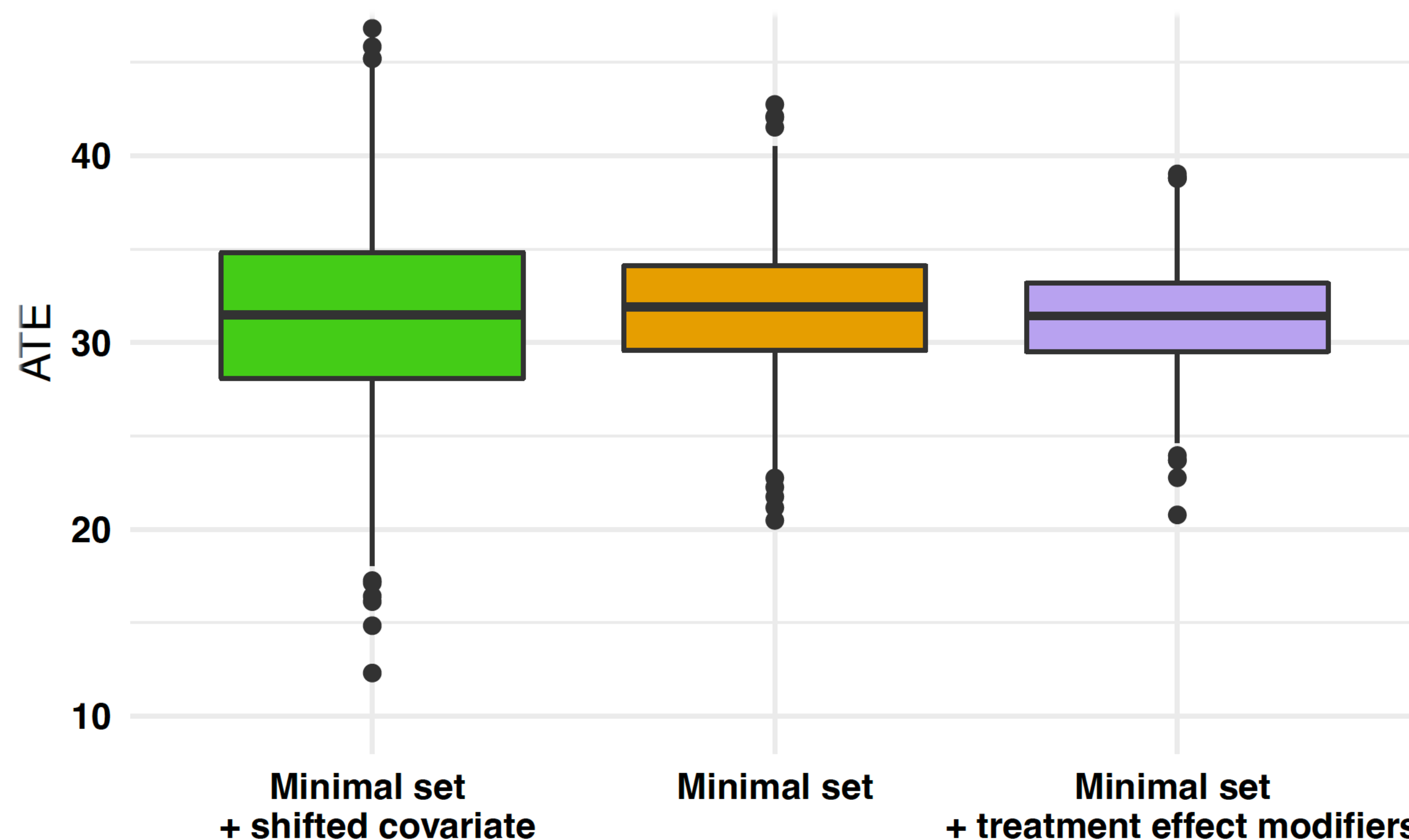
1. Re-weighting allows to recover the target effect
2. Two additional theoretical results not detailed above
 - Reducing variance when estimating the probability to be treated in the trial P_i ,
 - Re-weighted trial has not necessarily a larger variance.

Results from the semi-synthetic simulations (2)

Effect of non-necessary covariates on the variance

IPSW with $n = 3000$ and $m = 10000$ and 1,000 repetitions

- The addition of the covariate **GCS** increases the variance,
- while the addition of a **non-shifted treatment effect modifier** leads to an improvement in variance.



This simulation includes Z as the focus is on adding non-useful covariates

**Risk ratio, odds
ratio, risk
difference**

**Which causal measure is
easier to generalize?**



Contributions

1. A review of methods to combine experimental and observational data

— *Causal inference methods for combining randomized trials and observational studies: a review*, co-authored with Imke Mayer, Statistical Science, 2022

2. Consistency proofs and sensitivity analysis for generalisation

— *Causal effect on a target population: A sensitivity analysis to handle missing covariates*, Journal of Causal Inference, 2022

3. Properties of IPWS and discussion on covariates selection

— *Reweighting the RCT for generalization: finite sample error and variable selection*, in revision in JRRS-A

4. Extension of generalization to other causal measures than the difference

— *Risk ratio, odds ratio, risk difference... Which causal measure is easier to generalize?*, submitted to Stat. In Med.

Illustrative example

RCT from Cook and Sackett (1995)

- **Y** the observed binary outcome
- **A** binary treatment assignment
- **X** baseline covariates

Stroke after 5 years

11.1% Control — vs — **6.7 % Treated**

Usually referring to an *effect*, is related to how one *contrasts* those two

e.g. Ratio = $6.7/11.1 = 0.6$ or Diff = -0.04

Illustrative example

RCT from Cook and Sackett (1995)

- Y the observed binary outcome
- A binary treatment assignment
- X baseline covariates

Stroke after 5 years

11.1% Control — vs — 6.7 % Treated

— A variety of causal measures exist

Potential outcomes framework $\begin{matrix} \nearrow \mathbb{E}[Y^{(0)}] \\ \searrow \mathbb{E}[Y^{(1)}] \end{matrix}$

<p>Count the stroke</p> $\tau_{RR} = \frac{\mathbb{E}[Y^{(1)}]}{\mathbb{E}[Y^{(0)}]}$	<p>Count the non-stroke</p> $\tau_{SR} = \frac{1 - \mathbb{E}[Y^{(1)}]}{1 - \mathbb{E}[Y^{(0)}]}$
<p>Risk Difference</p> $\tau_{RD} = \mathbb{E}[Y^{(1)}] - \mathbb{E}[Y^{(0)}]$	<p>Number Needed to Treat</p> $\tau_{NNT} = \tau_{RD}^{-1}$
<p>Odds Ratio</p> $\tau_{OR} = \frac{\mathbb{E}[Y^{(1)}]}{1 - \mathbb{E}[Y^{(1)}]} \left(\frac{1 - \mathbb{E}[Y^{(0)}]}{1 - \mathbb{E}[Y^{(0)}]} \right)^{-1}$	

Note that for binary Y , $\mathbb{E}[Y(a)] = P(Y(a)=1)$

Computing all the measures on the illustrative clinical example

	τ_{RD}	τ_{RR}	τ_{SR}	τ_{NNT}	τ_{OR}
All (P_s)	-0.0452	0.6	1.05	22	0.57
X = 1	-0.006	0.6	1.01	167	0.6
X = 0	-0.08	0.6	1.1	13	0.545

$X = 1 \leftrightarrow$ low baseline risk

↕ Marginal effects τ
↕ Conditional effects $\tau(x)$

Computed from Cook & Sackett (1995)

 "Treated group has 0.6 times the risk of having a stroke outcome when compared with the placebo." **or** "The Number Needed to Treat is 22."

— leads to different impressions and heterogeneity patterns

The age-old question of how to report effects



Source: Wikipedia

“ We wish to decide whether we shall count the failures or the successes and whether we shall make relative or absolute comparisons”

— *Mindel C. Sheps, New England Journal of Medicine, in 1958*

The choice of the measure is still actively discussed

e.g. Spiegelman and VanderWeele, 2017; Baker and Jackson, 2018; Feng et al., 2019; Doi et al., 2022; Xiao et al., 2021, 2022; Huitfeldt et al., 2021; Lapointe-Shaw et al., 2022; Liu et al., 2022 ...

— **CONSORT** guidelines recommend to report all of them

A desirable property: collapsibility

i.e. population's effect is equal to a weighted sum of local effects



Discussed in Greenland, 1987; Hernàn et al. 2011; Huitfeldt et al., 2019; Daniel et al., 2020; Didelez and Stensrud, 2022 and many others.

A desirable property: collapsibility

i.e. population's effect is equal to a weighted sum of local effects



 Discussed in Greenland, 1987; Hernàn et al. 2011; Huitfeldt et al., 2019; Daniel et al., 2020; Didelez and Stensrud, 2022 and many others.

A very famous example: the Simpson paradox

(a) Overall population, $\tau_{OR} \approx 0.26$

	Y=0	Y=1
A=1	1005	95
A=0	1074	26

(b) $\tau_{OR|F=1} \approx 0.167$ and $\tau_{OR|F=0} \approx 0.166$

F=1	Y=0	Y=1	F=0	Y=0	Y=1
A=1	40	60	A=1	965	35
A=0	80	20	A=0	994	6

Marginal effect bigger than subgroups' effects

Toy example inspired from Greenland (1987).

— Unfortunately, not all measures are collapsible

Collapsibility and formalism

- Different definitions of collapsibility in the literature

Collapsibility and formalism

- Different definitions of collapsibility in the literature
- We propose three definitions encompassing previous works

1. Direct collapsibility $\mathbb{E} [\tau(X)] = \tau$

Collapsibility and formalism

- Different definitions of collapsibility in the literature
- We propose three definitions encompassing previous works

1. Direct collapsibility $\mathbb{E} [\tau(X)] = \tau$

2. Collapsibility $\mathbb{E} [w(X, P(X, Y^{(0)})) \tau(X)] = \tau$, **with** $w \geq 0$, **and** $\mathbb{E} [w(X, P(X, Y^{(0)}))] = 1$

e.g RR is collapsible, with

$$\mathbb{E} \left[\tau_{RR}(X) \frac{\mathbb{E} [Y^{(0)} | X]}{\mathbb{E} [Y^{(0)}]} \right] = \tau_{RR}$$

Collapsibility and formalism

- Different definitions of collapsibility in the literature
- We propose three definitions encompassing previous works
 1. Direct collapsibility $\mathbb{E} [\tau(X)] = \tau$
 2. Collapsibility $\mathbb{E} [w(X, P(X, Y^{(0)})) \tau(X)] = \tau$, **with** $w \geq 0$, **and** $\mathbb{E} [w(X, P(X, Y^{(0)}))] = 1$
 3. Logic-respecting $\tau \in \left[\min_x (\tau(x)), \max_x (\tau(x)) \right]$

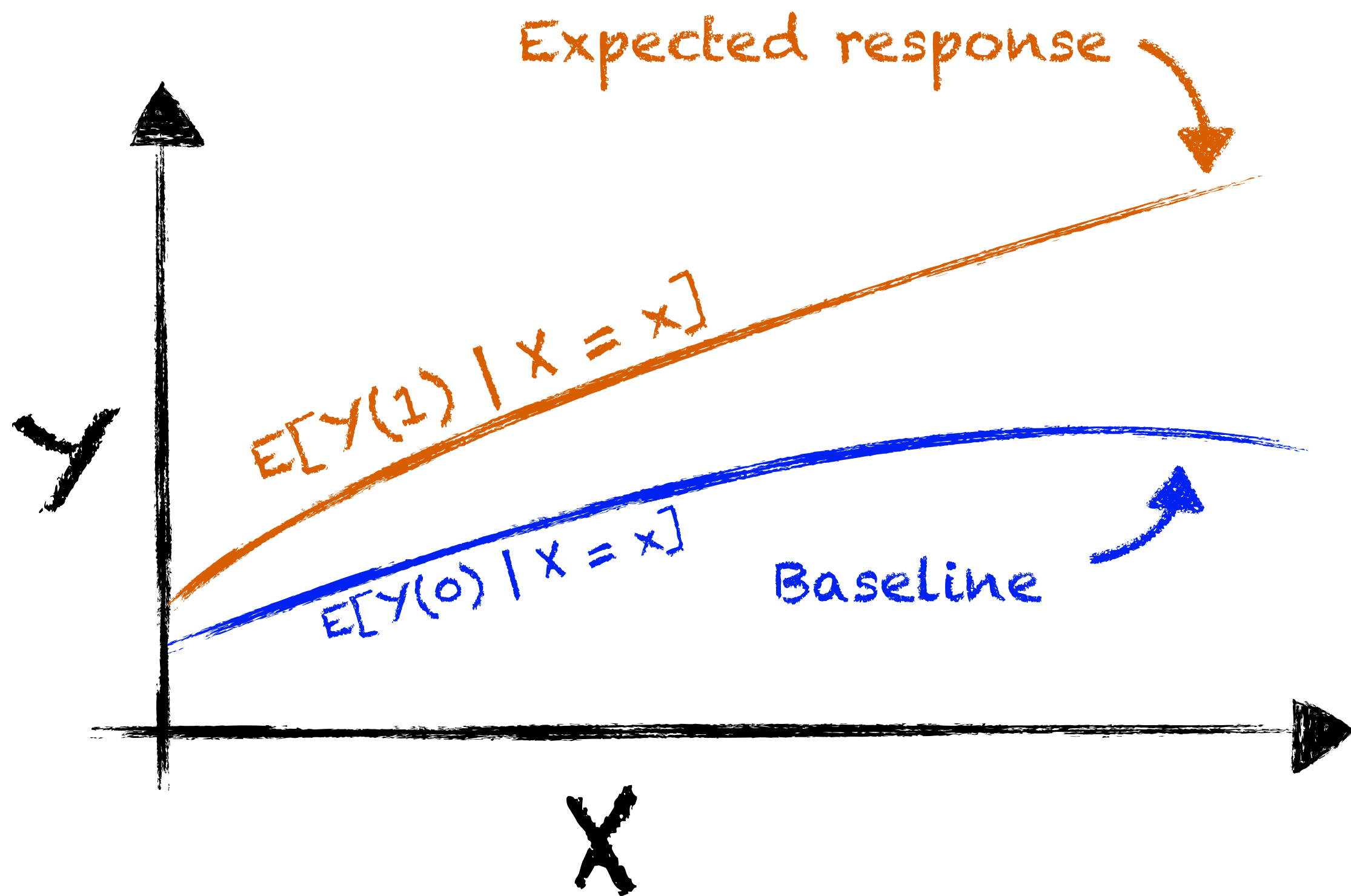
Collapsibility and formalism

- Different definitions of collapsibility in the literature
- We propose three definitions encompassing previous works
 1. Direct collapsibility $\mathbb{E} [\tau(X)] = \tau$
 2. Collapsibility $\mathbb{E} [w(X, P(X, Y^{(0)})) \tau(X)] = \tau$, **with** $w \geq 0$, **and** $\mathbb{E} [w(X, P(X, Y^{(0)}))] = 1$
 3. Logic-respecting $\tau \in \left[\min_x (\tau(x)), \max_x (\tau(x)) \right]$

Measure	Collapsible	Logic-respecting
Risk Difference (RD)	Yes	Yes
Number Needed to Treat (NNT)	No	Yes
Risk Ratio (RR)	Yes	Yes
Survival Ratio (SR)	Yes	Yes
Odds Ratio (OR)	No	No

Through the lens of non parametric generative models

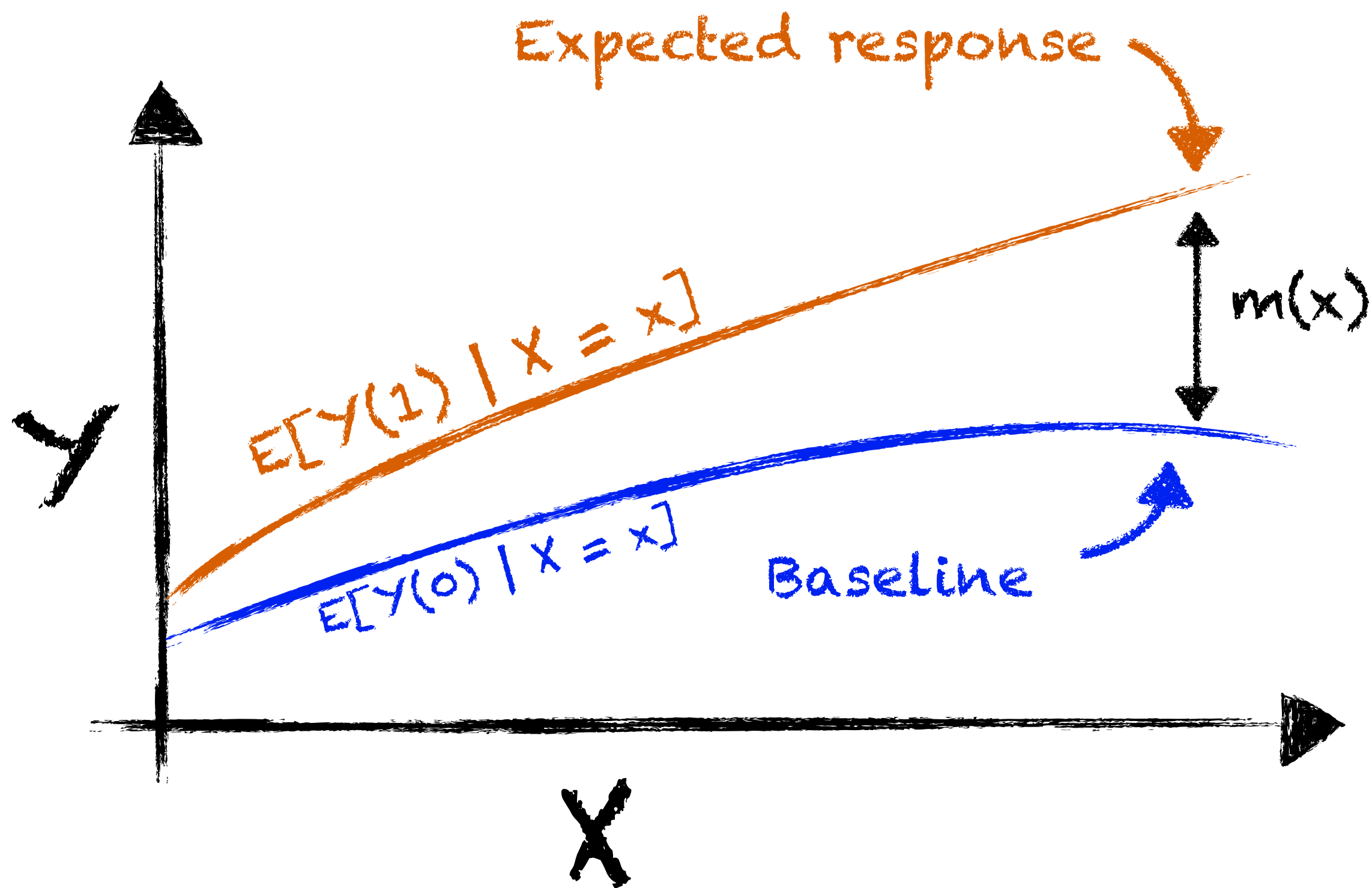
For Y continuous,



(*) This only assumes that conditional expected responses are defined for every x

Through the lens of non parametric generative models

For Y continuous,



Lemma*

There exist two functions $b(\cdot)$ and $m(\cdot)$ such that,

$$\mathbb{E} [Y^{(a)} | X] = b(X) + a m(X)$$

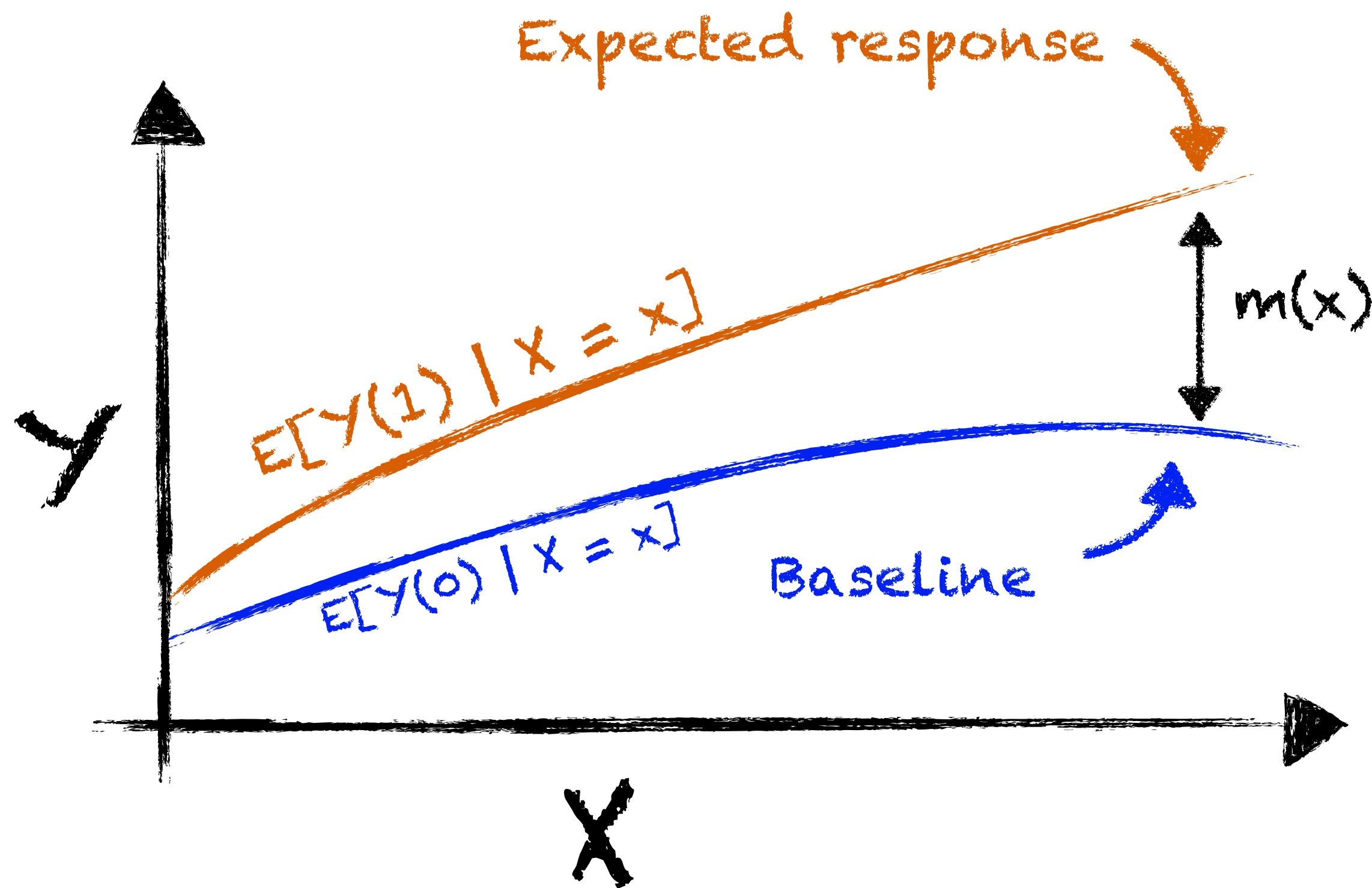
Additivity

Spirit of Robinson's decomposition (1988), further developed in Nie et al. 2020

(*) This only assumes that conditional expected responses are defined for every x

Through the lens of non parametric generative models

For Y continuous,



(*) This only assumes that conditional expected responses are defined for every x

Lemma*

There exist two functions $b(\cdot)$ and $m(\cdot)$ such that,

$$\mathbb{E} [Y^{(a)} | X] = b(X) + a m(X)$$

Additivity

Spirit of Robinson's decomposition (1988), further developed in Nie et al. 2020

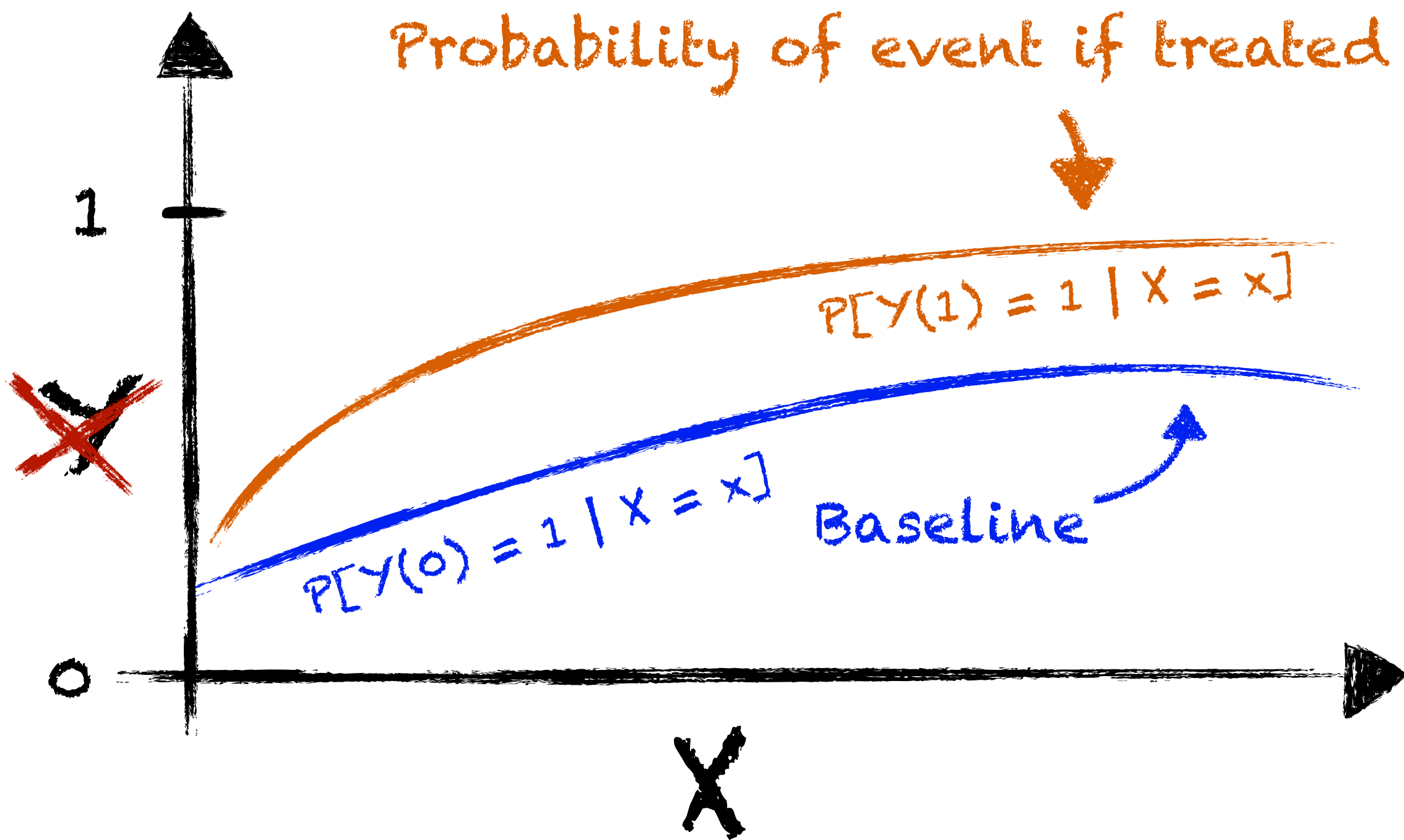
Linking generative functions with measures

$$\tau_{RR}(x) = 1 + m(x)/b(x) \quad \text{Entanglement}$$

$$\tau_{RD}(x) = m(x) \quad \text{No entanglement}$$

Through the lens of non parametric generative models

For Y binary,



~~Lemma~~

~~There exist two functions $b(\cdot)$ and $m(\cdot)$ such that,~~

$$\mathbb{E}[Y^{(a)} | X] = b(X) + a m(X)$$

~~Additivity~~

Adapted Lemma

There exist two functions $b(\cdot)$ and $m(\cdot)$ such that,

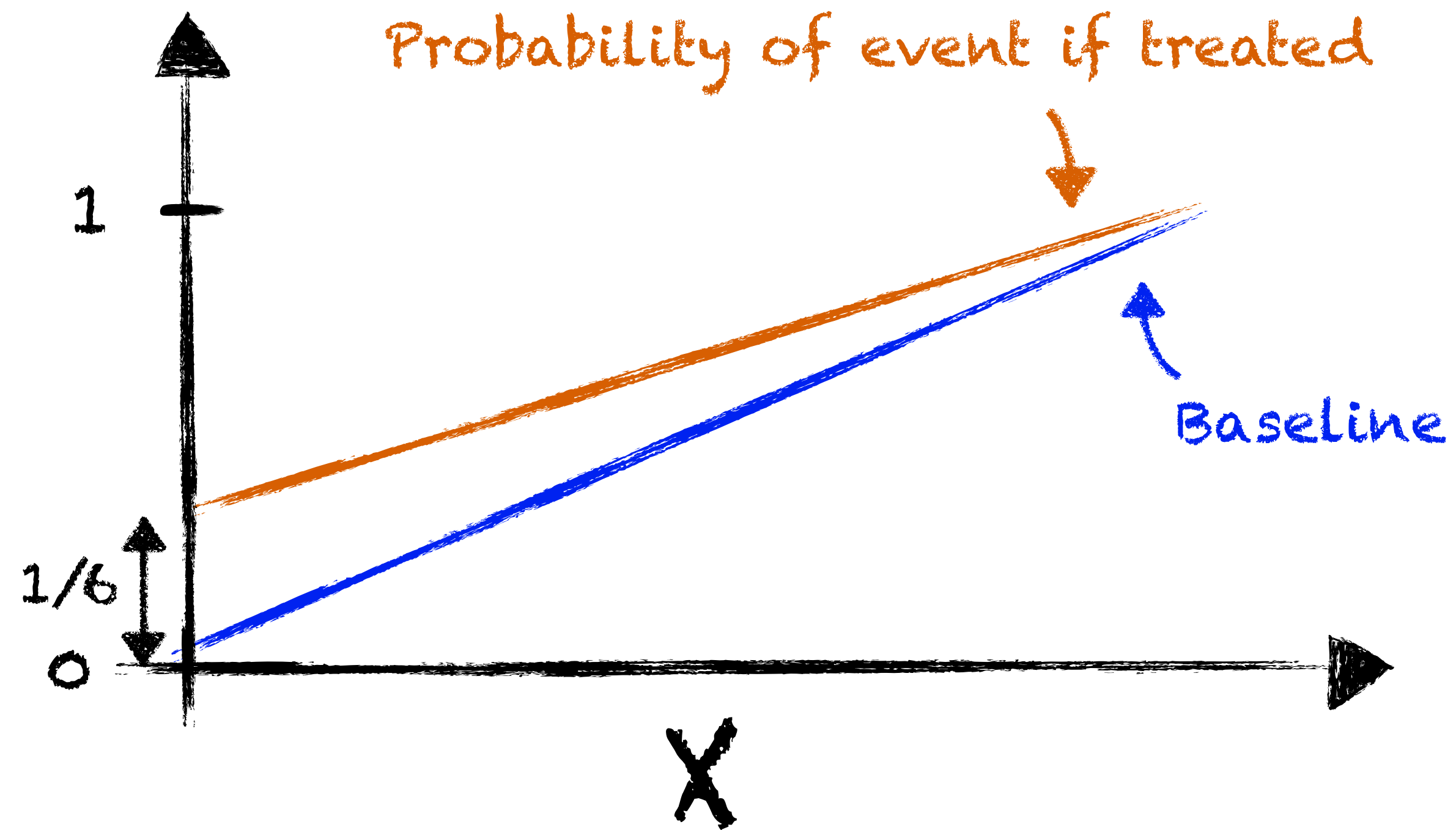
$$\ln \left(\frac{\mathbb{P}(Y^{(a)} = 1 | X)}{\mathbb{P}(Y^{(a)} = 0 | X)} \right) = b(X) + a m(X)$$

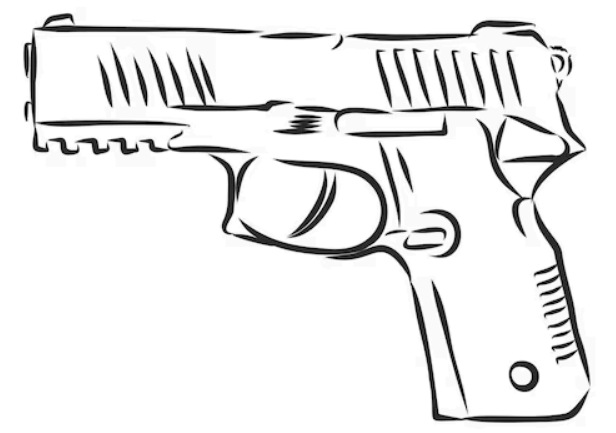
Harmful



The example of the Russian roulette

For Y binary,





The example of the Russian roulette

For Y binary,

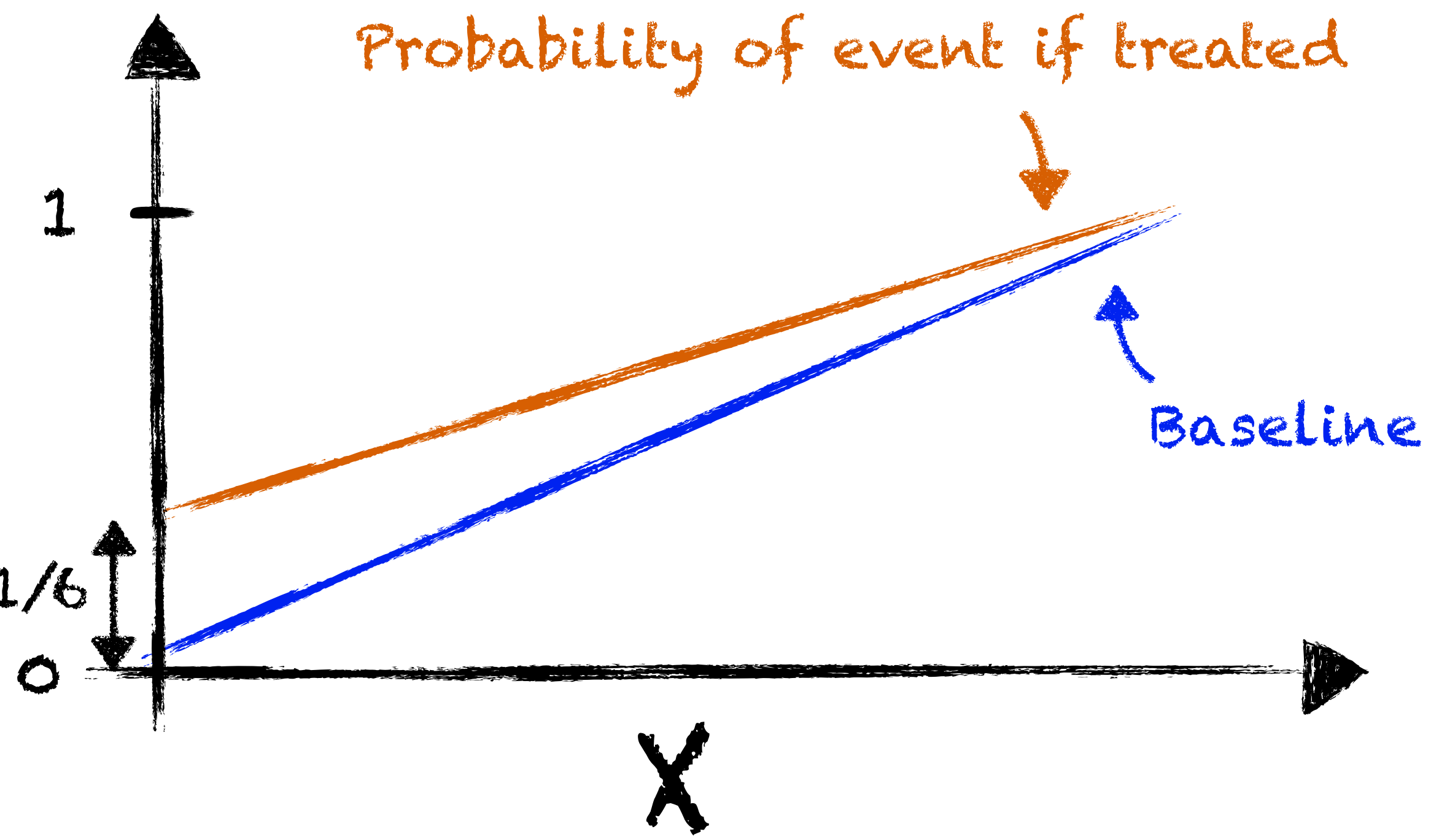
Lemma

There exist two functions $b(\cdot)$ and $m(\cdot)$ such that,

$$\mathbb{P} [Y^{(a)} = 1 | X] = b(X) + a (1 - b(X)) \frac{1}{6}$$

Simple additivity is not possible anymore

Probability of event if treated



The example of the Russian roulette



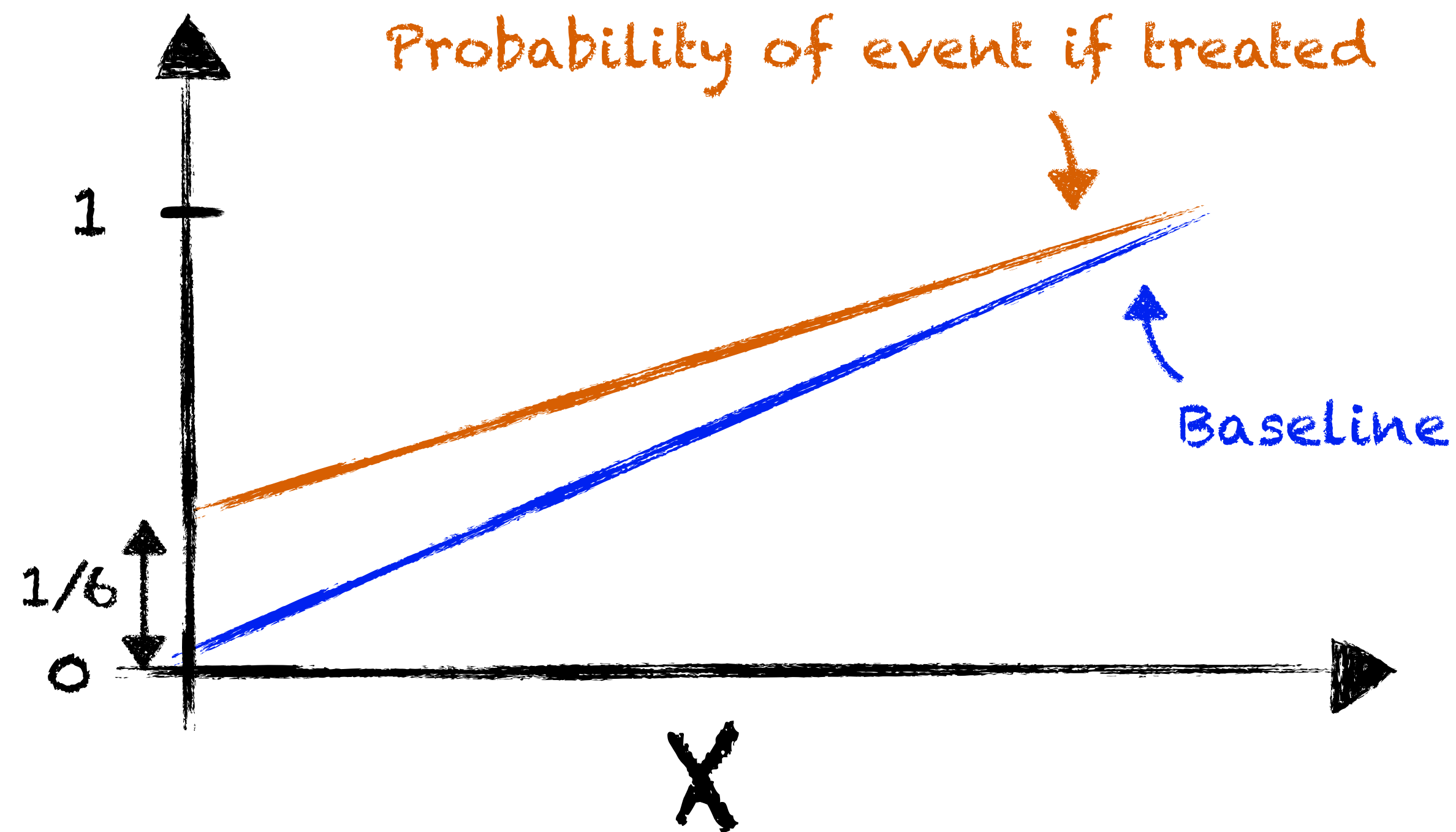
For Y binary,

Lemma

There exist two functions $b(\cdot)$ and $m(\cdot)$ such that,

$$\mathbb{P} [Y^{(a)} = 1 \mid X] = b(X) + a (1 - b(X)) \frac{1}{6}$$

Simple additivity is not possible anymore



Linking generative functions with measures

$$\tau_{RD}(x) = (1 - b(x)) \frac{1}{6} \quad \text{Entanglement}$$

$$\tau_{SR}(x) = 1 - \frac{1}{6} \quad \text{No entanglement}$$

The example of the Russian roulette

Harmful



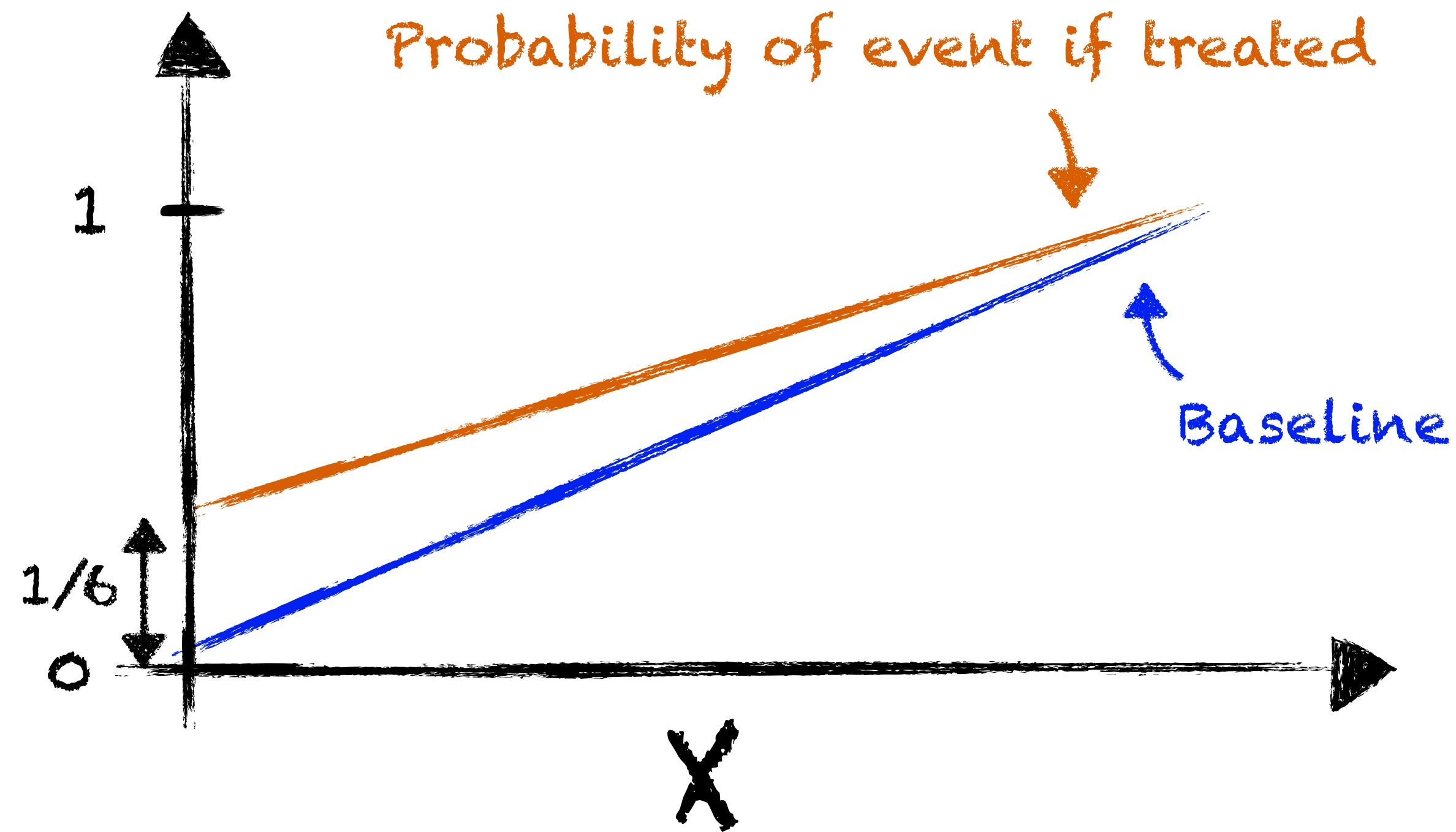
For Y binary,

Lemma

There exist two functions $b(\cdot)$ and $m(\cdot)$ such that,

$$\mathbb{P} [Y^{(a)} = 1 \mid X] = b(X) + a (1 - b(X)) \underline{m(X)}$$

Simple additivity is not possible anymore



Linking generative functions with measures

$$\tau_{RD}(x) = (1 - b(x)) \underline{m(x)} \quad \text{Entanglement}$$

$$\tau_{SR}(x) = 1 - \underline{m(x)} \quad \text{No entanglement}$$

Extension to all effect types (harmful and beneficial)

Considering a binary outcome, assume that

$$\forall x \in \mathbb{X}, \forall a \in \{0,1\}, \quad 0 < p_a(x) < 1, \quad \text{where } p_a(x) := \mathbb{P} [Y^{(a)} = 1 \mid X = x]$$

 Assumptions

Introducing,

$$m_g(x) := \mathbb{P} [Y^{(1)} = 0 \mid Y^{(0)} = 1, X = x] \quad \text{and} \quad m_b(x) := \mathbb{P} [Y^{(1)} = 1 \mid Y^{(0)} = 0, X = x],$$

Extension to all effect types (harmful and beneficial)

Considering a binary outcome, assume that

$$\forall x \in \mathbb{X}, \forall a \in \{0,1\}, \quad 0 < p_a(x) < 1, \quad \text{where } p_a(x) := \mathbb{P} [Y^{(a)} = 1 \mid X = x] \quad \begin{array}{c} \updownarrow \\ \text{Assumptions} \end{array}$$

Introducing,

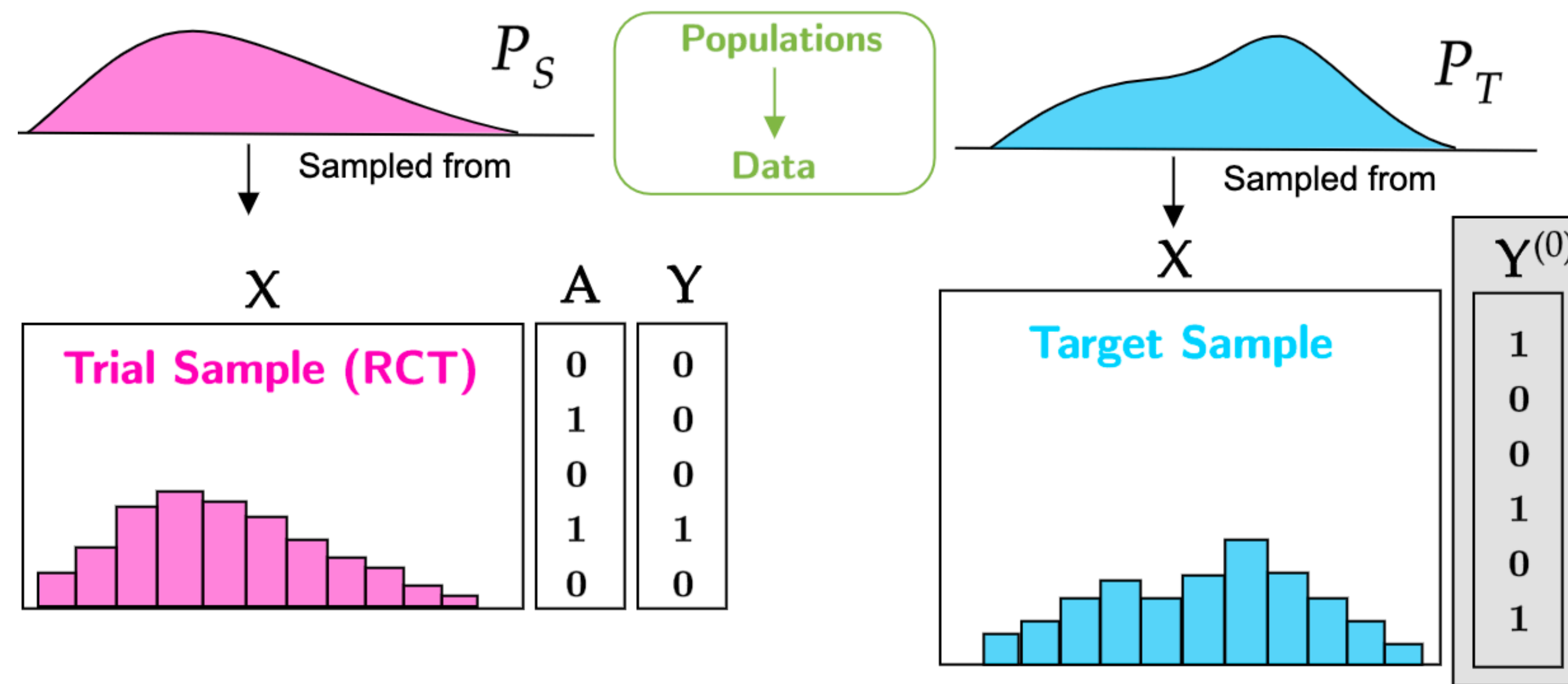
$$m_g(x) := \mathbb{P} [Y^{(1)} = 0 \mid Y^{(0)} = 1, X = x] \quad \text{and} \quad m_b(x) := \mathbb{P} [Y^{(1)} = 1 \mid Y^{(0)} = 0, X = x],$$

allows to have,

$$\mathbb{P} [Y^{(a)} = 1 \mid X = x] = b(x) + a \left(\underbrace{(1 - b(x)) m_b(x)}_{\substack{\uparrow \\ \text{More events}}} } - \underbrace{b(x) m_g(x)}_{\substack{\downarrow \\ \text{Less events}}} \right), \quad \text{where } b(x) := p_0(x).$$

Back to generalizability

Remember: we want to transport trial findings to a target population, using the trial data and a sample of the target population



We consider set-ups where control outcome is observed or not ↷

Two methods, two assumptions

S is the indicator of population's membership

Generalizing	Conditional potential outcomes	Local effects
Assumptions for RD	$\{Y^{(0)}, Y^{(1)}\} \perp\!\!\!\perp \underline{S} X$	$Y^{(1)} - Y^{(0)} \perp\!\!\!\perp \underline{S} X$
Unformal	All shifted prognostic covariates	All shifted <u>treatment effect modifiers</u> <i>Less covariates if homogeneity</i>
Identification		

Two methods, two assumptions

S is the indicator of population's membership

Generalizing	Conditional potential outcomes	Local effects
Assumptions for RD	$\{Y^{(0)}, Y^{(1)}\} \perp\!\!\!\perp \underline{S} X$	$Y^{(1)} - Y^{(0)} \perp\!\!\!\perp \underline{S} X$
Unformal	All shifted prognostic covariates	All shifted <u>treatment effect modifiers</u> <i>Less covariates if homogeneity</i>
Identification	$\mathbb{E}^T [Y^{(a)}] = \mathbb{E}^T \left[\mathbb{E}^R [Y^{(a)} X] \right]$	$\tau^T = \mathbb{E} \left[w(X, \boxed{Y^{(0)}}) \tau^R(X) \right]$ <i>Possible only if collapsible!</i>

— Depending on the assumptions, either conditional outcome or local treatment effect can be generalised

Generalizing local effect, the example of a binary Y and a beneficial effect

i.e. reducing number of events

Estimate using trial sample

$$\mathbb{E} \left[\frac{\tau_{RR}(X) \mathbb{E} [Y^{(0)} | X]}{\mathbb{E} [Y^{(0)}]} \right] = \tau_{RR}$$

Estimate using target sample

$$\tau_{RR}(x) = 1 - m_g(x)$$

Conditional RR only vary with the shifted treatment effect modulators

⚠ We need to have access to $Y(0)$!

A toy simulation

Introducing heterogeneities in the Russian roulette

- Probability to die varies
 - Stressed people can die from a heart attack
 - Executioner more merciful when facing women

A toy simulation

Introducing heterogeneities in the Russian roulette

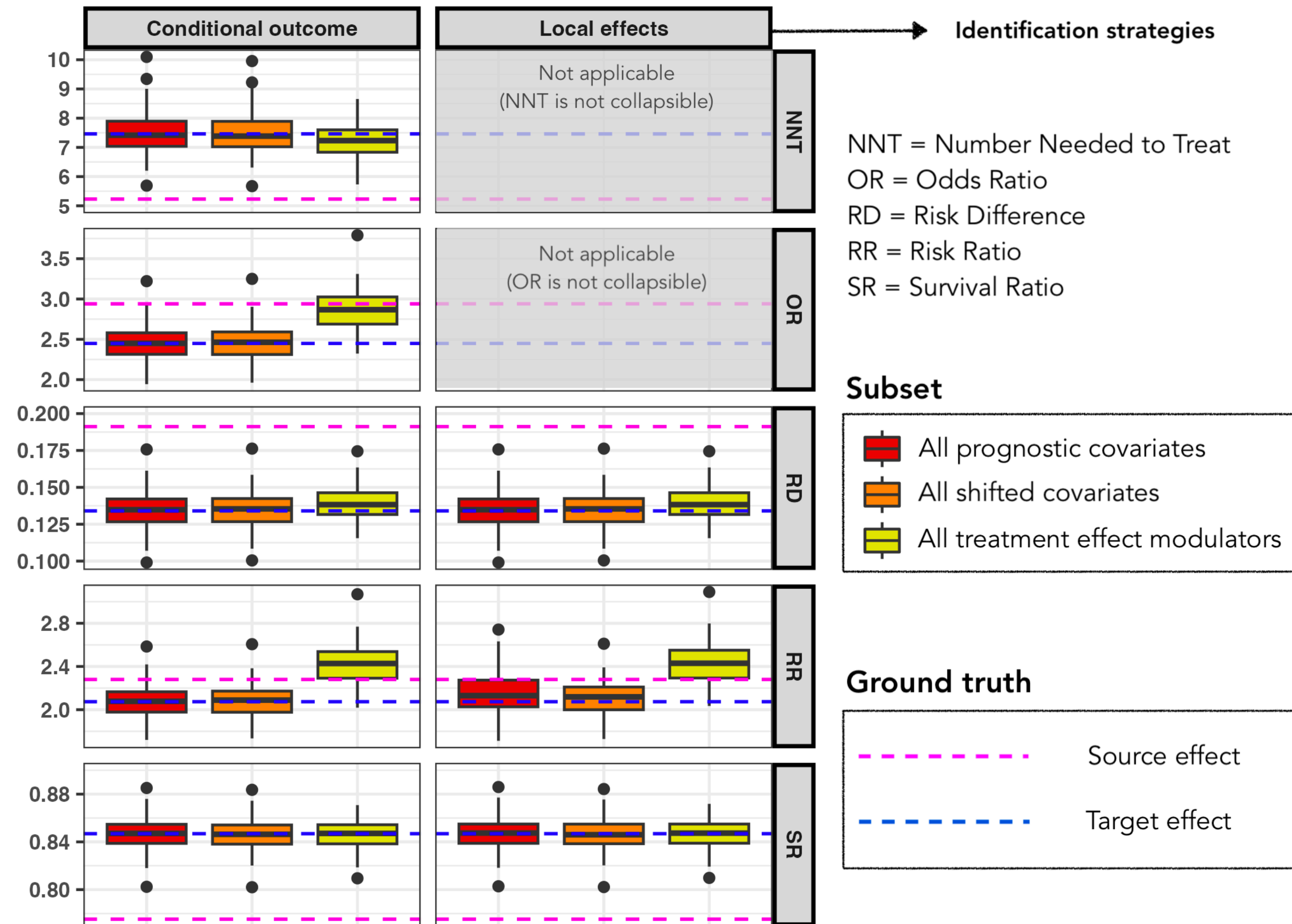
- Probability to die varies
 - Stressed people can die from a heart attack
 - Executioner more merciful when facing women

$$P[Y = 1 | X] = b(X_{1 \rightarrow 3}) + (1 - b(X_{1 \rightarrow 3})) m(X_{2 \rightarrow 3})$$

X_1 : Lifestyle general level

X_2 : stress

X_3 : gender (not shifted)



— Local SR can be generalised using only stress. All others measures requires lifestyle and stress.

Contributions

All started from motivating example from **critical care** and **two data samples** with CRASH-3 & Traumabase

This led us to tackle a the broader scope : trial's findings generalisation.

We realised **from application** that many challenges remain: missing covariates, covariate selection, consistency, impact of the causal measures, etc.

Our contribution is **to provide theoretical and methodological results** to strengthen the practice:

- Consistency proofs
- Sensitivity analysis
- Finite and large sample results of IPSW
- Characterisation of the impact of adding non necessary covariates on precision
- Impact of the causal measure on transported treatment effect identification

Future work

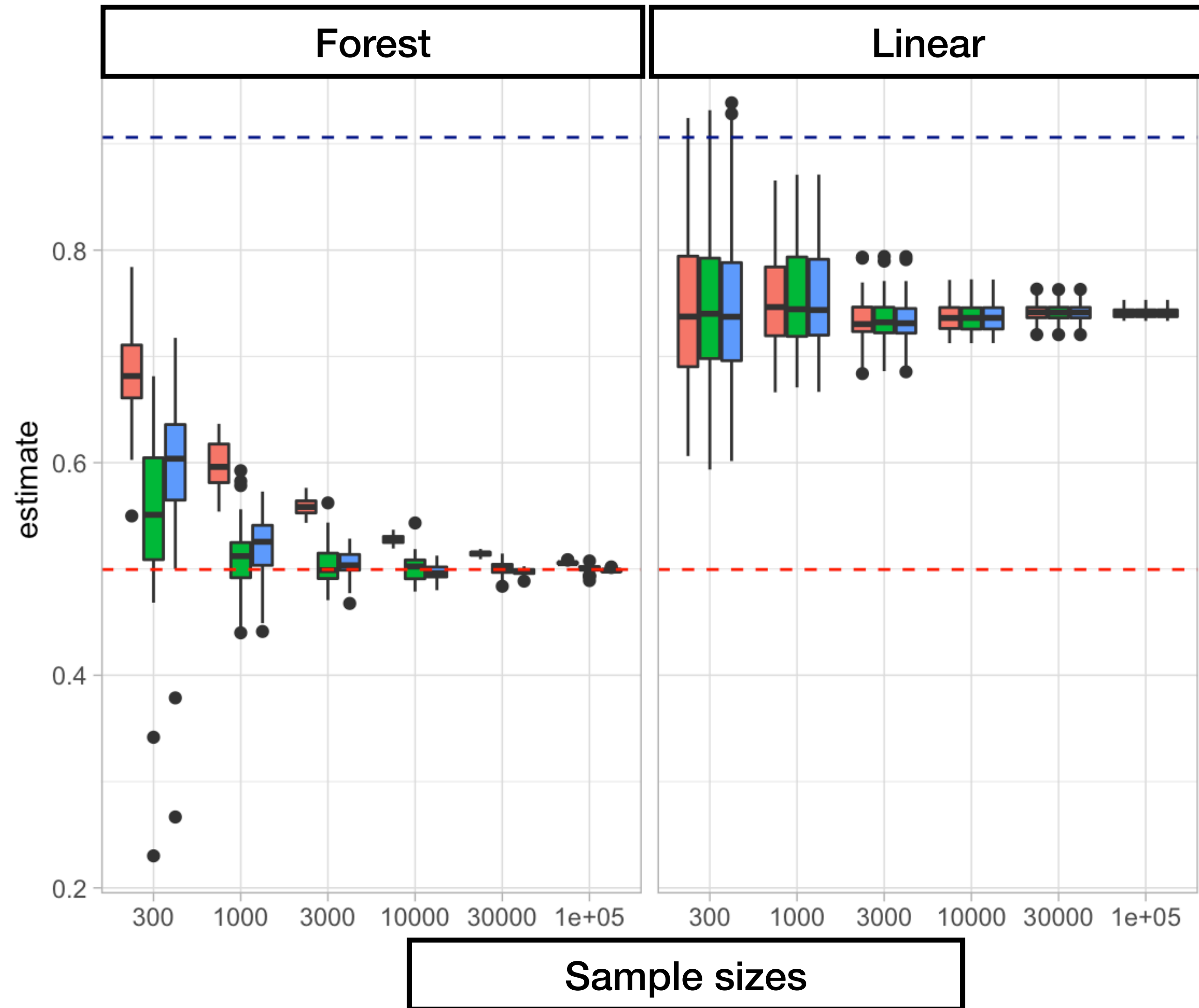
This work opens **new research questions**

- Extension of the theoretical finite and/or large sample results for
 - ▶ G-formula and AIPSW, and not only IPSW,
 - ▶ In a context where covariates are not categorical,
 - ▶ When the ratio is targeted
 - ◆ using local effect or
 - ◆ conditional outcomes re-weighting.
- Confront model with empirical data
 - ▶ Is the assumption of a completely beneficial or harmful effect valid in practice?
 - ▶ Using meta-analysis or different trials, investigate which causal measure is more or less dependent on the baseline level.

Why focusing on finite sample results? (1)

- (1) Usual sample sizes in medicine remains small
- (2) Results from simulations warned me and raised my interest

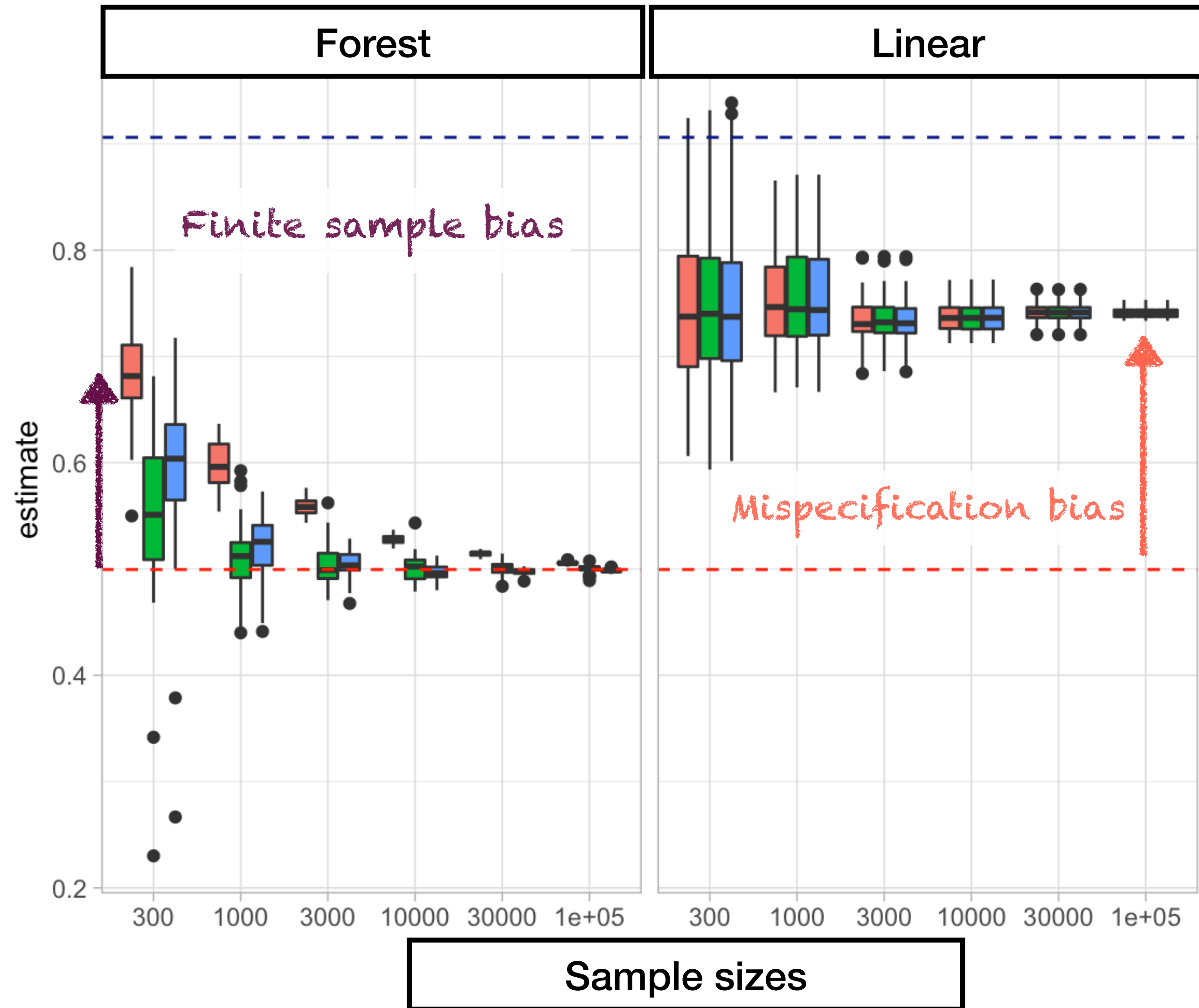
- Simulation set up from Nie and Wager
- Estimation with AIPW using either forest or linear models for nuisance parameters estimation



Why focusing on finite sample results? (1)

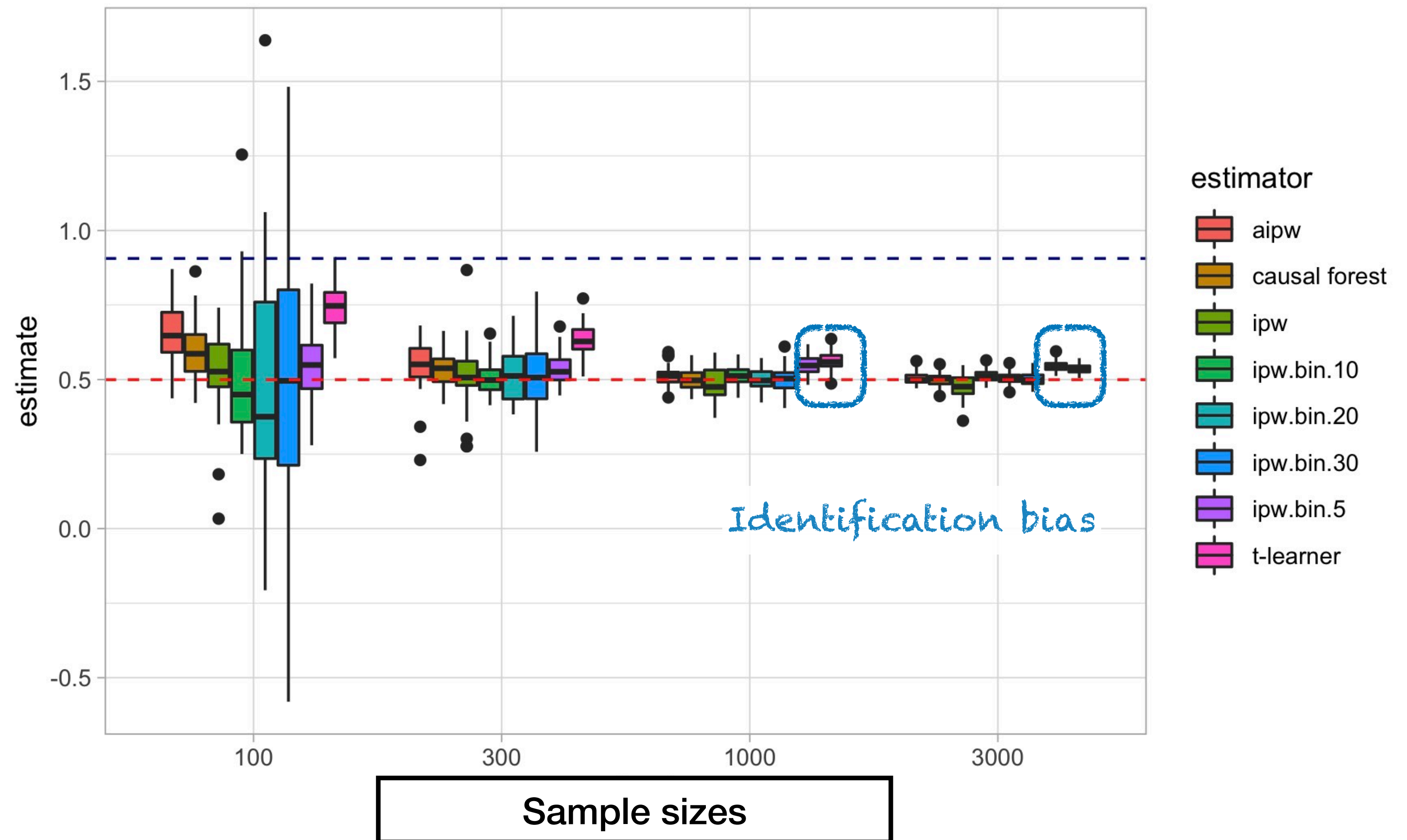
- (1) Usual sample sizes in medicine remains small
- (2) Results from simulations warned me and raised my interest

- Flexible estimation of the nuisance parameters guarantees large sample consistency...
- But at the cost of a **finite sample bias!**



Why focusing on finite sample results? (2)

- Flexible estimation of the nuisance parameters guarantees large sample consistency...
- But at the cost of a **finite sample bias**!
- Using a naive IPW with bins ensures a better finite sample risk than AIPW, at the cost of an **identification bias** that does not disappear with a bigger sample size.



Logistic regression and Russian roulette

Lemma 10 (Logit generative model for a binary outcome). *Considering a binary outcome Y , assume that*

$$\forall x \in \mathbb{X}, \forall a \in \{0, 1\}, \quad 0 < p_a(x) < 1, \quad \text{where } p_a(x) = \mathbb{P}(Y^{(a)} = 1 \mid X = x).$$

Then, there exist two functions $b, m : \mathcal{X} \rightarrow \mathbb{R}$ such that

$$\ln \left(\frac{\mathbb{P}(Y^{(a)} = 1 \mid X)}{\mathbb{P}(Y^{(a)} = 0 \mid X)} \right) = b(X) + a m(X).$$

Logistic regression and Russian roulette

Denoting $b_1(X)$ and $m_1(X)$ the functions for the intrication model, and $b_2(X)$ and $m_2(X)$ for the logistic model, one has:

$$b_2(X) = \ln \left(\frac{b_1(X)}{1 - b_1(X)} \right)$$

and

$$m_2(X) = \ln \left(\frac{(m_1(X) + b_1(X))(1 - b_1(X))}{1 - (m_1(X) + b_1(X))(1 - b_1(X))} \right) - \ln \left(\frac{b_1(X)}{1 - b_1(X)} \right)$$

Taking the case of the Russian Roulette, one has

$$b_1(X) := p_0(X), \quad m_1(X) = \frac{1}{6}$$

so that

$$b_2(X) := \ln \left(\frac{X}{1 - X} \right)$$

and

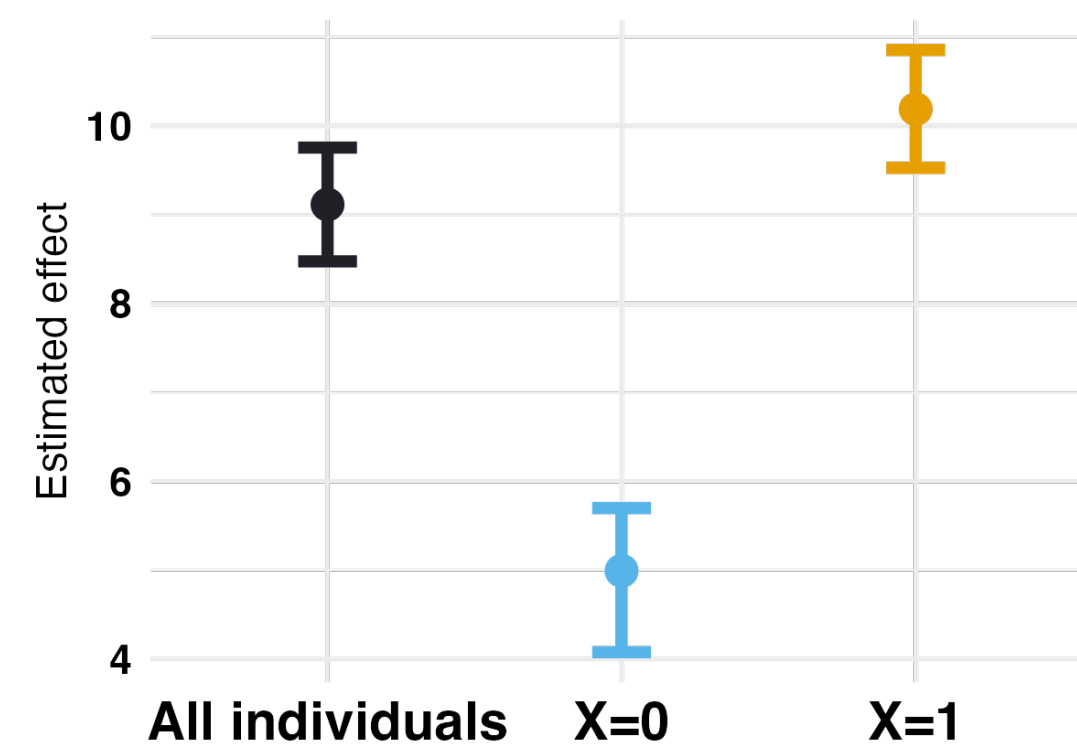
$$m_2(X) := \ln \left(\frac{\left(\frac{1}{6} + p_0(X)\right)}{1 - \left(\frac{1}{6} + p_0(X)\right)(1 - p_0(X))} \right) - \ln \left(\frac{p_0(X)}{1 - p_0(X)} \right).$$

Illustration on a toy simulation

Continuous outcome and binary baseline covariates X

	Target (\mathcal{P}_T)	Trial (\mathcal{P}_R)
$X = 1$	30%	75%
$X = 0$	70%	25%

Population's shift



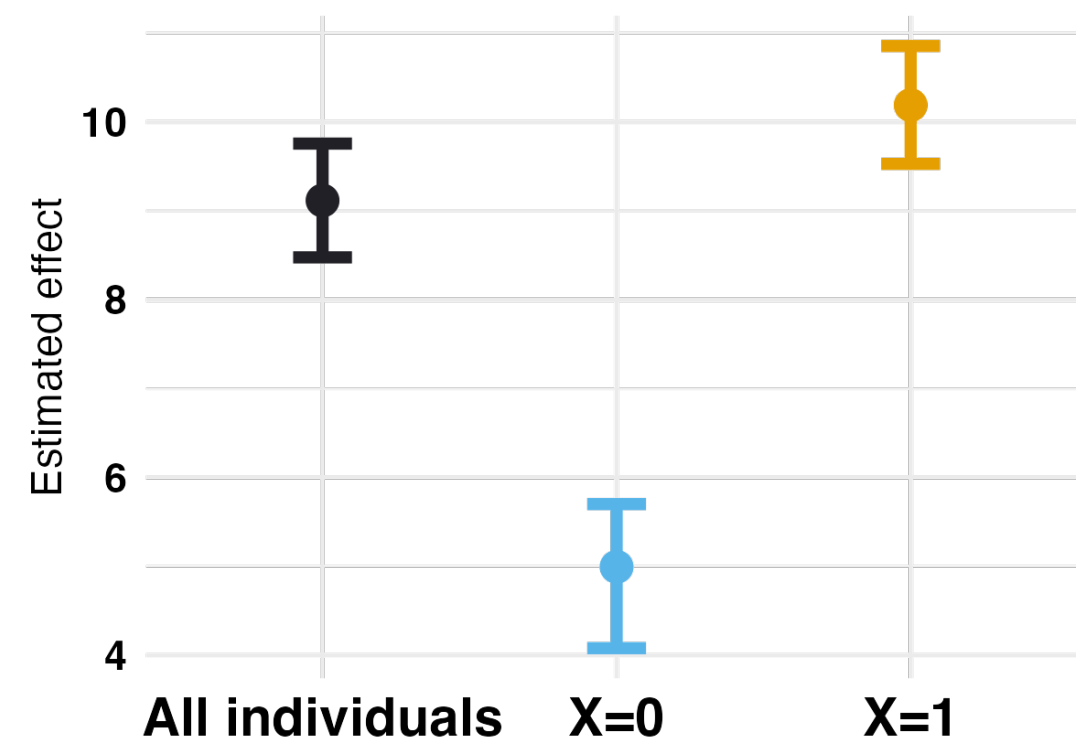
Hypothetical trial's results

Illustration on a toy simulation

Continuous outcome and binary baseline covariates X

	Target (\mathcal{P}_T)	Trial (\mathcal{P}_R)
$X = 1$	30%	75%
$X = 0$	70%	25%

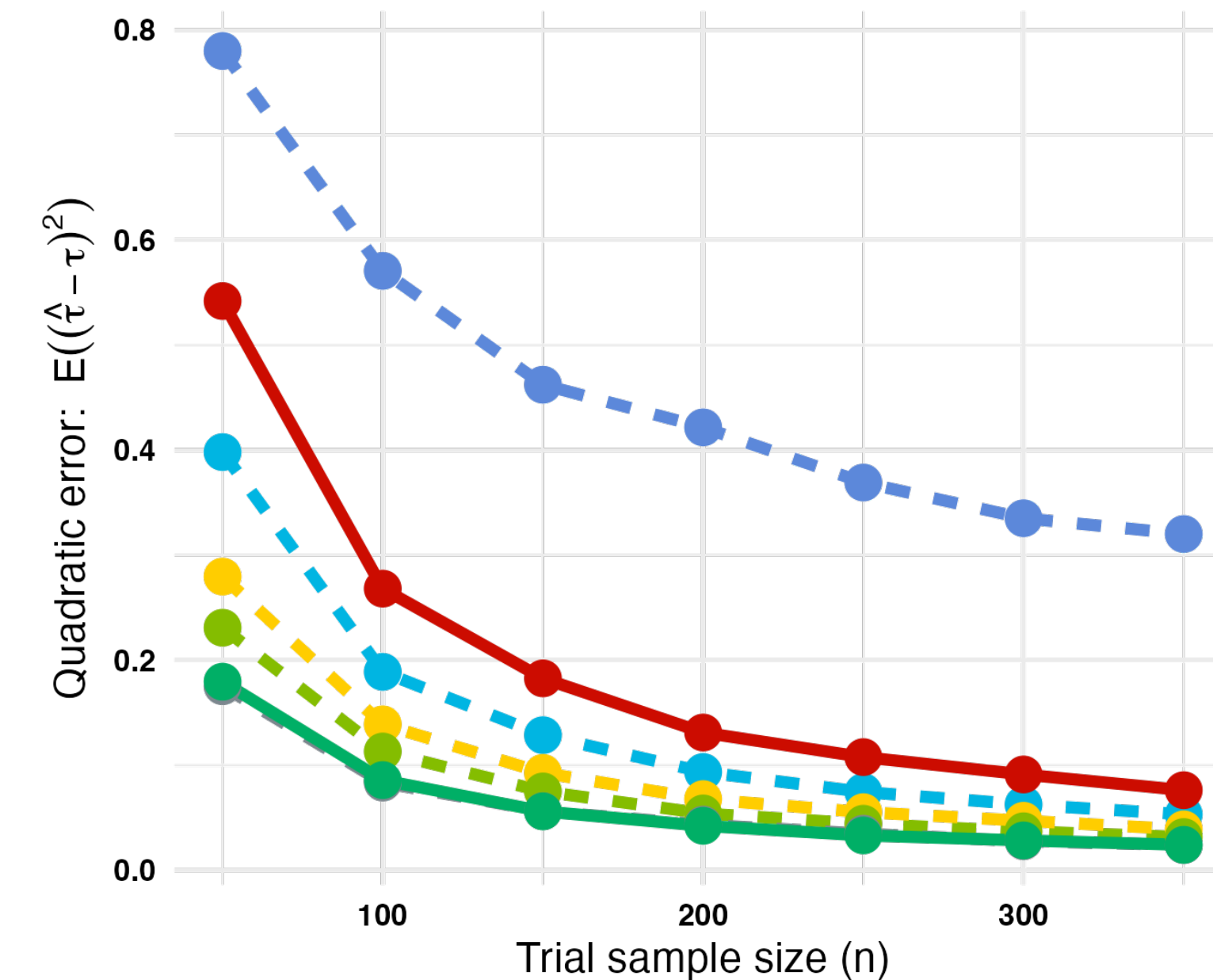
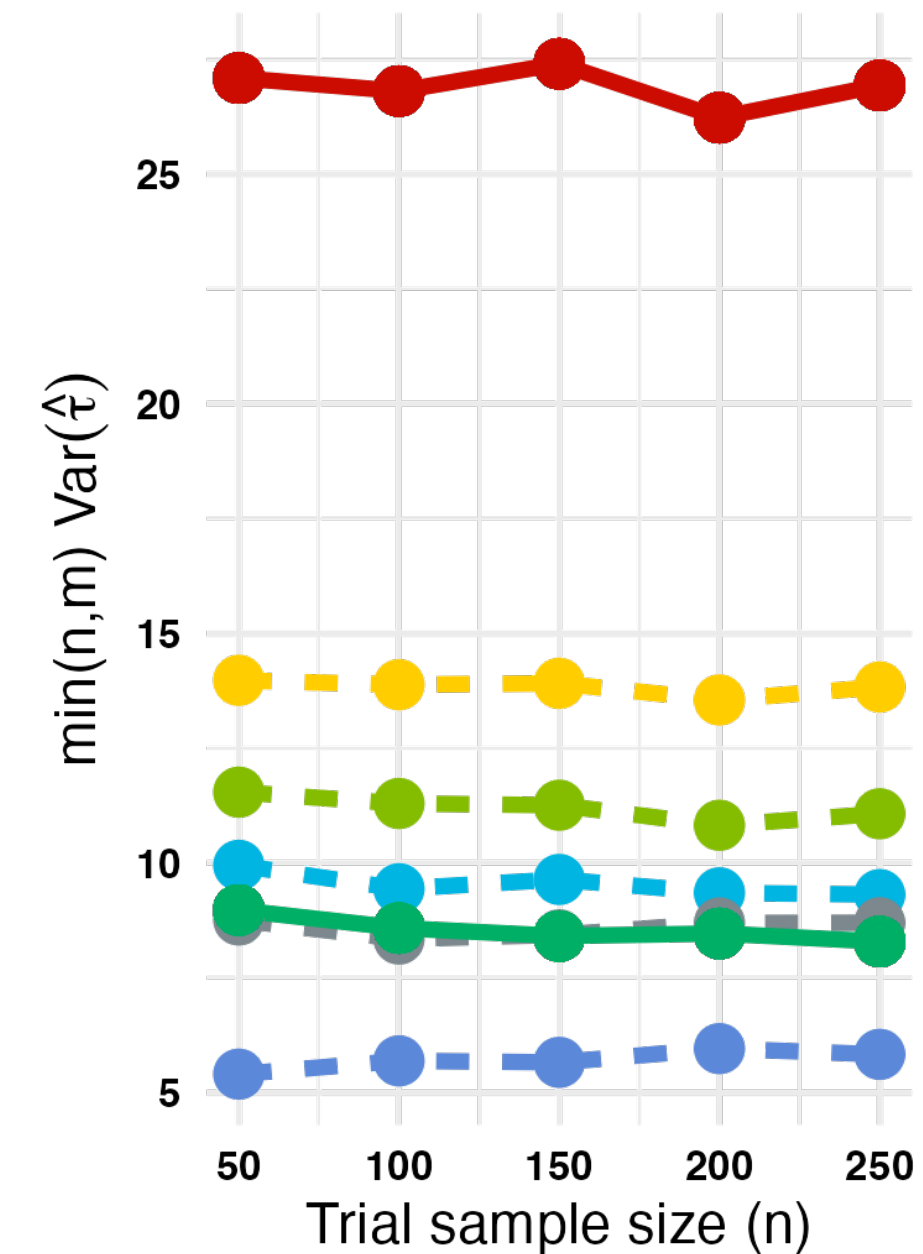
Population's shift



Hypothetical trial's results

IPSW
 — Completely or semi oracle
 - Estimated

Regime
 ● Completely-oracle
 ● IPSW: $m = \sqrt{n}$
 ● IPSW: $m = 2n$
 ● IPSW: $m = n$
 ● IPSW: $m = n^*n$
 ● IPSW: $m = n/2$
 ● Semi-oracle

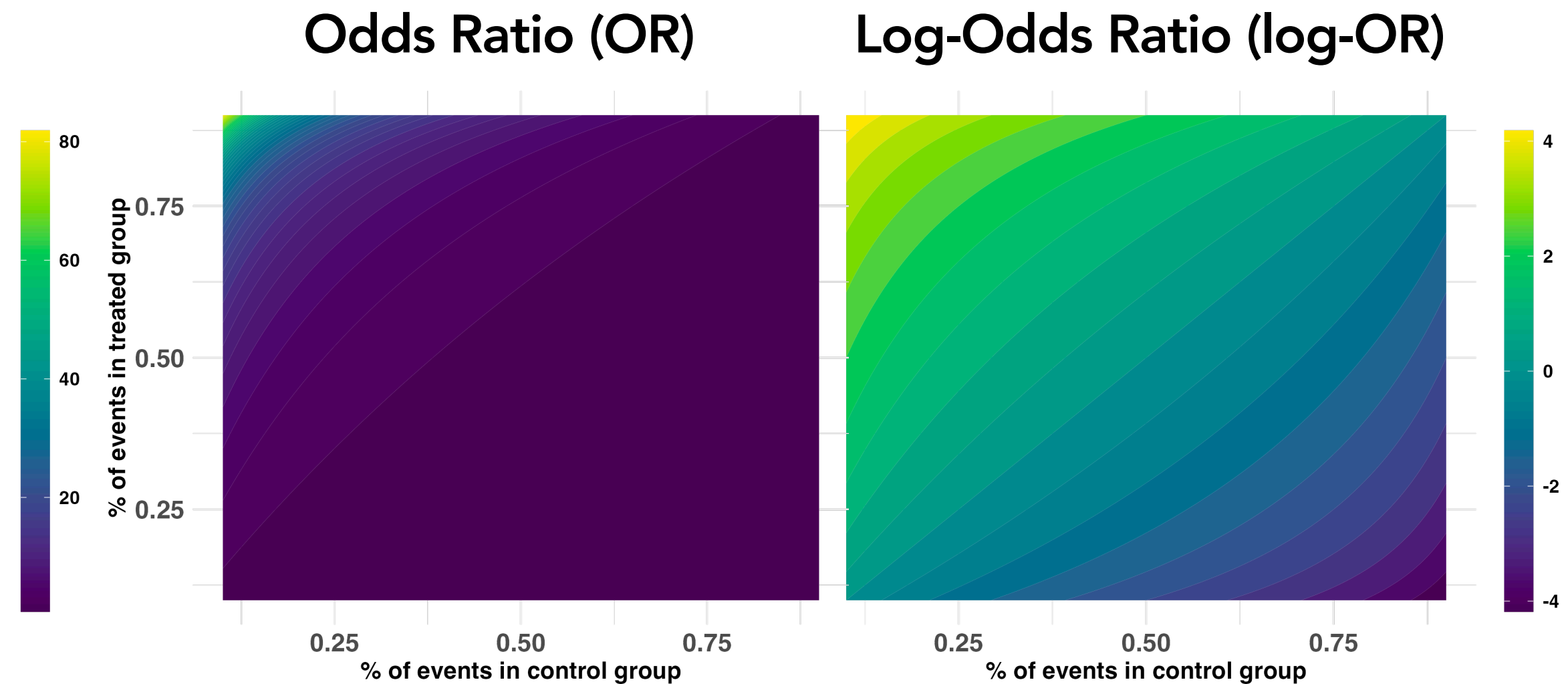
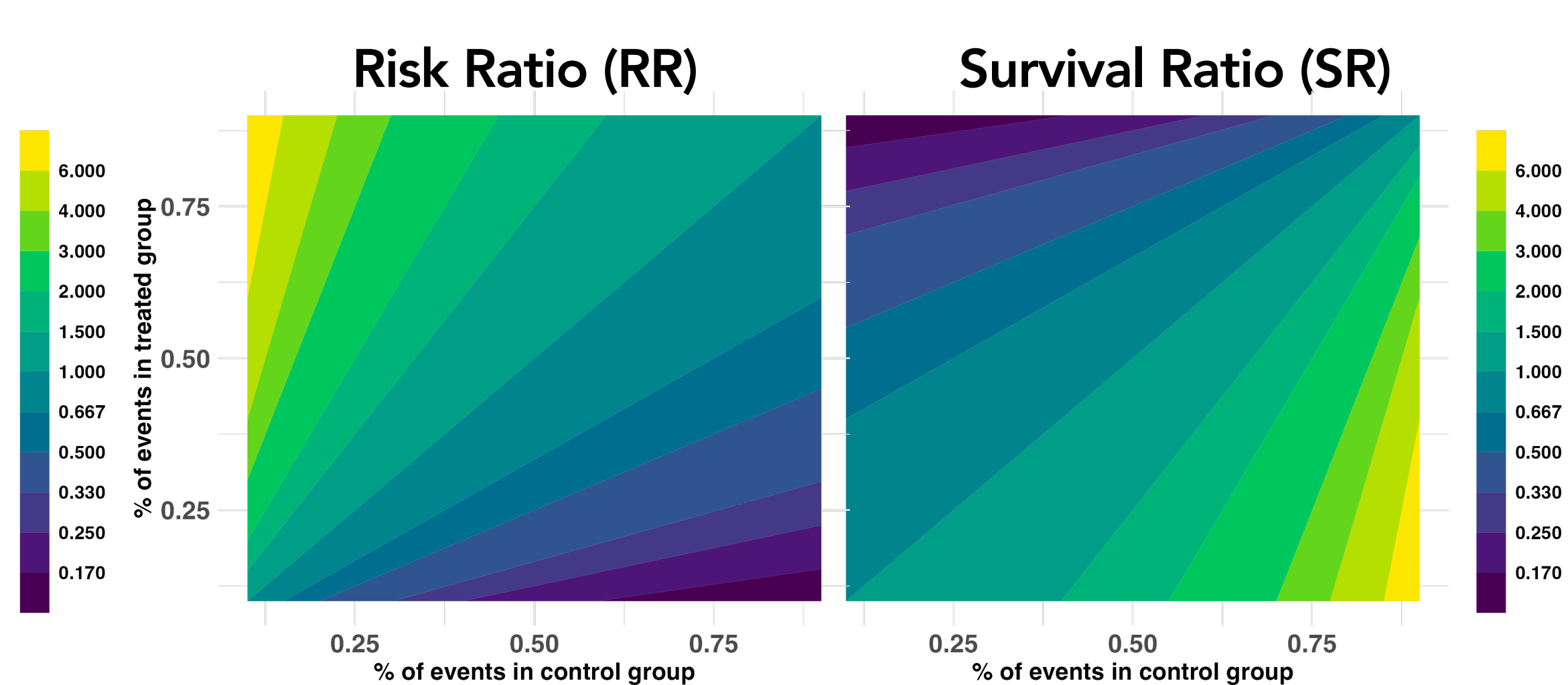
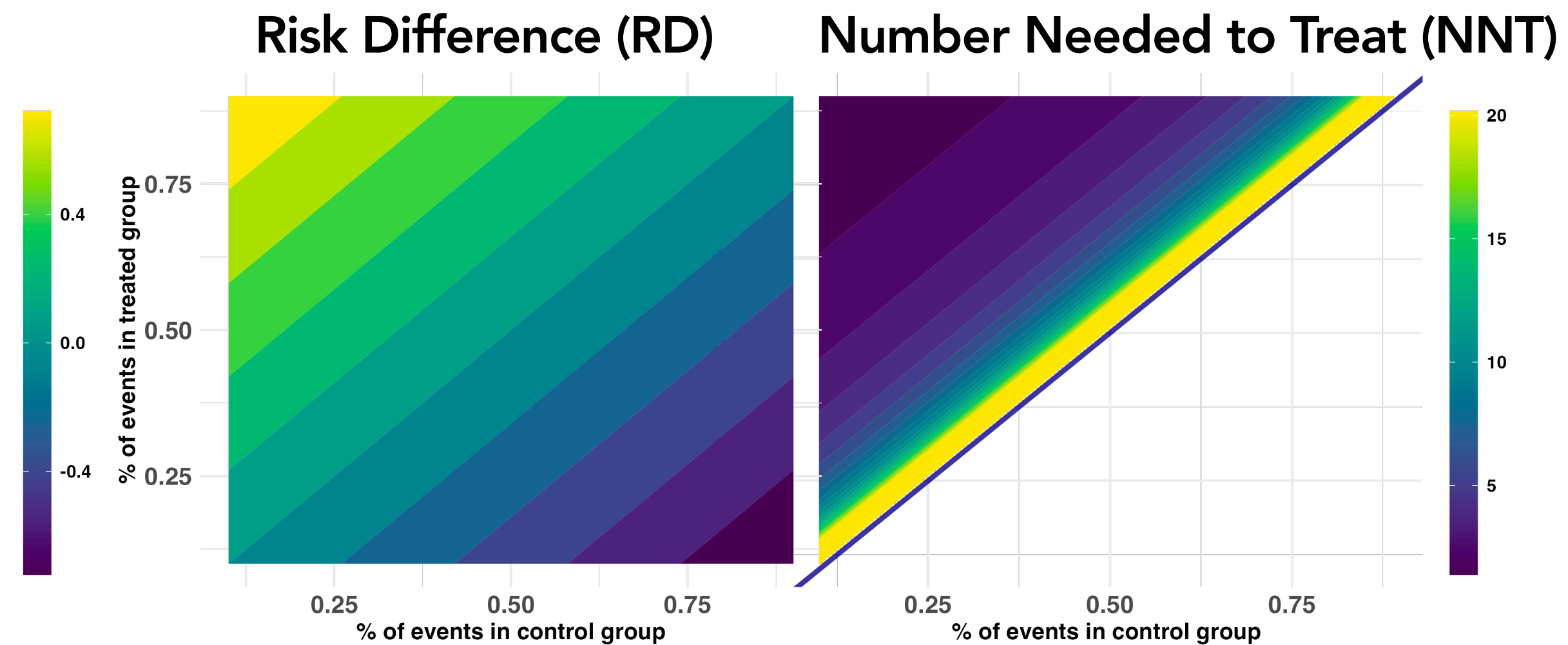
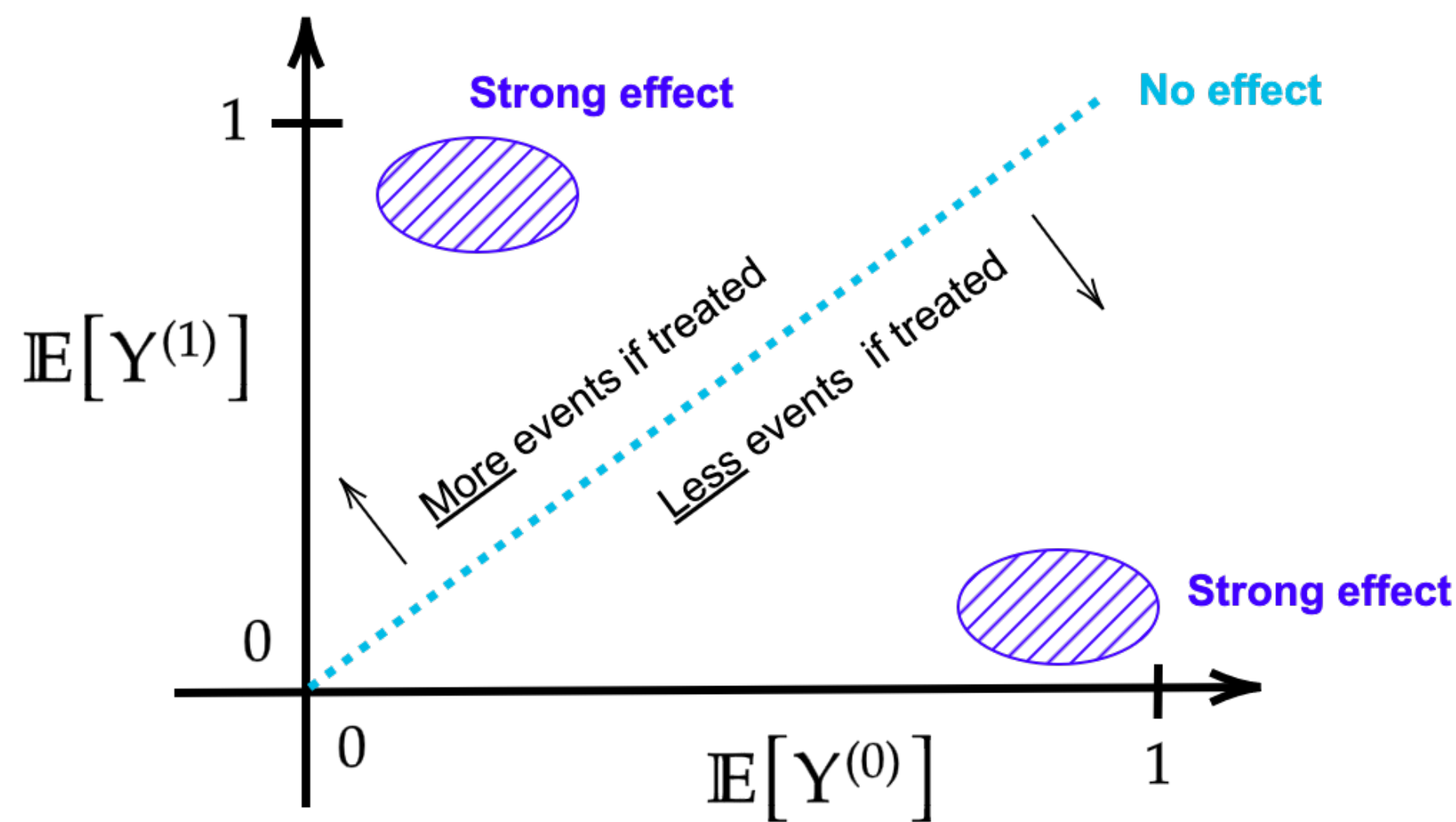


Empirical variance for different sizes n and m (6,000 repetitions for each dots) and different regimes

- Convergence speeds depend on the regime — i.e relative sizes of n and m ,
- Completely oracle IPSW has a bigger variance than the semi-oracle IPSW.

Ranges of effects

How to read plots



Common properties discussed

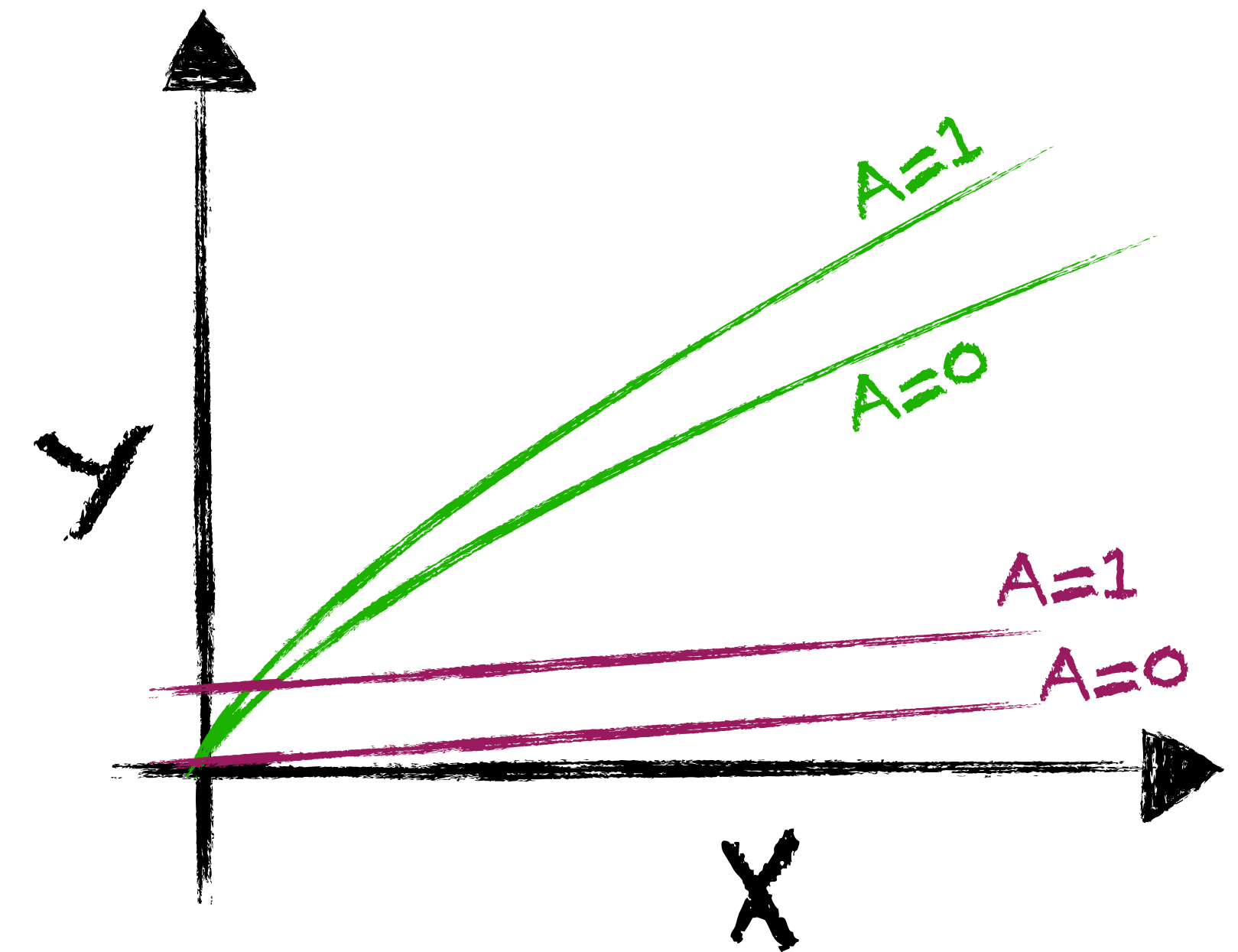
How the effect changes on sub-groups

Homogeneity $\forall x_1, x_2 \in \mathbb{X}, \tau(x_1) = \tau(x_2) = \tau$

Heterogeneity $\exists x_1, x_2 \in \mathbb{X}, \tau(x_1) \neq \tau(x_2)$

How the effect changes with labelling

e.g. Odds Ratio is symmetric, while Risk Ratio is not



! No non-zero effect can be homogeneous on all metrics

2+1 main approaches to generalize

1. **Re-weight** the trial individuals — *Inverse Propensity Sampling Weighting*
2. **Model the response** on the trial and impute the target sample — *plug-in G-formula*
3. **Combine the two** into a doubly robust approach — *A(ugmented) IPSW*

Consistency (Informal)

Considering that estimated surface responses are obtained following a cross-fitting estimation, then if IPSW or G-formula assumptions are ensured, then

$$\hat{\tau}_{AIPSW, n, m} \xrightarrow[n, m \rightarrow \infty]{L^1} \tau_T$$