

Generalizing a causal effect: review, sensitivity analysis, and missing covariates

ISCB 2021

Bénédicte Colnet, PhD student at Inria,

– advised by Julie Josse, Erwan Scornet, and Gaël Varoquaux,

– joint work with Imke Mayer.

Tuesday July, 21st

Motivation

Question of interest

Effect of acid tranexamic (TXA) on brain-injured related (TBI) deaths.

Data at hand

Randomized Controlled Trial - CRASH-3

- 29 different countries
- 9202 patients

Real World data - Traumabase

- 23 French Trauma centers
- 8270 patients

Is the RCT's estimate of the TXA effect the same for the Traumabase patients?

Outline for today's presentation

1. Review of the generalization estimators (<https://arxiv.org/abs/2011.08047>)
2. What if covariates from both data sets are different? (<https://arxiv.org/abs/2105.06435>)

Context

- Randomized Controlled Trials (RCT) : gold standard to estimate a treatment effect.

For example any new drug or treatment that receives an authorization usually has been assessed through several trials. This is part of the *evidence-based medicine*.

- Observational data have a higher representativeness but they can lack of internal validity.

The unconfoundedness assumption is unverifiable!




- Our motivation on tranexamic acid is part of a wider issue called **Generalization** or **Transportability**.

Remember the discussion on the Oxford-AstraZeneca vaccine's efficacy 

How can we leverage strengths of both type of data to gain information of a target population treatment efficacy?

Notations

Our notations take place in the Neyman-Rubin framework.

- A treatment of interest, 
- X covariates, 
- Y the outcome, 

The target quantity is the average population treatment effect,

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

External validity bias

We introduce s an indicator or eligibility in the trial ($s = 1$ corresponds to eligibility)

The distribution of covariates x is not the same in the target population and in the RCT,

$$f_X \neq f_{X|S=1}.$$

External validity

So that,

$$\tau_1 = \mathbb{E}[Y(1) - Y(0)|S = 1] \neq \mathbb{E}[Y(1) - Y(0)] = \tau.$$

Using two data sets?

	S	Set	Covariates			Treatment	Outcome
			X_1	X_2	X_3	A	Y
1	1	\mathcal{R}	1.1	20	F	1	1
	1	\mathcal{R}	-6	45	F	0	1
n	1	\mathcal{R}	0	15	M	1	0
$n+1$	0	\mathcal{O}	
	0	\mathcal{O}	-2	52	M	-	-
	1	\mathcal{O}	-1	35	M	-	-
$n+m$	0	\mathcal{O}	-2	22	M	-	-

This is called *generalization*, and is related to *data fusion*, *transportability*, and *covariate shift* problem.

Identification

But **first**, an important step is to ensure the identifiability of τ , and the two major assumptions are:

- **Ignorability assumption on trial participation**

$$Y(1) - Y(0) \perp S \mid X$$

💡 X contains all covariates that are *treatment effect modifiers* and with a distributional shift.

- **Positivity of trial participation**

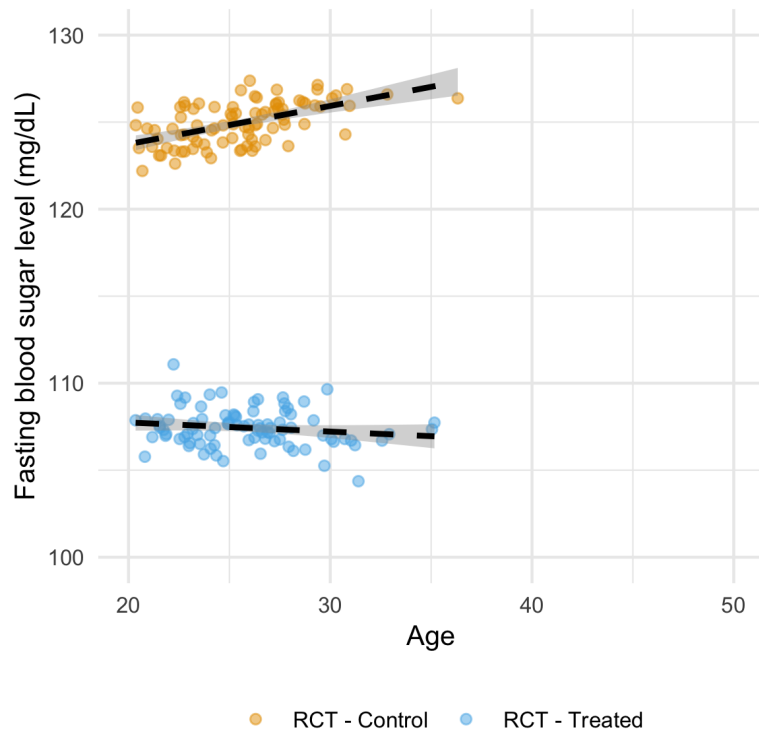
There exists a constant c such that for all x with probability $\mathbf{1}$, $\mathbb{P}(S = 1 \mid X = x) \geq c > 0$

💡 Each individual from the target population had a non-zero probability to be eligible in the trial.

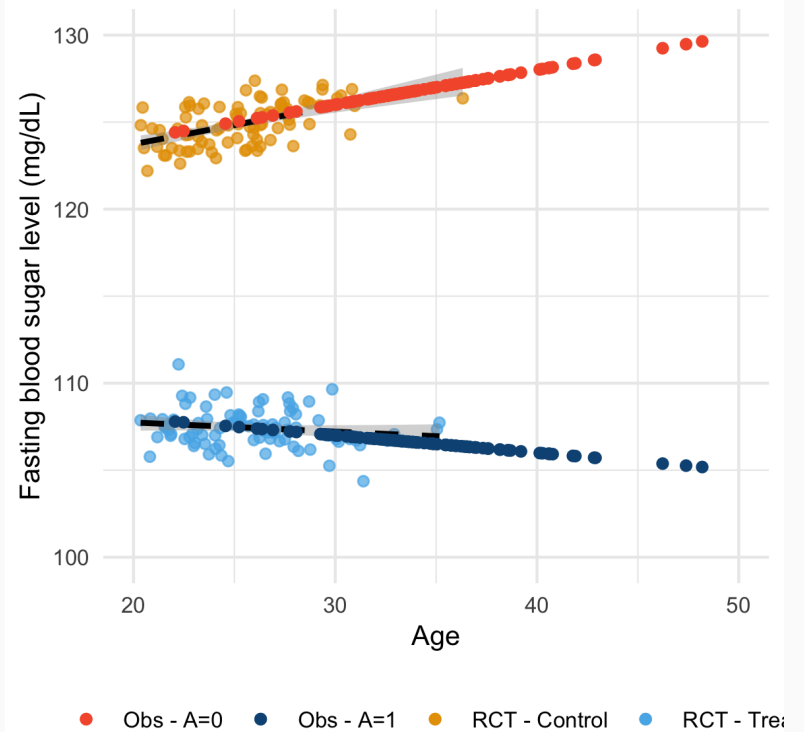
Outcome regression (G-formula)

Intuition

Step 1



Step 2



Outcome regression (G-formula)

Formalization

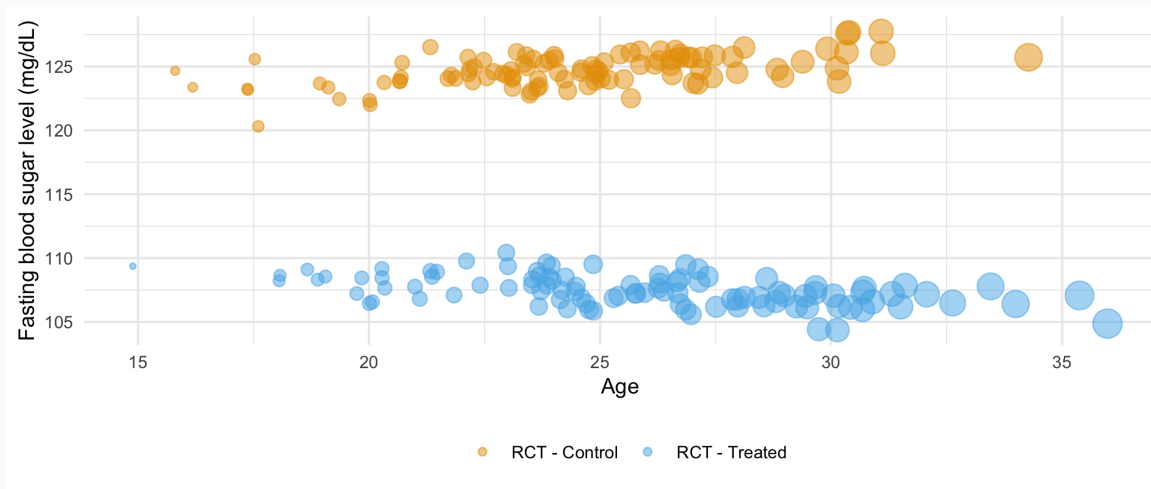
$$\hat{\tau}_{G,n,m} = \frac{1}{m} \sum_{i=n+1}^m \left(\hat{\mu}_{1,n}(X_i) - \hat{\mu}_{0,n}(X_i) \right),$$

where,

- $\mu_a(\boldsymbol{x}) \triangleq \mathbb{E}[Y(\boldsymbol{w}) | \boldsymbol{X} = \boldsymbol{x}, \boldsymbol{A} = \boldsymbol{a}]$ are the response surfaces,
- $\hat{\mu}_{a,n}(X_i)$ are estimated on the RCT sample.

Weighting (IPSW)

Intuition



Formalization

$$\hat{\tau}_{IPSW, n, m} = \frac{1}{m} \sum_{i=1}^n \frac{Y_i}{\hat{\alpha}_{n, m}(X_i)} \left(\frac{A_i}{e_1(X_i)} - \frac{1 - A_i}{1 - e_1(X_i)} \right),$$

where $\hat{\alpha}_{n, m}$ is an estimate of the odd ratio of the indicatrix of being in the RCT, and $e_1(X) = P(A = 1 \mid X = x, S = 1)$

Toward sensitivity analysis

What if a covariate is missing?

Mathematically, $\mathbf{X} = \mathbf{X}_{mis} \cup \mathbf{X}_{obs}$, and

$$Y(1) - Y(0) \not\perp S \mid \mathbf{X}_{obs}$$

⚠️ Such a missing covariate breaks the identifiability assumption.

What can we do?

- 💡 The intuition is that *a poorly shifted missing covariate and/or a weak treatment effect missing covariate will lead to a small bias.*
- 🤔 Is there a way to link the bias to these two characteristics, so that domain expert can help assess whether or not the missing covariate breaks any conclusion?
- 📖 Such an approach is called a sensitivity analysis

Key result - Model

Model

We admit there exist $\delta \in \mathbb{R}^p$, and $\sigma \in \mathbb{R}^+$ such that the semi-linear ^[1] model holds:

$$Y = g(X) + A\langle X, \delta \rangle + \varepsilon, \quad \text{where } \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

Assumption

The distribution of X is Gaussian, that is, $X \sim \mathcal{N}(\mu, \Sigma)$, and transportability of Σ is true, that is, $X \mid S = 1 \sim \mathcal{N}(\mu_{S=1}, \Sigma)$.

Theorem in a sentence

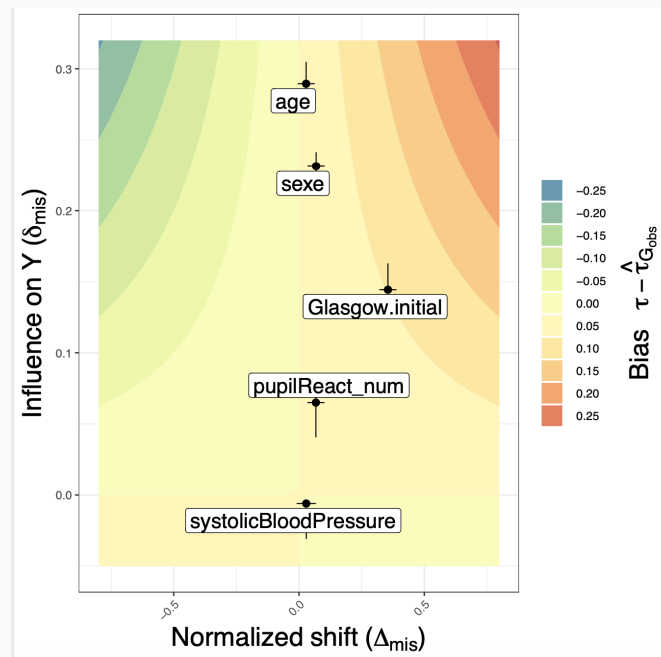
Under these assumptions, the bias of IPSW and G-formula are the same, that is:

$$B = \sum_{j \in mis} \delta_j \left(\mathbb{E}[X_j] - \mathbb{E}[X_j \mid S = 1] - \Sigma_{j,obs} \Sigma_{obs,obs}^{-1} (\mathbb{E}[X_{obs}] - \mathbb{E}[X_{obs} \mid S = 1]) \right).$$

Sensitivity analysis in the real-world

⚠ We had to take a surrogate outcome that is continuous! --> Disability Rating Scale (DRS)

Target population ATE estimation with the G-formula on the set of observed covariate: 0.17 [95% CI -0.34 - 0.29]) with bootstrap.



Confidence intervals are done with bootstrap.

Conclusion and perspectives

Contributions

1. Handling all missing covariate patterns,
2. Lighten the usual independency condition,
3. Insist on interpretability.

Also present in the paper

- Imputation?
- Proxy?

Place for improvement

- Lighten the semi-parametric and Gaussian assumption?
- Binary outcome.

Another related work



	Set	Covariates			Treatment	(Factual) Outcome
		X_1^*	X_2^*	X_3^*	A	Y
1	\mathcal{R}	1.1	20	NA	1	23.4
...	\mathcal{R}	
$n-1$	\mathcal{R}	-6	NA	8.3	0	26.3
n	\mathcal{R}	0	15	6.2	1	28.1
$n+1$	\mathcal{O}	-2	52	NA	NA	NA
$n+2$	\mathcal{O}	-1	NA	2.4	NA	NA
...	\mathcal{O}		...		NA	NA
$n+m$	\mathcal{O}	NA	NA	3.4	NA	NA

Generalizing with missing attributes?
(Mayer et al. 2021)

<https://arxiv.org/abs/2104.12639>

Thank you for listening!

Code: on Github [BenedicteColnet/unobserved-covariate](https://github.com/BenedicteColnet/unobserved-covariate)

? Questions / remarks / discussions / ideas are welcome : benedicte.colnet@inria.fr