

# Computationally Efficient Analysis of Randomized Trials with Non-Monotone Missing Binary Outcomes

Jaron J. R. Lee<sup>1</sup> (jaron.lee@jhu.edu), Daniel O. Scharfstein<sup>2</sup> (dscharf@jhu.edu) and Ilya Shpitser<sup>1,2</sup> (ilyas@cs.jhu.edu)

Department of Computer Science<sup>1</sup>, Department of Biostatistics<sup>2</sup>

## Summary

- Missing not at random (MNAR) models are the most realistic model class for studying non-monotone clinical trial data (Little and Rubin, 2014)
- (1) introduce an MNAR model where the probability of missingness at a visit depends on all unobserved outcomes prior to the visit, and all observed outcomes after the visit
- Problem:** This model is intractable for large  $K$  without further assumptions
- Idea:** Develop computational methods to efficiently estimate this model
  - Restrict to binary data
  - Introduce a Markov restriction on dependencies
  - Use directed acyclic graph (DAG) theory to derive a tractable recursive identification and estimation strategy

## The Unrestricted Robins Model

For  $i = 1, \dots, K$  time points,

- $Y_i^{(1)} \in \{0, 1\}$
- $R_i \in \{0, 1\}$
- $Y_i \equiv \begin{cases} Y_i^{(1)}, & R_i = 1 \\ ?, & R_i = 0 \end{cases}$
- $O_k = (R_k, Y_k^{(1)})$

- $\bar{Z}_k = (Z_1, \dots, Z_{k-1})$  for  $k = 2, \dots, K$ ,
- $Z_k = (Z_{k+1}, \dots, Z_K)$  for  $k = 1, \dots, K-1$ ,
- $\bar{Z}_k^m = (Z_{\max(1, k-m)}, \dots, Z_{k-1})$  for  $k = 2, \dots, K$ ,
- $Z_k^m = (Z_{k+1}, \dots, Z_{\min(k+m, K)})$  for  $k = 1, \dots, K-1$

### Model Assumptions:

$$\forall k \in 1, \dots, K$$

$$p(Y_k^{(1)} | R_k = 0, \bar{Y}^{(1)}_k, O_k) = p(Y_k^{(1)} | R_k = 1, \bar{Y}^{(1)}_k, O_k)$$

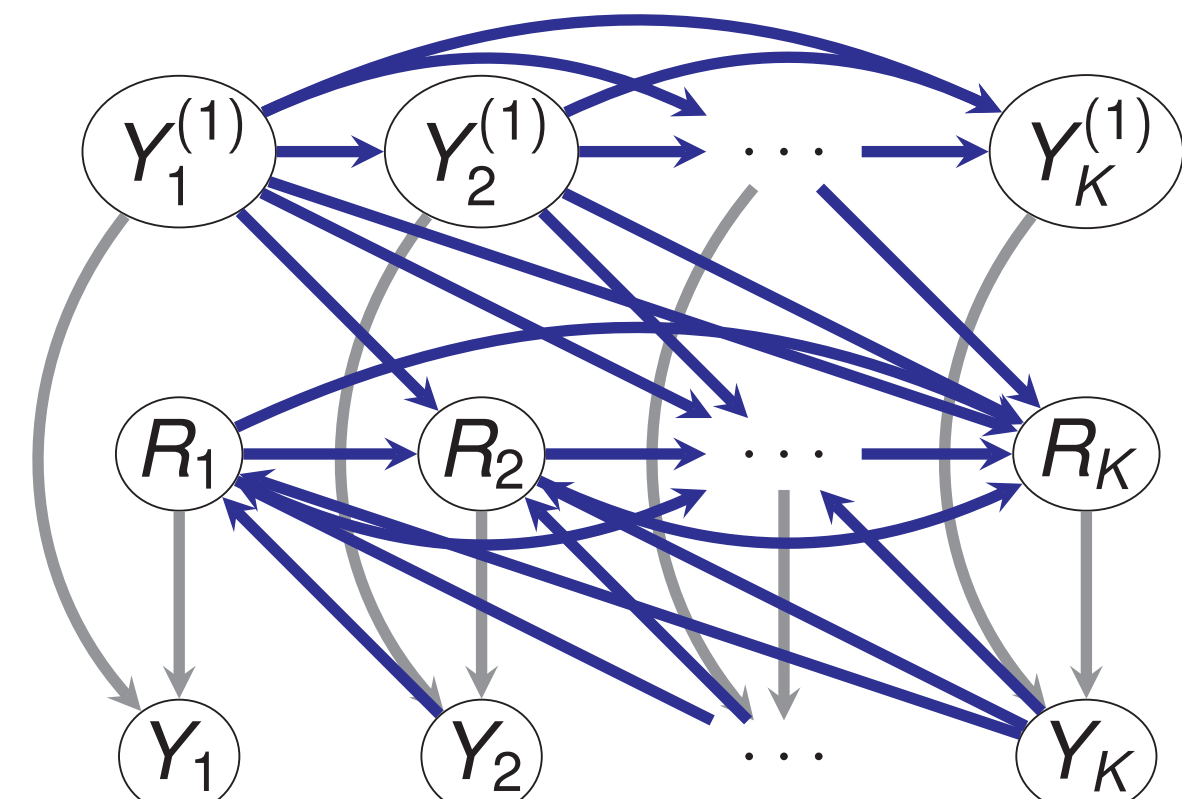


Figure 1: Unrestricted model

## Identification in the Unrestricted Model

**Theorem 1:**  $p(\bar{Y}^{(1)}_k)$  in the unrestricted model is identified.

- By induction we can show that  $p(\bar{Y}^{(1)}_k, O_k)$  is identified.
- For  $k = 0$ , this is  $p(O_1)$  which is observed.
- Suppose identified for  $k = s$ . Then,

$$p(\bar{Y}^{(1)}_s, O_s) = \sum_{j=0}^1 p(\bar{Y}^{(1)}_{s-1}, Y_s^{(1)}, R_s = j, O_s)$$

- $R_s = 1$  is identified by induction.
- $R_s = 0$  is identified since

$$p(\bar{Y}^{(1)}_{s-1}, Y_s^{(1)}, R_s = 0, O_s) = \frac{\text{identified by model assumption}}{p(Y_s^{(1)} | R_s = 0, \bar{Y}^{(1)}_{s-1}, O_s)}$$

$$\frac{p(\bar{Y}^{(1)}_{s-1}, R_s = 0, O_s)}{\text{identified by induction assumption}}$$

- Issue:** The proof relies on probability distributions over  $O(2^K)$  elements ( $K$  binary variables) - inference for large  $K$  intractable!

## The Restricted Robins Model

- Solution:** Use Markov restrictions to avoid estimating high-dimensional distributions - related to efficient calculation of causal effects in (2)
- Let  $m$  denote the Markov restriction. At time  $t$ ,

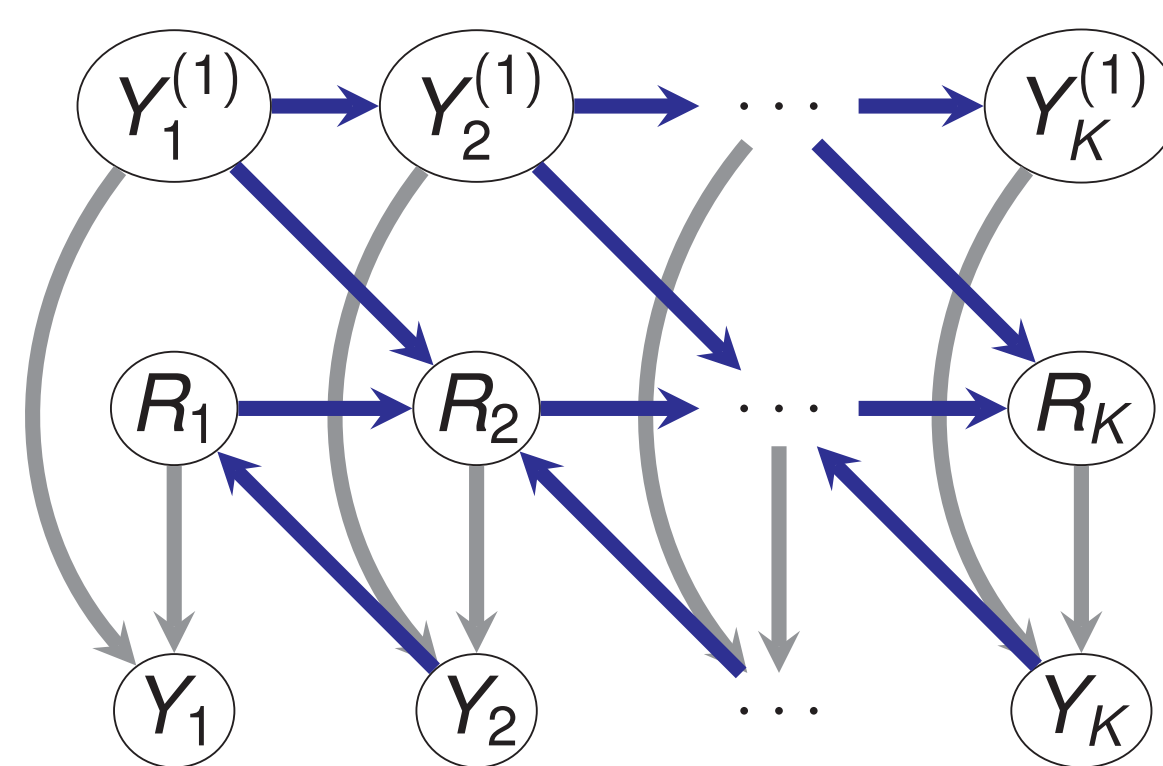


Figure 2:  $m = 1$

### Model Assumptions:

$$\forall k \in 1, \dots, K$$

$$p(Y_k^{(1)} | R_k = 0, \bar{Y}^{(1)}_k^m, O_k^m) = p(Y_k^{(1)} | R_k = 1, \bar{Y}^{(1)}_k^m, O_k^m)$$

$$p(Y_k^{(1)} | \bar{Y}^{(1)}_k) = p(Y_k^{(1)} | \bar{Y}^{(1)}_k^m)$$

- Remark:** Since the unrestricted model is identified, the restricted model must also be identified.
- Remark:**  $p(R_i | \bar{Y}^{(1)}_i^m, O_i^m)$  is a probability distribution over  $O(2^{2m})$  elements - over large  $K$  is tractable for fixed  $m$ .

### Theorem 2:

$$R_k \perp\!\!\!\perp Y_k^{(1)} | \bar{Y}^{(1)}_k^m, O_k^m \quad (2)$$

## Estimation of the Restricted Model

The identification of the unrestricted model offers no insight into the estimation of the restricted model that is linear in  $K$  for a fixed  $m$ .

- We derive a *recursive* estimation strategy.
- For fixed  $m$ , linear in  $K$ .

### Base case:

$$p(\bar{Y}^{(1)}_2^{m+1}, O_1^{m+1}) = p(O_{m+2} | Y_1^{(1)}, O_1^m) p(Y_1^{(1)} | O_1^m) p(O_1^m)$$

$$= \underbrace{p(O_{m+2} | R_1 = 1, Y_1^{(1)}, O_1^m) p(Y_1^{(1)} | R_1 = 1, O_1^m)}_{\text{by Theorem 2}} p(O_1^m)$$

**Inductive case:** Assume that  $p(\bar{Y}^{(1)}_{k+1}^{m+1}, O_k^{m+1})$  is identified. Then,

$$p(\bar{Y}^{(1)}_{k+1}^{m+1}, O_k^{m+1}) = p(O_{k+m+1} | \bar{Y}^{(1)}_{k+1}^{m+1}, O_k^m) p(Y_k^{(1)} | \bar{Y}^{(1)}_k^m, O_k^m) p(\bar{Y}^{(1)}_k^m, O_k^m)$$

$$= \underbrace{p(O_{k+m+1} | R_k = 1, \bar{Y}^{(1)}_{k+1}^{m+1}, O_k^m) p(Y_k^{(1)} | R_k = 1, \bar{Y}^{(1)}_k^m, O_k^m) p(\bar{Y}^{(1)}_k^m, O_k^m)}_{\text{function of } p(\bar{Y}^{(1)}_2^{m+1}, O_1^{m+1})}$$

We want to show that  $p(\bar{Y}^{(1)}_m^m, O_{k-1}^{m+2})$  is a function of  $p(\bar{Y}^{(1)}_{k+1}^{m+1}, O_k^{m+1})$ .

$$p(\bar{Y}^{(1)}_k^m, O_{k-1}^{m+1}) = \begin{cases} p(\bar{Y}^{(1)}_{k+1}^{m+1}, O_k^{m+1}), & 2 \leq k \leq m+1 \\ \sum_{Y_{k-m-1}^{(1)}} p(\bar{Y}^{(1)}_k^{m+1}, O_{k-1}^{m+1}), & m+1 < k \leq K+1 \end{cases}$$

$$p(\bar{Y}^{(1)}_m^m, O_{k-1}^{m+2}) = \begin{cases} \text{by Theorem 2} \\ p(O_{k+m+1} | \bar{Y}^{(1)}_k^m, O_{k-1}^{m+1}, \bar{R}_k^m = \mathbf{1}) p(\bar{Y}^{(1)}_k^m, O_{k-1}^{m+1}), & k \leq K-m-1 \\ p(\bar{Y}^{(1)}_k^m, O_{k-1}^{m+1}), & k > K-m-1 \end{cases}$$

## Application: Drug Use Abatement Dataset

- Study:** National Institute of Drug Abuse Study No. CTN-0044 - evaluating effectiveness of computer tool in outpatient substance abuse treatment. Patient dropout is non-monotone, 20-25% missing.
  - $N = 500$ ,  $K = 24$
  - $Y$ : negative drug test (abstinence)
- Question:** Evaluate effectiveness of treatment vs. control
  - $T = \sum_{k=1}^{24} \mathbb{E}[Y_k^{(1)}]$
- Estimation:** Discrete probability distributions estimated using random forests
  - Use random forests to model  $p(O_K | O_K^m)$ ,  $p(O_{K-1} | O_{K-1}^m)$ ,  $\dots$ ,  $p(O_1)$
  - Tune random forests to avoid positivity violations

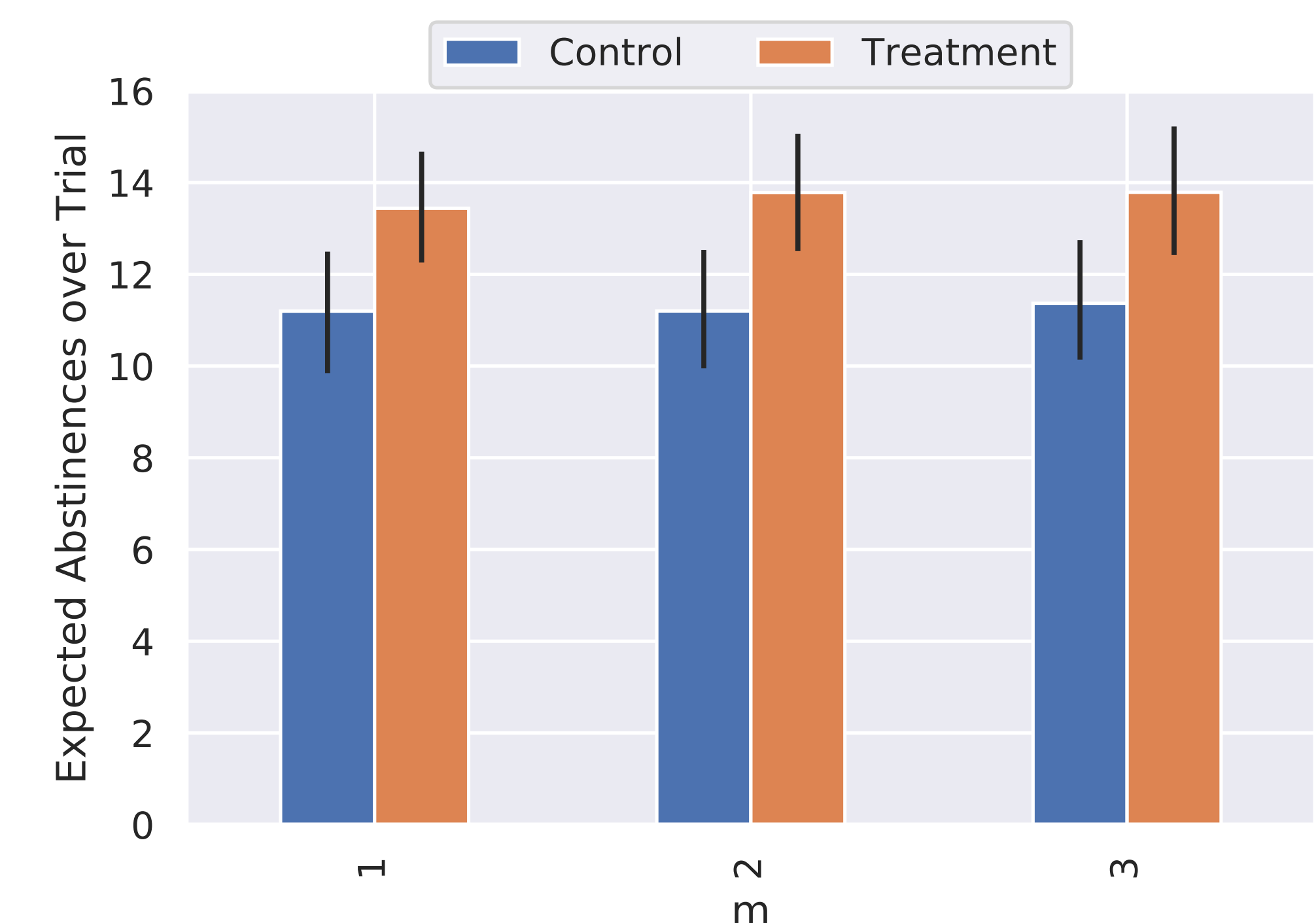


Figure 3: Expected abstinences over trial for varying  $m$ , with bootstrapped 95% CI ( $N_{\text{bootstrap}} = 500$ )

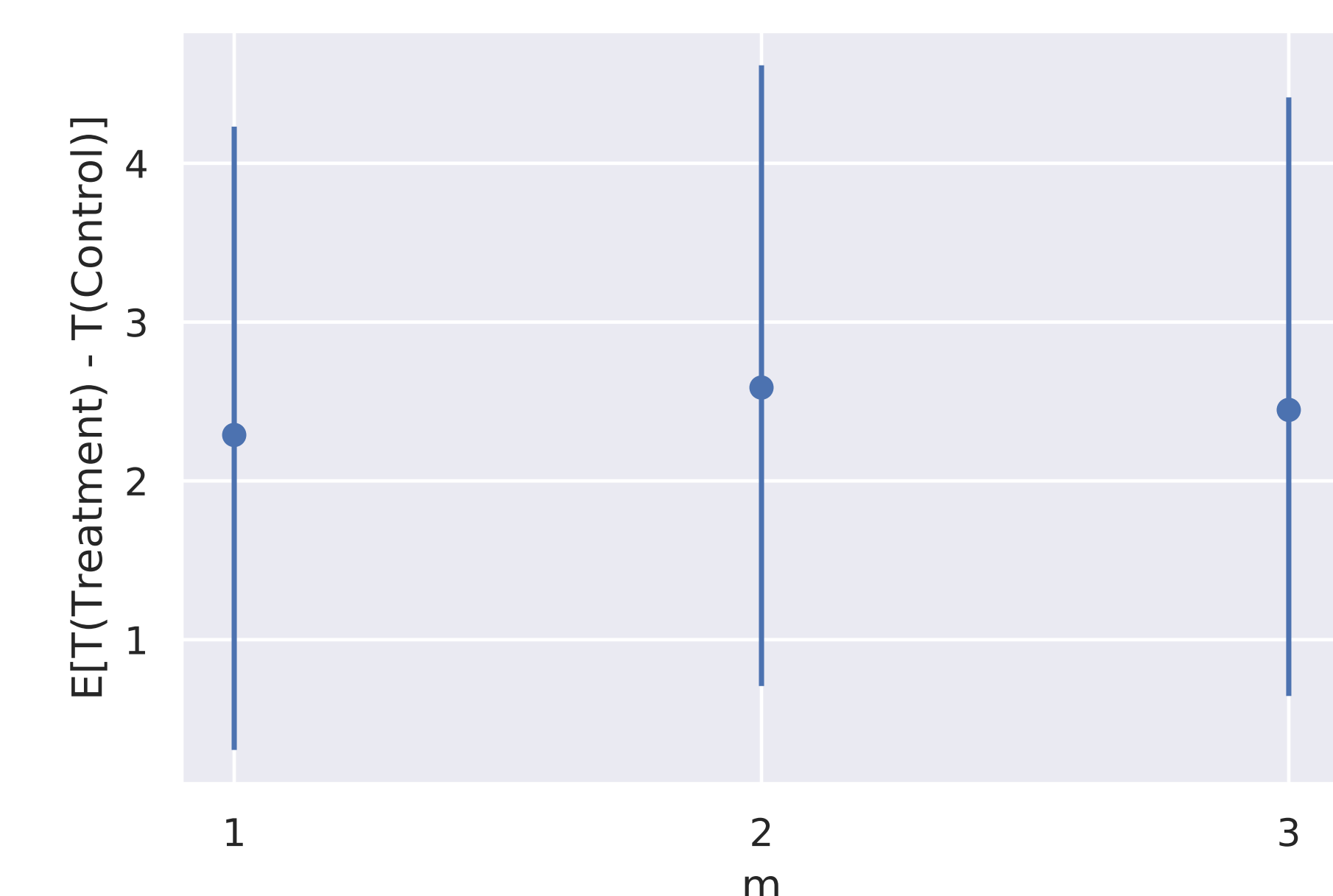


Figure 4: Expected treatment effect for varying  $m$ , with bootstrapped 95% CI ( $N_{\text{bootstrap}} = 500$ )

## References

- ROBINS, J. M. Non-Response Models for the Analysis of Non-Monotone Non-Ignorable Missing Data. *Statistics in Medicine* 16, 1 (1997), 21–37.
- SHPIITSER, I., RICHARDSON, T. S., AND ROBINS, J. M. An Efficient Algorithm for Computing Interventional Distributions in Latent Variable Causal Models. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence* (2011), UAI'11, AUAI Press, pp. 661–670.