## When Causal Forests Mislead **Evaluating the precision of Confidence Intervals**

This study tackles an important issue in evaluating the reliability of confidence intervals in causal forests by examining how data characteristics and hyperparameters influence actual coverage rates compared to theoretical benchmarks. A primary discovery is the identification of a practical limit for reliable confidence interval coverage: When the sum of confounders and effect modifiers exceeds 4, coverage rates drop considerably below 80%, even for simple treatment effect functions.



Within-group correction Between-group variance Treatment effect estimates are asymptotically Gaussian

$$\hat{\tau}(x) \pm 1.96 \sqrt{\hat{V}_{BLB}(x)}$$

## **Related Literature**

Leo Breiman. "Random forests". In: Machine learning 45.1 (2001), pp. 5–32.

Stefan Wager and Susan Athey. "Estimation and inference of heterogeneous treatment effects using random forests". In: Journal of the American Statistical Association 113.523 (2018), pp. 1228–1242. Susan Athey, Julie Tibshirani, and Stefan Wager. "Generalized random forests". In: Annals of Statistics 47.2 (2019), pp. 1148–1178.

Author Rares lordan r.iordan@student.tudelft.nl Supervisor **Rickard Karlsson** Professor Jesse Krijthe

## Methodology

Data Generating Process: Polynomial datagenerating process that generates synthetic

Representation with maximum interaction order of 2 to evaluate sensitivity of coverage rates

## 05 Findings

- Critical Threshold Identified
- When the combined number of confounders and effect modifiers exceeds 4, coverage rates decline dramatically below 80% even for the simplest treatment effect function. This threshold appears robust, as increasing computational resources provided only marginal improvements.
- Most Influential Parameters

• The number of confounders, their interaction with effect modifiers, and effect modifiers are the dominant factors (sensitivity indices  $\approx$  $0.28 \pm 0.04$ ,  $\approx 0.14 \pm 0.03$ , and  $\approx 0.10 \pm 0.02$ )

- Hyperparameter Recommendations
  - max\_depth: Leave unset for best results
  - max\_samples: Increase to 0.5 for best coverage rate performance
  - min\_balancedness\_tol: Set to 0.5 for optimal coverage
  - n\_estimators: Use  $\ge$  2400 trees for best performance
  - min\_impurity\_decrease: Keep at default (0.0) to avoid degrading inference quality

	0	1	1 2 3 4 5 6 9					
0		0.945 ±0.001	0.907 ±0.015	0.756 ±0.048	0.573 ±0.001	0.531 ±0.006	0.394 ±0.030	0.333 ±0.018
	0.939	0.931	0.788	0.655	0.462	0.433	0.362	0.308
	±0.001	±0.000	±0.013	±0.024	±0.050	±0.050	±0.016	±0.005
2	0.917	0.780	0.618	0.463	0.448	0.372	0.365	0.308
	±0.006	±0.013	±0.010	±0.033	±0.007	±0.017	±0.009	±0.013
ers	0.819	0.661	0.540	0.422	0.372	0.358	0.343	0.290
3	±0.022	±0.007	±0.052	±0.012	±0.016	±0.028	±0.004	±0.016
nodifi	0.655	0.536	0.423	0.364	0.334	0.317	0.303	0.279
4	±0.005	±0.013	±0.034	±0.018	±0.004	±0.011	±0.020	±0.013
ffect_r	0.543	0.421	0.400	0.332	0.314	0.314	0.298	0.290
5	±0.015	±0.030	±0.007	±0.015	±0.019	±0.005	±0.002	±0.009
e e	0.461	0.448	0.371	0.354	0.344	0.318	0.291	0.278
	±0.011	±0.011	±0.014	±0.008	±0.020	±0.015	±0.006	±0.001
6	0.366	0.313	0.316	0.285	0.278	0.269	0.262	0.244
	±0.010	±0.006	±0.009	±0.006	±0.007	±0.002	±0.009	±0.005