

1 **Benchmarking Six Causal Inference Methods for ATE**
2 **Estimation: A Comprehensive Simulation Study**

3 Zung-Ru Lin*

University of Pennsylvania, Philadelphia, PA, USA

Email: lin61302@gmail.com

4 April 5, 2025

5 **Abstract**

6 Estimating the average treatment effect (ATE) from observational data is a fundamental
7 challenge in medical and policy research. Numerous methods have been proposed to address
8 this task, ranging from classical interpretable techniques (e.g., propensity score matching and
9 weighting) to more flexible machine learning approaches. However, the relative performance
10 and trade-offs of these methods in realistic settings remain incompletely understood. In this
11 study, we conduct an extensive simulation-based comparison of six commonly used ATE esti-
12 mation approaches, including two propensity score matching strategies, two propensity score
13 weighting schemes (one using logistic regression and one using a boosted decision tree model),
14 and two doubly robust estimation methods (targeted maximum likelihood estimation (TMLE)
15 and its doubly robust extension, DRTMLE). The simulation design emulates a plausible clin-
16 ical scenario with a known true ATE, enabling unbiased evaluation of each method's accuracy,
17 precision, and coverage across repeated trials. Our results reveal substantial discrepancies
18 among the methods. Approaches that incorporate flexible modeling of the outcome or propen-
19 sity score (such as the use of boosting within the propensity score estimation or nonparametric
20 outcome regression in TMLE) tend to achieve lower bias and more stable estimates compared
21 to simpler parametric or matching methods. At the same time, the estimates produced by these
22 advanced techniques remain interpretable as causal effects, which is crucial for real-world
23 decision-making. This study highlights the importance of simulation-based benchmarking to
24 understand method performance and provides practical guidance on selecting appropriate, in-
25 terpretable causal inference tools.

26 **Keywords:** causal inference, average treatment effect, propensity score, matching, doubly robust
27 estimation, targeted maximum likelihood estimation

*This research was conducted independently and was not supported by the author's affiliated institution.

28 **1 Introduction**

29 In medical and epidemiological research, it is crucial to accurately determine whether a treatment
30 truly improves patient outcomes, especially when randomized trials are infeasible or unethical. In
31 previous research, we analyzed observational health data and identified smoking as a dominant
32 risk factor for costly and deadly conditions. Smoking contributes to diseases such as emphy-
33 sema, chronic bronchitis, and lung cancer, which carry a high risk of mortality. Meanwhile, some
34 individuals have unscrupulously claimed that their medications can cure such serious diseases.
35 According to *The Guardian*, “*More than 110 countries have reported more than 2,000 cases of*
36 *bad drugs over WHO’s global surveillance and monitoring system*,” noted Michael Deats.¹ In
37 other words, driven by tremendous potential profit, some actors promote ineffective or harmful
38 treatments—jeopardizing lives and eroding public trust in the medical environment. In such an en-
39 vironment, rigorous evaluation of treatment effectiveness is critical to ensure that only genuinely
40 beneficial interventions are advocated.

41 This tension motivates our research comparing different methodologies to examine treatment
42 effects. We illustrate this approach in the context of a lung cancer scenario, chosen for its high
43 mortality and well-understood risk factors (such as smoking) to ground our simulation in a realis-
44 tic clinical setting. Specifically, we explore the ability to estimate the true ATE using the following
45 methods: nearest neighbor matching on propensity scores, full matching on propensity scores,
46 propensity score weighting by logistic regression, propensity score weighting by a generalized
47 boosted model, targeted maximum likelihood estimation, and doubly robust targeted maximum
48 likelihood estimation. Given the true ATE from simulation, we can compare the accuracy of each
49 method. The outcomes give credence to our research, yielding meaningful insights for practition-
50 ers.

51 In this study, we present a rigorous simulation-based comparison of these six interpretable
52 methods under a unified, realistic scenario. By combining matching, weighting, and doubly robust
53 estimators in one framework, our evaluation highlights each method’s relative strengths, common
54 pitfalls, and biases in causal effect estimation. With the true ATE known from the simulated pop-
55 ulation, we directly assess the bias and variability of each estimator. The results provide insight
56 into how each technique performs under realistic conditions and form the basis for practical rec-
57 ommendations. In summary, by evaluating matching, weighting, and doubly robust strategies side
58 by side, we offer a comprehensive comparison of popular ATE estimation methods and actionable
59 guidance on choosing among them in applied research.

60 **2 Methods**

61 **2.1 Data Generating Process**

62 We simulate data to reflect a real-world scenario of patients suffering from lung cancer. Each
63 patient is monitored for five years to record whether they survive or not. We consider smoking
64 as a strong predictor of death for lung cancer patients, based on prior analysis. However, we aim
65 to quantify the effect of smoking on the odds of death. Our simulation incorporates covariates
66 inspired by those in Luque-Fernandez’s tutorial on targeted maximum likelihood estimation for a
67 binary treatment,⁵ combined with additional factors found in online sources, which are believed to