

Recap of Fall 2023

Fall 2023 started with the goal of showing that $\hat{\theta}_{prop}$ in Equation 1 is (or is not) an optimal estimator.

$$\begin{aligned} \hat{\theta}_{prop} = & \quad (1) \\ & n^{-1} \sum_{i=1}^n E[g_i | X_i] + n^{-1} \sum_{i=1}^n \frac{\delta_{1+}}{\pi_{1+}} (E[g_i | X_i, Y_{1i}] - E[g_i | X_i]) \\ & + n^{-1} \sum_{i=1}^n \frac{\delta_{2+}}{\pi_{2+}} (E[g_i | X_i, Y_{2i}] - E[g_i | X_i]) \\ & + n^{-1} \sum_{i=1}^n \frac{\delta_{11}}{\pi_{11}} (g_i - E[g_i | X_i, Y_{1i}] - E[g_i | X_i, Y_{2i}] + E[g_i | X_i]). \end{aligned}$$

While this estimator performs quite well overall, it is outperformed by $\hat{\theta}_\delta$ (see Table 1) in a simulation where the number of observations in segments A_{10} differed from A_{01} . Studying this estimator led to the creation of a class of estimators in Equation 2 with specific examples in Table 1.

$$\hat{\theta} = \frac{\delta_{11}}{\pi_{11}} g(Z) + \beta_0(\delta, c_0) E[g(Z) | X] + \beta_1(\delta, c_1) E[g(Z) | X, Y_1] + \beta_2(\delta, c_2) E[g(Z) | X, Y_2]. \quad (2)$$

Table 1: Specific examples of estimators from the larger class in Equation 2.

Estimator	$\beta_0(\delta, c_0)$	$\beta_1(\delta, c_1)$	Implemented
$\hat{\theta}_{prop}$	$\left(1 - \frac{(\delta_{10} + \delta_{11})}{(\pi_{10} + \pi_{11})} - \frac{(\delta_{01} + \delta_{11})}{(\pi_{01} + \pi_{11})} + \frac{\delta_{11}}{\pi_{11}}\right)$	$\left(\frac{\delta_{10} + \delta_{11}}{\pi_{10} + \pi_{11}} - \frac{\delta_{11}}{\pi_{11}}\right)$	✓
$\hat{\theta}_{prop}^{ind}$	$\left(1 - \frac{(\delta_{10})}{(\pi_{10})} - \frac{(\delta_{01})}{(\pi_{01})} + \frac{\delta_{11}}{\pi_{11}}\right)$	$\left(\frac{\delta_{10}}{\pi_{10}} - \frac{\delta_{11}}{\pi_{11}}\right)$	✓
$\hat{\theta}_c$	$c_0 \left(1 - \frac{(\delta_{10} + \delta_{11})}{(\pi_{10} + \pi_{11})} - \frac{(\delta_{01} + \delta_{11})}{(\pi_{01} + \pi_{11})} + \frac{\delta_{11}}{\pi_{11}}\right)$	$c_1 \left(\frac{\delta_{10} + \delta_{11}}{\pi_{10} + \pi_{11}} - \frac{\delta_{11}}{\pi_{11}}\right)$	✓
$\hat{\theta}_{c, ind}$	$c_0 \left(1 - \frac{(\delta_{10})}{(\pi_{10})} - \frac{(\delta_{01})}{(\pi_{01})} + \frac{\delta_{11}}{\pi_{11}}\right)$	$c_1 \left(\frac{\delta_{10}}{\pi_{10}} - \frac{\delta_{11}}{\pi_{11}}\right)$	✓
$\hat{\theta}_\delta$	$c_0 \left(\frac{\delta_{11}}{\pi_{11}} - \frac{\delta_{00}}{\pi_{00}}\right)$	$c_1 \left(\frac{\delta_{11}}{\pi_{11}} - \frac{\delta_{10}}{\pi_{10}}\right)$	✓

The challenge with working with this class is the β coefficients. These are functional coefficients and so it is difficult for me to understand how to show that a particular class is optimal. Dispite this difficulty, we were able to show that the proposed estimator is not optimal via simulation. When this estimator was compared to other estimators in Table 2, it did not perform the best when the segments were unbalanced and the estimating function was nonlinear. This is shown in Table 2

Table 2: True Values: $\theta = 10$, $\rho = 0.5$. This simulation assesses the bias and standard deviation (SD) of different estimators of the function $\theta = E[Y_1^2 Y_2]$ where $Y_1 = x + \varepsilon_1$, $Y_2 = \beta + x + \varepsilon_2$, $x \perp (\varepsilon_1, \varepsilon_2)$, and $(\varepsilon_1, \varepsilon_2)$ come from a mean zero bivariate normal distribution with unit variance and covariance of ρ . This simulation uses a sample size of $N = 1000$ and $B = 3000$ Monte Carlo simulations. Each observation has the following independent probabilities of landing into a respective segment: $\pi_{11} = 0.2$, $\pi_{10} = 0.4$, $\pi_{01} = 0.1$, and $\pi_{00} = 0.3$. The columns of Tstat and P-value give the test statistic and p-value from a two sample test for unbiasedness.

Algorithm	Bias	SD	Tstat	P-value
Oracle	0.007	0.529	0.741	0.229
CC	0.036	1.190	1.667	0.048
IPW	0.029	1.372	1.170	0.121
$\hat{\theta}_{prop}$	0.004	0.694	0.292	0.385
$\hat{\theta}_c$	0.007	0.670	0.582	0.280
$\hat{\theta}_\delta$	0.018	0.675	1.478	0.070

Current Work

As discussed in a related note, we are currently trying to understand the loss of efficiency of using a semiparametric estimator with the model is correctly specified. I am having difficulties with the current simulation setup but this work is contained in `efficiencyloss_semi.tex`.

Potential Ideas to Pursue

- Tsiatis (2006) describes a recursive method to get the optimal semiparametric estimator for non-monotone data. This is difficult and confusing. However, we can pose the estimation as the solution to an integral equation. While people know that one can do this, it has not been done. The real work would be estimating the variance.
- Often in simulations, the optimal value for c_i would be 1 or -1 . Since the optimal semiparametric estimator is a double projection (see Chapter 10 of Tsiatis (2006)), I hypothesize that it may be possible to construct a test to see how close to the optimal we are by projecting the β coefficients onto some constants c .