

# Conservative estimation of Difference-in-Means estimator's variance under random assignment\*

Thomas Leavitt

The expression for the variance of the Difference-in-Means estimator,  $\text{Var}_\Omega[\hat{\tau}]$ , is

$$(1) \quad \text{Var}_\Omega[\hat{\tau}] = \frac{1}{n-1} \left( \frac{n_C \sigma_n^2(\mathbf{y}_T)}{n_T} + \frac{n_T \sigma_n^2(\mathbf{y}_C)}{n_C} + 2\sigma_n(\mathbf{y}_C, \mathbf{y}_T) \right),$$

where

$$(2) \quad \sigma_n^2(\mathbf{y}_T) = \left( \frac{1}{n} \right) \sum_{i=1}^n \left( y_{Ti} - \frac{1}{n} \sum_{i=1}^n y_{Ti} \right)^2$$

$$(3) \quad \sigma_n^2(\mathbf{y}_C) = \left( \frac{1}{n} \right) \sum_{i=1}^n \left( y_{Ci} - \frac{1}{n} \sum_{i=1}^n y_{Ci} \right)^2$$

$$(4) \quad \sigma_n(\mathbf{y}_C, \mathbf{y}_T) = \left( \frac{1}{n} \right) \sum_{i=1}^n \left( y_{Ci} - \frac{1}{n} \sum_{i=1}^n y_{Ci} \right) \left( y_{Ti} - \frac{1}{n} \sum_{i=1}^n y_{Ti} \right).$$

Absent additional assumptions, we cannot unbiasedly estimate  $\text{Var}_\Omega[\hat{\tau}]$ . We can unbiasedly estimate the first two terms of  $\text{Var}_\Omega[\hat{\tau}]$  in Equation (1),  $\sigma_n^2(\mathbf{y}_T)$  and  $\sigma_n^2(\mathbf{y}_C)$ , but not the third term,  $\sigma_n(\mathbf{y}_C, \mathbf{y}_T)$ , because no two potential outcomes of the same unit are observable. To conservatively estimate  $\text{Var}_\Omega[\hat{\tau}]$ , we will derive a quantity we can unbiasedly estimate that is always at least as great as the true variance,  $\text{Var}_\Omega[\hat{\tau}]$ .

---

\*This is a live document that is subject to updating at any time.

The Cauchy-Schwarz inequality implies that

$$\sigma_n(\mathbf{y}_C, \mathbf{y}_T) \leq \sqrt{\sigma_n^2(\mathbf{y}_C)\sigma_n^2(\mathbf{y}_T)}$$

and the AM-GM inequality further implies that

$$\sqrt{\sigma_n^2(\mathbf{y}_C)\sigma_n^2(\mathbf{y}_T)} \leq \frac{\sigma_n^2(\mathbf{y}_C) + \sigma_n^2(\mathbf{y}_T)}{2}.$$

Hence, it follows that  $2\sigma_n(\mathbf{y}_C, \mathbf{y}_T) \leq \sigma_n^2(\mathbf{y}_C) + \sigma_n^2(\mathbf{y}_T)$ . Substituting  $\sigma_n^2(\mathbf{y}_C) + \sigma_n^2(\mathbf{y}_T)$  for  $2\sigma_n(\mathbf{y}_C, \mathbf{y}_T)$  therefore yields a tight upper bound for the true variance of the Difference-in-Means estimator:

$$\begin{aligned} \text{Var}_\Omega[\hat{\tau}] &= \frac{1}{n-1} \left( \frac{n_T \sigma_n^2(\mathbf{y}_C)}{n_C} + \frac{n_C \sigma_n^2(\mathbf{y}_T)}{n_T} + 2\sigma_n(\mathbf{y}_C, \mathbf{y}_T) \right) \\ &\leq \frac{1}{n-1} \left( \frac{n_T \sigma_n^2(\mathbf{y}_C)}{n_C} + \frac{n_C \sigma_n^2(\mathbf{y}_T)}{n_T} + \sigma_n^2(\mathbf{y}_C) + \sigma_n^2(\mathbf{y}_T) \right) \\ &= \frac{1}{n-1} \left( \frac{n_T \sigma_n^2(\mathbf{y}_C)}{n_C} + \frac{n \sigma_n^2(\mathbf{y}_T) - n_T \sigma_n^2(\mathbf{y}_T)}{n_T} + \sigma_n^2(\mathbf{y}_C) + \sigma_n^2(\mathbf{y}_T) \right) \\ &= \frac{1}{n-1} \left( \frac{n_T \sigma_n^2(\mathbf{y}_C)}{n_C} + \frac{n \sigma_n^2(\mathbf{y}_T)}{n_T} - \frac{n_T \sigma_n^2(\mathbf{y}_T)}{n_T} + \sigma_n^2(\mathbf{y}_C) + \sigma_n^2(\mathbf{y}_T) \right) \\ &= \frac{1}{n-1} \left( \frac{n_T \sigma_n^2(\mathbf{y}_C)}{n_C} + \frac{n \sigma_n^2(\mathbf{y}_T)}{n_T} - \sigma_n^2(\mathbf{y}_T) + \sigma_n^2(\mathbf{y}_C) + \sigma_n^2(\mathbf{y}_T) \right) \\ &= \frac{1}{n-1} \left( \frac{n_T \sigma_n^2(\mathbf{y}_C)}{n_C} + \frac{n \sigma_n^2(\mathbf{y}_T)}{n_T} + \sigma_n^2(\mathbf{y}_C) \right) \\ &= \frac{1}{n-1} \left( \frac{n_T \sigma_n^2(\mathbf{y}_C)}{n_C} + \frac{n_C \sigma_n^2(\mathbf{y}_C)}{n_C} + \frac{n \sigma_n^2(\mathbf{y}_T)}{n_T} \right) \\ &= \frac{1}{n-1} \left( \frac{n_T \sigma_n^2(\mathbf{y}_C) + n_C \sigma_n^2(\mathbf{y}_C)}{n_C} + \frac{n \sigma_n^2(\mathbf{y}_T)}{n_T} \right) \\ &= \frac{1}{n-1} \left( \frac{n_T \sigma_n^2(\mathbf{y}_C) + n \sigma_n^2(\mathbf{y}_C) - n_T \sigma_n^2(\mathbf{y}_C)}{n_C} + \frac{n \sigma_n^2(\mathbf{y}_T)}{n_T} \right) \\ &= \frac{1}{n-1} \left( \frac{n \sigma_n^2(\mathbf{y}_C)}{n_C} + \frac{n \sigma_n^2(\mathbf{y}_T)}{n_T} \right) \\ &= \frac{n}{n-1} \left( \frac{\sigma_n^2(\mathbf{y}_C)}{n_C} + \frac{\sigma_n^2(\mathbf{y}_T)}{n_T} \right). \end{aligned}$$

We can unbiasedly estimate the two unknown parameters  $\sigma_n^2(\mathbf{y}_C)$  and  $\sigma_n^2(\mathbf{y}_T)$ . Following [Cochran](#)

(1977, Theorem 2.4), unbiased estimators of  $\sigma_n^2(\mathbf{y}_C)$  and  $\sigma_n^2(\mathbf{y}_T)$  are

$$\begin{aligned}\hat{\sigma}_n^2(\mathbf{y}_C) &= \left( \frac{n-1}{n(n_C-1)} \right) \sum_{i=1}^n (1-Z_i) \left( y_{Ci} - \frac{1}{n_C} \sum_{i=1}^n (1-Z_i) Y_i \right)^2 \\ \hat{\sigma}_n^2(\mathbf{y}_T) &= \left( \frac{n-1}{n(n_T-1)} \right) \sum_{i=1}^n Z_i \left( y_{Ti} - \frac{1}{n_T} \sum_{i=1}^n Z_i Y_i \right)^2.\end{aligned}$$

Therefore, the conservative estimator of the variance of the Difference-in-Means estimator is

$$(5) \quad \widehat{\text{Var}}_{\Omega}[\hat{\tau}] = \frac{n}{n-1} \left( \frac{\hat{\sigma}_n^2(\mathbf{y}_T)}{n_T} + \frac{\hat{\sigma}_n^2(\mathbf{y}_C)}{n_C} \right)$$

This estimator is conservative in that  $\text{E}_{\Omega} \left[ \widehat{\text{Var}}_{\Omega}[\hat{\tau}] \right] \geq \text{Var}_{\Omega}[\hat{\tau}]$ . However, the estimator is unbiased when individual causal effects are constant, which implies that  $2\sigma_n(\mathbf{y}_C, \mathbf{y}_T) = \sigma_n^2(\mathbf{y}_C) + \sigma_n^2(\mathbf{y}_T)$ .

We can also see that the conservative variance estimator is unbiased when individual causal effects are constant by examining the alternative expression for the Difference-in-Means estimator's variance given by

$$(6) \quad \text{Var}_{\Omega}[\hat{\tau}] = \frac{S_n^2(\mathbf{y}_T)}{n_T} + \frac{S_n^2(\mathbf{y}_C)}{n_C} - \frac{S_n^2(\boldsymbol{\tau})}{n},$$

where

$$(7) \quad S_n^2(\mathbf{y}_T) = \left( \frac{1}{n-1} \right) \sum_{i=1}^n \left( y_{Ti} - \frac{1}{n} \sum_{i=1}^n y_{Ti} \right)^2$$

$$(8) \quad S_n^2(\mathbf{y}_C) = \left( \frac{1}{n-1} \right) \sum_{i=1}^n \left( y_{Ci} - \frac{1}{n} \sum_{i=1}^n y_{Ci} \right)^2$$

$$(9) \quad S_n^2(\boldsymbol{\tau}) = \left( \frac{1}{n-1} \right) \sum_{i=1}^n \left( \tau_i - \frac{1}{n} \sum_{i=1}^n \tau_i \right)^2.$$

The conservative variance estimator equivalent to Equation (5) is

$$(10) \quad \widehat{\text{Var}}_{\Omega}[\hat{\tau}] = \frac{1}{n_T} \widehat{S}_n^2(\mathbf{y}_T) + \frac{1}{n_C} \widehat{S}_n^2(\mathbf{y}_C),$$

where

$$\widehat{S}_n^2(\mathbf{y}_T) = \left( \frac{1}{n_T-1} \right) \sum_{i=1}^n Z_i \left( y_{Ti} - \frac{1}{n_T} \sum_{i=1}^n Z_i Y_i \right)^2$$

$$\widehat{S}_n^2(\mathbf{y}_C) = \left( \frac{1}{n_C - 1} \right) \sum_{i=1}^n (1 - Z_i) \left( y_{Ci} - \frac{1}{n_C} \sum_{i=1}^n (1 - Z_i) Y_i \right)^2,$$

assuming that  $n_T \geq 2$  and  $n_C \geq 2$ .

This estimator in Equation (10) is conservative in that  $E_{\Omega} \left[ \widehat{\text{Var}}_{\Omega} [\hat{\tau}] \right] \geq \text{Var}_{\Omega} [\hat{\tau}]$ . However, this estimator of the Difference-in-Means estimator's variance is unbiased when individual causal effects are constant, i.e.,  $S_n^2(\boldsymbol{\tau}) = 0$ .

## References

Cochran, W. G. (1977). *Sampling Techniques* (3rd ed.). Hoboken, NJ: John Wiley & Sons. 2