

Notes on Simulating Power

Prepared by Mauricio Cáceres

September 12, 2016

1 Parametric Power

We typically consider the model

$$(1) \quad Y_i = \alpha + \beta T_i + \gamma X_i + \varepsilon_i$$

for some treatment T_i at the individual level and the OLS estimator $\hat{\beta}_{OLS}$ the difference in outcomes for the treatment and control groups. For $H_0 : \beta = 0$, consider the rejection probability function

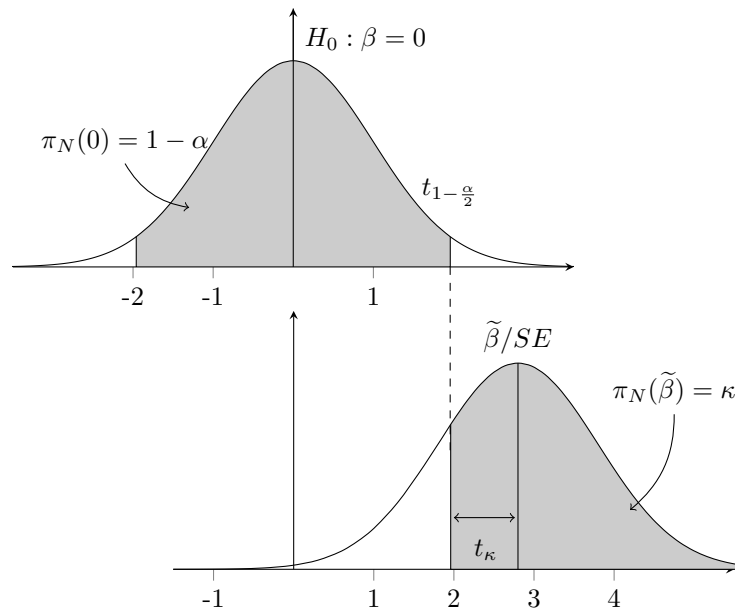
$$(2) \quad \pi_N(\beta) = P\{\text{reject } H_0 | \beta\}$$

For $\beta = 0$, this is α the probability of Type I error or *significance*: How likely are we to make a mistake? For $\beta = \tilde{\beta} \neq 0$ this is *power*, the probability of rejecting the null: How likely are we to get it right? $\hat{\beta}_{OLS}$ is \sqrt{n} -consistent, that is,

$$\sqrt{n}(\hat{\beta}_{OLS} - \beta) \xrightarrow{D} N(0, V_{\hat{\beta}})$$

where β is the true mean. Relying on large-sample asymptotics, we can visualize $\pi_N(\beta)$,

Figure 1.1: Power of a Test



A *parametric* approach to power estimates $\tilde{\beta}$ given α, κ, N, SE and terms it the *minimum detectable effect*, MDE, or it estimates N given α, κ, SE, MDE .

2 Simulated Confidence Interval

However, it is possible to follow a *non-parametric* approach to estimating power. Note there are $C = \binom{N}{NP}$ ways to treat PN individuals. If we estimate $\hat{\beta}_{OLS}$ for each $c = 1, \dots, C$, then we would know the exact distribution of our estimator for the treatment effect under the null given the data. Thus we could compute an exact p -value and determine whether to reject the null.

Even for modestly-sized data, C will be intractably large. Hence we simulate K draws from the possible treatment-control arrangements, T_{ik} such that $\sum_{i=1}^N T_{ik} = PN$, and estimate

$$(3) \quad Y_i = \alpha + \beta_k T_{ik} + \gamma X_i + \varepsilon_i$$

Here $\hat{\beta}_k$ will be distributed around 0 and a $1 - \alpha$ CI under the null is given by

$$(4) \quad \widehat{CI}_{1-\alpha} = \left(\hat{F}^{-1}(\alpha/2), \hat{F}^{-1}(1 - \alpha/2) \right)$$

with $\hat{F}(\hat{\beta}_k)$ the empirical cdf of $\hat{\beta}_k$. This approach is appealing compared to a parametric approach because it naturally takes into account the correlation structure of the errors, whereas a parametric approach requires making an assumption about $V_{\hat{\beta}}$. If we had historical data on our study population (we can think of historical data as data on a population where our treatment had no effect, or the counterfactual of what would happen were we to treat our population with no effect) then we can simulate the CI above, and say that we expect to be able to reject effects outside the confidence interval.

3 Simulated Power

Note, however, that this says nothing about power. In fact, power at either end of the confidence interval should be about 0.5 (if $\hat{\beta}_{OLS}$ were to be symmetrically distributed). Typically we look for a power level of 0.8 or 0.9. We could assume that the true effect of T_i is $\tilde{\beta}$, and estimate

$$(5) \quad \begin{aligned} \tilde{Y}_{ik} &= Y_i + \tilde{\beta} T_{ik} \\ \tilde{Y}_{ik} &= \alpha + \beta_k T_{ik} + \gamma X_i + \varepsilon_i \end{aligned}$$

In this case, $\hat{\beta}_k$ will be distributed around $\tilde{\beta}$ but the shape of the distribution would not have changed. Thus we can estimate power as

$$(6) \quad \hat{\kappa} = \frac{1}{K} \sum_k 1 \left(\hat{\beta}_k \notin \widehat{CI}_{1-\alpha} \right)$$

and we can search for $\tilde{\beta}$ such that $\hat{\kappa} \approx \kappa$ for some desired power level κ (note $\hat{\kappa} \xrightarrow{P} P(\beta_k \notin CI_{1-\alpha}) = \kappa$). That is, we look for a $\tilde{\beta}$ that causes us to reject the null κ portion of the time. This is more complicated if Y_i is binary. The approach works to obtain a CI under the null but the subsequent search does not map trivially. One idea is to randomly swap successes to failures (or the converse) based on $\tilde{\beta}$. Consider

$$(7) \quad \begin{aligned} \tilde{Y}_{ik} &= Y_i(1 - T_{ik}) + (Y_i + S_{ik})T_{ik} = Y_i + S_{ik}T_{ik} \\ \tilde{Y}_{ik} &= \alpha + \beta_k T_{ik} + \gamma X_i + \varepsilon_i \end{aligned}$$

where S_{ik} is constructed as follows

- Let

$$\begin{aligned}
T_k &= \sum_i T_{ik} & S_k &= \sum_i T_{ik} Y_i \\
\tilde{\beta} &\in \left[-\frac{S_k}{T_k}, \frac{T_k - S_k}{T_k} \right] \\
S_k^1 &= \tilde{\beta} T_k \\
S_k^2 &= \begin{cases} \frac{S_k}{T_k} - S_{1k} & \tilde{\beta} > 0 \\ \frac{T_k - S_k}{T_k} - S_{1k} & \tilde{\beta} < 0 \end{cases}
\end{aligned}$$

- Construct the set ς_k with S_k^1 entries equal to $1(\tilde{\beta} > 0) - 1(\tilde{\beta} < 0)$ and S_k^2 entries equal to 0.
- For i such that $T_{ik} = 1$, draw s from ς_k *without* replacement and set $S_{ik} = s$ ($S_{ik} = 0$ otherwise).

Note that $T_k^{-1} \sum_i S_{ik} T_{ik} = \tilde{\beta}$, hence

$$\begin{aligned}
E\hat{\beta}_k &= E \left[\tilde{Y}_{ik} | T_{ik} = 1 \right] - E \left[\tilde{Y}_{ik} | T_{ik} = 0 \right] \\
&= E[Y_i + S_{ik} | T_{ik} = 1] - E[Y_i | T_{ik} = 0] \\
&= EY_i + E[S_{ik} | T_{ik} = 1] - EY_i \\
&= \tilde{\beta}
\end{aligned}$$

So $\hat{\beta}_k$ will be distributed around $\tilde{\beta}$. Now we can outline a general simulation procedure:

1. Estimate a $1 - \alpha$ CI for $\hat{\beta}_{OLS}$ under $H_0 : \beta = 0$ using [Equation \(3\)](#) and [Equation \(4\)](#).
2. Choose a starting MDE, $\tilde{\beta}$, and estimate $\hat{\beta}_k$ for $k = 1, \dots, K$ using [Equation \(5\)](#) if Y_i is continuous or [Equation \(7\)](#) if Y_i is binary.
3. Estimate power using [Equation \(6\)](#).
4. If $\hat{\kappa} < \kappa$ then increase $\tilde{\beta}$; if $\hat{\kappa} > \kappa$ then decrease $\tilde{\beta}$.
5. Continue until $|\hat{\kappa} - \kappa| < \epsilon$ for ϵ small.

4 Monte Carlo Simulations

Consider some data-generating process (DGP) for $w_i = (y_i, x_i, T_i)$,

$$x_i \sim F_x \quad T_i \sim \text{Bernoulli}(P) \quad T_i \perp\!\!\!\perp x_i \quad y_i | x_i, T_i \sim F_y$$

where

$$\beta = E(y_i | T_i = 1, x_i) - E(y_i | T_i = 0, x_i)$$

The aim is to compute the power of testing $\hat{\beta}_{OLS}$ against our simulated confidence interval. For $m = 1, \dots, M$:

1. Generate two draws from the DGP, w_{1m} with $P = 0$ and w_{2m} $P \in (0, 1)$.

2. Compute $\hat{\kappa}_{1m}, \hat{\kappa}_{2m}$ from our power simulation procedure and $\hat{\beta}_{1m}, \hat{\beta}_{2m}$ from OLS.
3. Construct a $1 - \alpha$ confidence interval for β , $\widehat{CI}_{1-\alpha}^M$, using [Equation \(4\)](#) and compute

$$\bar{\kappa}^M = \frac{1}{M} \sum_m 1 \left(\hat{\beta}_{2m} \notin \widehat{CI}_{1-\alpha}^M \right)$$

4. It should be the case that $\hat{\kappa}_{1m}$ are distributed around α and $\hat{\kappa}_{2m}$ are distributed around $\bar{\kappa}^M$.