

Diseño PO

Daniel Czarniewicz

2017

Estrategia de selección

El diseño PO es un diseño que permite probabilidades de inclusión distintas para cada elemento de la población, y puede verse como una generalización del diseño BE . Dada una población $U = \{1; \dots; k; \dots; N\}$, consideremos $\pi_1; \dots; \pi_k; \dots; \pi_N$ valores predeterminados y no necesariamente iguales, tales que $0 < \pi_k \leq 1 \ \forall k \in U$. Luego, sean $\varepsilon_1; \dots; \varepsilon_k; \dots; \varepsilon_N$ iid $Unif(0; 1)$. La muestra se conforma de la siguiente manera:

$$s = \{k : \varepsilon_k < \pi_k; \ k \in U\}$$

El diseño PO está dado por:

$$p(s) = \prod_{k \in S} \pi_k \prod_{k \in U-S} (1 - \pi_k) \ \forall s \in \mathcal{S}_{PO}$$

Probabilidades de inclusión

Como las indicadores son independientes e $I_k \sim Ber(\pi_k) \ \forall k \in U$, las probabilidades de inclusión que induce el diseño son:

$$\begin{aligned} \star P(k \in s) &= \pi_k \ \forall k \in U \\ \star P(k; l \in s) &= \pi_{kl} = \pi_k \pi_l \ \forall k \neq l \in U \end{aligned}$$

Además:

$$\star \Delta_{kl} = \pi_{kl} - \pi_k \pi_l = \begin{cases} 0 & \text{si } k \neq l \\ \pi_k(1 - \pi_l) & \text{si } k = l \end{cases}$$

Elección de los π_k

Un criterio razonable sería elegir los π_k de forma tal que minimicen la variarianza del estimador \hat{t}_π , sujetos a un tamaño de muestra esperado fijo, n . Se debe resolver el siguiente problema de optimización:

$$\begin{aligned} &\min_{\pi_k} \{V_{PO}(\hat{t}_\pi)\} \\ \text{s.a. } &\sum_U \pi_k = n \end{aligned}$$

Como $V_{PO}(\hat{t}_\pi) = \sum_U \left(\frac{1}{\pi_k} - 1 \right) y_k^2 = \sum_U \frac{y - k^2}{\pi_k} - \sum_U y_k^2$ y, además, $\sum_U y_k^2$ y $\sum_U \pi_k = n$ están dadas, el problema es equivalente a resolver:

$$\min_{\pi_k} \left\{ \left(\sum_U \frac{y - k^2}{\pi_k} - \sum_U y_k^2 \right) \left(\sum_U \pi_k \right) \right\} \approx \min_{\pi_k} \left\{ \left(\sum_U \frac{y - k^2}{\pi_k} \right) \left(\sum_U \pi_k \right) \right\}$$

Por la desigualdad de Cauchy-Schwartz:¹

$$\left(\sum_U \frac{y - k^2}{\pi_k} \right) \left(\sum_U \pi_k \right) = \left[\sum_U \left(\frac{y_k}{\sqrt{\pi_k}} \right)^2 \right] \left[\sum_U (\sqrt{\pi_k})^2 \right] \geq \left[\sum_U \left(\frac{y_k}{\sqrt{\pi_k}} (\sqrt{\pi_k}) \right) \right]^2 = \left(\sum_U y_k \right)^2$$

¹Desigualdad de Cauchy-Schwartz:

$$\left(\sum_U x_k z_k \right)^2 \leq \left(\sum_U x_k \right) \left(\sum_U z_k \right) \text{ y la igualdad se cumple } \Leftrightarrow x_k = \lambda z_k \ \forall k \in U, \ \lambda \text{ fijo.}$$

La igualdad se cumple $\Leftrightarrow \frac{y_k}{\sqrt{\pi_k}} = \lambda \sqrt{\pi_k} \Rightarrow \pi_k = \frac{y_k}{\lambda} \quad \forall k \in U$.

Ahora bien, como $\sum_U \pi_k = \sum_U \frac{y_k}{\lambda} = \frac{t_y}{\lambda} = n \Rightarrow \lambda = \frac{t_y}{n}$, por lo que, la mejor elección de los π_k , sujeto a que $0 < \pi_k \leq 1 \quad \forall k \in U$, viene dada por:

$$\star \pi_k = n \frac{y_k}{t_y} = \frac{n y_k}{\sum_U y_k} \quad \forall k \in U$$

Lo anterior no es asequible dado que y_k no es conocido para toda la población. Por lo tanto, se debe trabajar con una variable auxiliar que cumpla que:

- x_k es conocida para toda la población.
- $x_k > 0$ para toda la población.
- $x_k \doteq c y_k$ para toda la población.

En este caso, $\pi_k = n \frac{x_k}{t_x} = \frac{n x_k}{\sum_U x_k} \quad \forall k \in U$

No tener información auxiliar podría pensarse como una situación en la que $x_k = 1 \quad \forall k \in U$, con lo que se obtendría que $\pi_k = \frac{n}{N}$. Esto último es lo que justifica el uso de diseños con probabilidades de inclusión iguales para todos los elementos de la población, cuando no se dispone de información auxiliar.

El estimador \hat{t}_π

$$\begin{aligned} \star \hat{t}_\pi &= \sum_s y_k^\vee = \sum_s \frac{y_k}{\pi_k} \\ \star V_{PO}(\hat{t}_\pi) &= \sum \sum_U \Delta_{kl} y_k^\vee y_l^\vee = \sum_U \pi_k (1 - \pi_k) y_k^{\vee 2} = \sum_U \left(\frac{1}{\pi_k} - 1 \right) y_k^2 \\ \star \hat{V}_{PO}(\hat{t}_\pi) &= \sum_s (1 - \pi_k) \frac{y_k^2}{\pi_k^2} = \sum_s \frac{1}{\pi_k} \left(\frac{1}{\pi_k} - 1 \right) y_k^2 \end{aligned}$$

Si los π_k fueron elegidos de forma óptima, entonces tendremos que:

$$\star \hat{t}_\pi = \sum_s \frac{y_k}{\pi_k} = \sum_s \frac{y_k}{\frac{n y_k}{\sum_U y_k}} = \sum_s \left(\sum_U \frac{y_k}{n} \right) = \left(\sum_U \frac{y_k}{n} \right) \left(\sum_s 1 \right) = \frac{n_S}{n} t_y$$

Con lo que la varianza de \hat{t}_π dependerá únicamente de n_S .

El estimador \hat{t}_{alt}

$$\star \hat{t}_{alt} = N \frac{\sum_s y_k^\vee}{\sum_s \frac{1}{\pi_k}} = N \frac{\hat{t}_\pi}{\hat{N}}$$

Tamaño muestral

El tamaño de muestra en un diseño PO es aleatorio. Dado que $n_S = \sum_U I_k$:

$$\star E_{PO}(n_S) = \sum_U \pi_k \quad \star V_{PO}(n_S) = \sum_U \pi_k (1 - \pi_k)$$