

# Diseño SY

Daniel Czarniewicz

2017

## Estrategia de selección

Se fija un número  $a$  llamado *intervalo de muestreo*, de forma tal que  $n = \lceil N/a \rceil \Rightarrow N = na + c$  con  $0 \leq c < a$ . Luego se sortea  $r \sim Unif(1; \dots; a)$  llamado *arranque aleatorio*. La muestra queda conformada por:

$$s = \{k : k = r + (j-1)a \leq N; \ j = 1; \dots; n_S\}$$

El espacio muestral será:  $\mathcal{S}_{SY} = \{s_1; \dots; s_a\}$  donde  $s_i \cap s_j = \emptyset \ \forall i \neq j$  y  $\bigcup_{i=1}^a s_i = U$ .

El diseño muestral será entonces:

$$p(s) = \begin{cases} 1/a & \text{si } s \in \mathcal{S}_{SY} \\ 0 & \text{en otro caso} \end{cases}$$

Dada la estrategia de selección, las muestras posibles serán:

Muestra	$s_1$	$\dots$	$s_r$	$\dots$	$s_a$
U	$y_1$	$\dots$	$y_k$	$\dots$	$y_a$
	$y_{1+a}$	$\dots$	$y_{r+a}$	$\dots$	$y_{2a}$
	$\vdots$		$\vdots$		$\vdots$
	$y_{1+(n-1)a}$	$\dots$	$y_{r+(n-1)a}$	$\dots$	$y_{na}$
Total	$t_{s_1}$	$\dots$	$t_{s_r}$	$\dots$	$t_{s_a}$
Media	$\bar{y}_{s_1}$	$\dots$	$\bar{y}_{s_r}$	$\dots$	$\bar{y}_{s_a}$

## Tamaño muestral

El tamaño muestral será entonces:

$$n_S = \begin{cases} n+1 & \text{si } 0 < r \leq c \\ n & \text{si } c < r \leq a \end{cases}$$

Con probabilidades:

$$p(n_S) = \begin{cases} P(n_S = n+1) & = P(r \leq c) & = c/a \\ P(n_S = n) & = P(r > c) & = 1 - c/a \end{cases}$$

## Control del tamaño muestral

### ■ Intervalo de muestreo fraccionario:

Sea  $a = N/n$  con  $n \in \mathbb{N}$  el tamaño deseado, y sea  $\varepsilon \sim Unif(0; a)$  el arranque aleatorio. La muestra se forma con las etiquetas  $k$  que cumplen:

$$s = \{k : k-1 < \varepsilon + (j-1)a \leq k; \ j = 1; \dots; n\}$$

De igual forma, si  $r = \varepsilon n \sim Unif(0; n)$ , la muestra se forma por las etiquetas:

$$s = \{k : (k-1)n < r + (j-1)N \leq kn; \ j = 1; \dots; n\}$$

Las probabilidades de inclusión serán:  $\pi_k = P(k \in s) = n \frac{1}{N} = \frac{n}{N} = \frac{1}{a} \ \forall k \in U$

■ **Muestreo sistemático circular**

La población se ordena de forma circular, de forma que al elemento  $N$  le sigue el elemento 1. Se sortea  $r \sim Unif(1; N)$ . Sea  $a = \lceil N/n \rceil$ . La muestra será:

$$S = \left\{ k : k = \begin{cases} r + (j-1)a & \text{si } r + (j-1)a \leq N \\ r + (j-1)a - N & \text{si } r + (j-1)a > N \end{cases} ; j = 1; \dots; n \right\}$$

## Probabilidades de inclusión

$$\star \pi_k = P(k \in s) = P(\text{"seleccionar la muestra } s_r") = 1/a \quad \forall k \in U$$

$$\star \pi_{kl} = P(k, l \in s) = \begin{cases} 1/a & \text{si } k, l \in s_r \in \mathcal{S}_{SY} \\ 0 & \text{en otro caso} \end{cases}$$

Por lo tanto, el diseño no es medible por lo que no existen estimadores insesgados de  $V_{SY}(\hat{t}_\pi)$ .

## El estimador $\hat{t}_\pi$

$$\star t_y = \sum_U y_k = \sum_{r=1}^a t_{s_r} = \sum_{r=1}^a \sum_{s_r} y_k$$

$$\star \hat{t}_\pi = \sum_s y_k^\vee = \sum_s \frac{y_k}{1/a} = a \sum_s y_k = a t_s = \frac{N}{n} t_s = N \bar{y}_s$$

Esto puede pensarse como si tuviéramos una población  $U_t = \{t_{s_1}; \dots; t_{s_r}; \dots; t_{s_a}\}$  y se quiere estimar  $t = \sum_{U_t} t_{s_r}$ . Se toma una muestra bajo un diseño  $SI$  de tamaño  $n = 1$ , así  $\hat{t}_\pi = t_s^\vee = \frac{t_s}{1/a} = a t_s$ .

$$\star Rec(\hat{t}_\pi) = \{a t_{s_1}; \dots; a t_{s_a}\}$$

$$\star P(\hat{t}_\pi = a t_{s_r}) = 1/a \quad \forall r = 1; \dots; a$$

$$= a \sum_U E_{SY}(I_k) y_k = a \sum_U \frac{1}{a} y_k = \sum_U y_k = t_y$$

$$\star E_{SY}(\hat{t}_\pi) = E_{SY}(a t_s) = a E_{SY}(t_s) = a E_{SY}\left(\sum_s y_k\right) = a E_{SY}\left(\sum_U I_k y_k\right) =$$

$$= a \sum_U E(I_k) y_k = a \sum_U \frac{1}{a} y_k = \sum_U y_k = t_y$$

$$\star E_{SY}(\hat{t}_\pi) = \sum_{r=1}^a \frac{\hat{t}_\pi}{a} = \sum_{r=1}^a a \frac{t_{s_r}}{a} = \sum_{r=1}^a t_{s_r} = t_y$$

$$\star V_{SY}(\hat{t}_\pi) = \sum \sum_U \Delta_{kl} y_k^\vee y_l^\vee = \sum \sum_U (\pi_{kl} - \pi_k \pi_l) \frac{y_k}{\pi_k} \frac{y_l}{\pi_l} =$$

$$= \sum \sum_U \left( \frac{\pi_{kl}}{\pi_k \pi_l} \right) y_k y_l - \sum \sum_U y_k y_l =$$

Luego como  $\pi_{kl} = 0$  si  $k$  y  $l$  no pertenecen a la misma muestra,  $\sum \sum_U$  puede cambiarse por  $\sum_{r=1}^a \left( \sum \sum_{s_r} \right)$ .

$$= \sum_{r=1}^a \left( \sum \sum_{s_r} a y_k y_l \right) - \underbrace{\left( \sum_U y_k \right)^2}_{t_y} = a \sum_{r=1}^a \left( \underbrace{\sum_{s_r} y_k}_{t_{s_r}} \right)^2 =$$

$$\begin{aligned}
&= a \sum_{r=1}^a t_{s_r}^2 - (a\bar{t})^2 = a \left( \sum_{r=1}^a t_{s_r}^2 - a\bar{t}^2 \right) = a \sum_{r=1}^a (t_{s_r} - \bar{t})^2 = a \left( \frac{a-1}{a-1} \right) \sum_{r=1}^a (t_{s_r} - \bar{t})^2 = \\
&= a(a-1) \underbrace{\frac{1}{a-1} \sum_{r=1}^a (t_{s_r} - \bar{t})^2}_{S_t^2} = a(a-1)S_t^2
\end{aligned}$$

Si  $t_{s_1} = t_{s_2} = \dots = t_{s_a} = t/a \Rightarrow V_{SY}(\hat{t}_\pi) = 0$ .  $V_{SY}(\hat{t}_\pi)$  depende de cómo se ordene la población.

$$\begin{aligned}
\star V_{SY}(\hat{t}_\pi) &= a \sum_{r=1}^m (t_{s_r} - \bar{t})^2 = a \sum_{r=1}^m \left( n\bar{y}_{s_r} - \frac{t}{a} \right)^2 = \\
&= a \sum_{r=1}^m \left( \frac{N}{a} \bar{y}_{s_r} - \frac{N}{a} \bar{y}_U \right)^2 = \frac{N^2}{a} \sum_{r=1}^m (\bar{y}_{s_r} - \bar{y}_U)^2 = N n \sum_{r=1}^m (\bar{y}_{s_r} - \bar{y}_U)^2 = N \times SSB
\end{aligned}$$

Como  $SST$  es fija, aumentar  $SSW$  implica disminuir  $SSB$ . Conviene que cada una de las muestras posibles sean muy heterogéneas, de forma de que  $SSW$  sea grande. Medimos la homogeneidad mediante:

$$\star \delta = 1 - \frac{N-1}{N-a} \frac{SSW}{SST} = 1 - \frac{S_{yw}^2}{S_{yu}^2} \text{ donde } S_{yw}^2 = \frac{SSW}{N-a} \text{ y } S_{yu}^2 = \frac{SST}{N-1}$$

- $\delta_{\text{máx}} = 1 \Leftrightarrow S_{yw}^2 = 0 \Rightarrow SSW = 0 \Rightarrow$  los grupos son lo más homogéneos posible  $\Rightarrow V_{SY}(\hat{t}_\pi)$  es la mayor posible.
- $\delta_{\text{mín}} = -\frac{a-1}{N-a} \Leftrightarrow SSW = SST \Rightarrow SSB = 0 \Rightarrow$  los grupos son lo más heterogéneos posible  $\Rightarrow V_{SY}(\hat{t}_\pi) = 0$ .

## Efecto diseño

Una forma alternativa de escribir la varianza del estimador  $\pi$  es:

$$\begin{aligned}
\star V_{SY}(\hat{t}_\pi) &= N \times SSB = N(SST - SSW) = \\
&= N \left[ (N-1)S_{yu}^2 - SSW \frac{SST}{SST} \frac{N-a}{N-a} \right] = N \left[ (N-1)S_{yu}^2 - SSW \frac{(N-1)S_{yu}^2}{SST} \frac{N-a}{N-a} \right] = \\
&= N \left[ (N-1)S_{yu}^2 - (N-a)S_{yu}^2 \underbrace{\frac{SSW}{SST} \frac{N-1}{N-a}}_{1-\delta} \right] = N \left[ (N-1)S_{yu}^2 + (\delta-1) \underbrace{(N-a)S_{yu}^2}_{N-\frac{N}{n}} \right] = \\
&= N \left[ (N-1)S_{yu}^2 + (\delta-1) \underbrace{\left( N - \frac{N}{n} \right) S_{yu}^2}_{\frac{N}{n}(n-1)} \right] = N \left[ (N-1)S_{yu}^2 + (\delta-1) \left( \frac{N}{n}(n-1) \right) S_{yu}^2 \right] = \\
&= N \left[ (N-1)S_{yu}^2 + \delta \frac{N}{n}(n-1)S_{yu}^2 - \frac{N}{n}(n-1)S_{yu}^2 \right] = \\
&= N \left[ \left( N-1 - \frac{N}{n}(n-1) \right) S_{yu}^2 + \delta \frac{N}{n}(n-1)S_{yu}^2 \right] = \\
&= N \left[ \left( N-1 - N + \frac{N}{n} \right) S_{yu}^2 + \delta \frac{N}{n}(n-1)S_{yu}^2 \right] =
\end{aligned}$$

$$\begin{aligned}
&= N \left[ \frac{N}{n} \left( 1 - \frac{n}{N} \right) S_{y_U}^2 + \delta \frac{N}{n} (n-1) S_{y_U}^2 \right] = \\
&= N \left[ \frac{N}{n} (1-f) S_{y_U}^2 + \delta \frac{N}{n} (n-1) S_{y_U}^2 \right] = \\
&= \frac{N^2}{n} S_{y_U}^2 \left[ (1-f) + \delta(n-1) \right] \text{ con } f = \frac{n}{N} = \frac{1}{a} \\
\star \text{Def}f(SY; \hat{t}_\pi) &= \frac{\frac{N^2}{n} S_{y_U}^2 \left[ (1-f) + \delta(n-1) \right]}{\frac{N^2}{n} (1-f) S_{y_U}^2} = \frac{(1-f) + \delta(n-1)}{1-f} = 1 + \frac{\delta(n-1)}{1-f}
\end{aligned}$$

$$\frac{n-1}{1-f} > 0 \Rightarrow \begin{cases} \text{si } \delta > 0 & \Rightarrow SI \text{ es más eficiente que } SY \\ \text{si } \delta = 0 & \Rightarrow SI \text{ y } SY \text{ son igualmente eficientes} \\ \text{si } \delta < 0 & \Rightarrow SY \text{ es más eficiente que } SI \end{cases}$$

$\delta < 0 \Leftrightarrow S_{y_W}^2 > S_{y_U}^2$ , esto ocurre cuando los grupos son suficientemente heterogéneos.

## Estimación de la varianza

Dado que no existe un estimador insesgado para  $V_{SY}(\hat{t}_\pi)$ , se emplean las siguientes tácticas:

1. Si existe razón para creer que  $V_{SY} \leq V_{SI}$ , se utiliza la varianza del  $SI$ , por lo que:

$$\hat{V}_{SY}(\hat{t}_\pi) = \frac{N^2}{n} (1-f) S_{y_{s_r}}^2 \text{ donde } S_{y_{s_r}}^2 = \frac{1}{n-1} \sum_{s_r} (y_k - \bar{y}_{s_r})^2$$

2. Tomar más de un arranque aleatorio,  $m > 1$ , con intervalo de muestreo  $ma$ . Ahora cada  $SY$  contribuye una fracción  $n/m$  de la muestra.

Sean  $r_1; \dots; r_m$  los diferentes arranques y, por simplicidad, se supone que  $n/m$  y  $a$  son enteros. La muestra será:

$$s = \{k : k = r_i + (j-1)ma; \ i = 1; \dots; m; \ j = 1; \dots; n/m\}$$

Las probabilidades de inclusión serán:

$$\begin{aligned}
\star \pi_k &= \frac{m}{ma} = \frac{n}{N} \quad \forall k \in U \\
\star \pi_{kl} &= \begin{cases} \frac{n}{N} & \text{si } k, l \text{ pertenecen a la misma muestra} \\ \frac{n}{N} \frac{m-1}{ma-1} & \text{si } k, l \text{ pertenecen a distintas muestras} \end{cases}
\end{aligned}$$