

# Muestreo en dos etapas

Daniel Czarniewicz

2017

En un muestreo en dos etapas la población  $U = \{1; \dots; k; \dots; N\}$  es particionada en  $N_I$  grupos o conglomerados llamados *PSU* (primary sampling units), de forma tal que:  $U_I = \{U_1; \dots; U_i; \dots; U_{N_I}\}$ . Cada conglomerado  $U_i$  es de tamaño  $N_i$ . Luego entonces  $N = \sum_{U_I} N_i$ .

## Mecanismo de selección

En una primera etapa se toma una muestra  $S_I$  de  $U_I$  según el diseño  $p_I(\cdot)$ . El número de *PSUs* seleccionado lo anotamos como  $n_{S_I}$  si el diseño de primera etapa es de tamaño aleatorio, o  $n_I$  si el diseño de primera etapa es de tamaño fijo.

Luego, en una segunda etapa, para cada uno de los  $i$  conglomerados seleccionados en la primera etapa, se toma una muestra  $S_i$  de  $U_i$  según el diseño  $p_i(\cdot|s_I)$ . El número de elementos en  $S_i$  se anota como  $n_{S_i}$  si el diseño de segunda etapa es de tamaño aleatorio, o  $n_i$  si el diseño de segunda etapa es de tamaño fijo. A los elementos seleccionados en la segunda etapa se los llama *SSU* (secondary sampling units).

Como resultado se obtiene una muestra  $s$  de  $U$  tal que  $s = \bigcup_{i \in S_I} s_i$ . El número total de elementos seleccionados será  $n_s = \sum_{i \in S_I} n_i$ .

## Invarianza e Independencia

En lo que sigue se asume que se cumplen las siguientes propiedades:

1. **Invarianza:**  $\forall i \in U_I$  y  $\forall s_I \ni i$ , se tiene que  $p_i(\cdot|s_I) = p_i(\cdot)$ . Esto implica que, sea cual sea la muestra de *PSU* en la primera etapa, si sale el conglomerado  $i$ , el diseño de muestreo de segunda etapa para dicho conglomerado será siempre  $p_i(\cdot)$ .
2. **Independencia:**  $p\left(\bigcup_{i \in s_I} s_i | s_I\right) = p(s|s_I) = \prod_{i \in s_I} p_i(\cdot|s_I)$ . Esto implica que el diseño de muestreo llevado a cabo en una *PSU* es independiente del llevado a cabo en cualquier otra *PSU*.

Asumiendo que se cumplen las propiedades de invarianza e independencia, y asumiendo que las *SSU* son elementos y no clusters, tenemos que:

$$p(s) = p\left(\bigcup_{i \in s_I} s_i | s_I\right) = p(s|s_I) = \prod_{i \in s_I} p_i(s_i|s_I) = \prod_{i \in s_I} p_i(\cdot)$$

## Probabilidades de inclusión

### Probabilidades de inclusión para la primera etapa

$$\begin{aligned} \star \pi_{I_i} &= P(\text{"el cluster } i \text{ fue seleccionado en la primera etapa"}) = P(i \in s_I) = P(U_i \in s_I) \\ \star \pi_{I_{ij}} &= P(\text{"los clusters } i \text{ y } j \text{ fueron seleccionados en la primera etapa"}) = \\ &= P(i, j \in s_I) = P(U_i, U_j \in s_I) \end{aligned}$$

Luego entonces:

$$\star \Delta_{I_{ij}} = \begin{cases} \pi_{I_{ij}} - \pi_{I_i} \pi_{I_j} & \text{si } i \neq j \\ \pi_{I_i}(1 - \pi_{I_i}) & \text{si } i = j \end{cases} \quad \star \Delta_{I_{ij}}^{\checkmark} = \frac{\Delta_{I_{ij}}}{\pi_{I_{ij}}}$$

## Probabilidades de inclusión para la segunda etapa

$$\star \pi_{k|i} = P(\text{"seleccionar el elemento } k, \text{ dado que se seleccionó el cluster } i") = P(k \in s|i \in s_I)$$

$$\star \pi_{kl|i} = P(\text{"seleccionar los elementos } k \text{ y } l, \text{ dado que se seleccionó el cluster } i") = P(k, l \in s|i \in s_I)$$

Luego entonces:

$$\star \Delta_{kl|i} = \begin{cases} \pi_{kl|i} - \pi_{k|i} \pi_{l|i} & \text{si } k \neq l \\ \pi_{k|i}(1 - \pi_{k|i}) & \text{si } k = l \end{cases} \quad \star \Delta_{kl|i}^{\checkmark} = \frac{\Delta_{kl|i}}{\pi_{kl|i}}$$

## Probabilidades de inclusión de los elementos

De las propiedades de invarianza e independencia se desprende que:

$$P(k \in s) = P(i \in s_I \cap k \in s_i) = P(k \in s_i|i \in s_I) P(i \in s_I) = \pi_{I_i} \pi_{k|i}$$

Por lo tanto

$$\star \pi_k = \pi_{I_i} \pi_{k|i} \quad \forall k \in U_i$$

Luego entonces

$$\begin{aligned} P(k, l \in s) &= P(i, j \in s_I \cap k \in s_i \cap l \in s_j) = \\ &= P(k \in s_i|i, j \in s_I) P(l \in s_j|i, j \in s_I) P(i, j \in s_I) = \pi_{I_{ij}} \pi_{k|i} \pi_{l|j} \end{aligned}$$

Por lo tanto:

$$\star \pi_{kl} = \begin{cases} \pi_{I_i} \pi_{k|i} & \text{si } k = l \in U_i \\ \pi_{I_i} \pi_{kl|i} & \text{si } k, l \in U_i \\ \pi_{I_{ij}} \pi_{k|i} \pi_{l|j} & \text{si } k \in U_i \text{ y } l \in U_j \end{cases}$$

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

$$\begin{aligned} \star \hat{t}_\pi &= \sum_s y_k^{\checkmark} = \sum_{s_I} \sum_{s_i} \frac{y_{k|i}}{\pi_k} = \sum_{s_I} \sum_{s_i} \frac{y_{k|i}}{\pi_{I_i} \pi_{k|i}} = \sum_{s_I} \frac{1}{\pi_{I_i}} \sum_{s_i} \frac{y_{k|i}}{\pi_{k|i}} = \\ &= \sum_{s_I} \frac{1}{\pi_{I_i}} \sum_{s_i} y_{k|i}^{\checkmark} = \sum_{s_I} \frac{\hat{t}_{\pi_i}}{\pi_{I_i}} \end{aligned}$$

$$\text{donde } \hat{t}_{\pi_i} = \sum_{s_i} y_{k|i}^{\checkmark} \text{ estima } t_{y_i} = \sum_{U_i} y_{k|i}$$

$$\star E_{p_i(\cdot|s_I)}(\hat{t}_{\pi_i}) = E_{p_i(\cdot|s_I)}\left(\sum_{s_i} y_{k|i}^{\checkmark}\right) = \sum_{U_i} E_{p_i(\cdot|s_I)}(I_{k|i}) \frac{y_{k|i}}{\pi_{k|i}} = \sum_{U_i} \pi_{k|i} \frac{y_{k|i}}{\pi_{k|i}} = \sum_{U_i} y_{k|i} = t_{y_i}$$

$$\star V_i(\hat{t}_{\pi_i}) = \sum \sum_{U_i} \Delta_{kl|i} y_{k|i}^{\checkmark} y_{l|i}^{\checkmark}$$

$$\star \hat{V}_i(\hat{t}_{\pi_i}) = \sum \sum_{s_i} \Delta_{kl|i}^{\checkmark} y_{k|i}^{\checkmark} y_{l|i}^{\checkmark}$$

Si  $p_i(\cdot)$  es de tamaño fijo:

$$\star V_i(\hat{t}_{\pi_i}) = -\frac{1}{2} \sum \sum_{U_i} \Delta_{kl|i} \left(y_{k|i}^{\checkmark} - y_{l|i}^{\checkmark}\right)^2$$

$$\star \hat{V}_i(\hat{t}_{\pi_i}) = -\frac{1}{2} \sum \sum_{s_i} \Delta_{kl|i}^{\checkmark} \left(y_{k|i}^{\checkmark} - y_{l|i}^{\checkmark}\right)^2$$

$$\begin{aligned}
\star V_{2ST}(\hat{t}_\pi) &= V_{p_I(\cdot)} \left[ \underbrace{E_{p_i(\cdot|s_I)}(\hat{t}_\pi|s_I)}_{\sum_{s_I} t_{y_i}^\vee} \right] + E_{p_I(\cdot)} \left[ \underbrace{V_{p_i(\cdot|s_I)}(\hat{t}_\pi|s_I)}_{\sum_{s_I} \frac{V_i}{\pi_{I_i}^2}} \right] = \\
&= \underbrace{V_{p_I(\cdot)} \left( \sum_{s_I} t_{y_i}^\vee \right)}_{\text{varianza del estimador } \pi \text{ de } t_{y_i} \text{ dado } s_I} + \underbrace{E_{p_I(\cdot)} \left( \sum_{s_I} \frac{V_i}{\pi_{I_i}^2} \right)}_{\text{valor esperado de la varianza del estimador } \pi \text{ de } t_{y_i} \text{ dado } s_I} = \sum \sum_{U_I} \Delta_{I_{ij}} t_{y_i}^\vee t_{y_j}^\vee + \frac{1}{\pi_{I_i}^2} \sum_{U_I} E_{p_I(\cdot)}(I_i) V_i = \\
&= \sum \sum_{U_I} \Delta_{I_{ij}} t_{y_i}^\vee t_{y_j}^\vee + \frac{1}{\pi_{I_i}^2} \sum_{U_I} \pi_{I_i} V_i = \sum \sum_{U_I} \Delta_{I_{ij}} t_{y_i}^\vee t_{y_j}^\vee + \sum_{U_I} \frac{V_i}{\pi_{I_i}}
\end{aligned}$$

Debido a la propiedad de independencia:

$$\begin{aligned}
E(\hat{t}_{\pi_i} \hat{t}_{\pi_j}) &= \begin{cases} t_{y_i}^2 + V_i & \text{si } i = j \\ t_{y_i} t_{y_j} & \text{si } i \neq j \end{cases} \\
\star V_{PSU} &= \sum \sum_{U_I} \Delta_{I_{ij}} \hat{t}_{\pi_i}^\vee \hat{t}_{\pi_j}^\vee \\
\star \hat{V}_{PSU} &= \sum \sum_{U_I} \Delta_{I_{ij}}^\vee \hat{t}_{\pi_i}^\vee \hat{t}_{\pi_j}^\vee - \sum_{s_I} \frac{1}{\pi_{I_i}} \left( \frac{1}{\pi_{I_i}} - 1 \right) \hat{V}_i \\
\star E_{2ST} \left( \sum \sum_{s_I} \Delta_{I_{ij}}^\vee \hat{t}_{\pi_i}^\vee \hat{t}_{\pi_j}^\vee \right) &= E_{2ST} \left( \sum \sum_{s_I} \Delta_{I_{ij}}^\vee \frac{\hat{t}_{\pi_i} \hat{t}_{\pi_j}}{\pi_{I_i} \pi_{I_j}} \right) = \\
&= E_{p_I(\cdot)} \left[ E_{p_i(\cdot|s_I)} \left( \sum \sum_{s_I} \Delta_{I_{ij}}^\vee \frac{\hat{t}_{\pi_i} \hat{t}_{\pi_j}}{\pi_{I_i} \pi_{I_j}} \middle| s_I \right) \right] = E_{p_I(\cdot)} \left[ \sum \sum_{s_I} \Delta_{I_{ij}}^\vee \frac{1}{\pi_{I_i} \pi_{I_j}} E_{p_i(\cdot|s_I)}(\hat{t}_{\pi_i} \hat{t}_{\pi_j} | s_I) \right] = \\
&= E_{p_I(\cdot)} \left( \sum \sum_{s_I} \Delta_{I_{ij}}^\vee \frac{t_{y_i}}{\pi_{I_i}} \frac{t_{y_j}}{\pi_{I_j}} \right) + E_{p_I(\cdot)} \left( \sum_{s_I} \Delta_{I_{ij}}^\vee \frac{V_i}{\pi_{I_i}^2} \right) = \\
&= \sum \sum_{s_I} \Delta_{I_{ij}} \frac{t_{y_i}}{\pi_{I_i}} \frac{t_{y_j}}{\pi_{I_j}} + \sum_{s_I} \left( \frac{1}{\pi_{I_i}} - 1 \right) V_i = V_{PSU} + \sum_{s_I} \left( \frac{1}{\pi_{I_i}} - 1 \right) V_i \\
\star E_{2ST} \left[ - \sum_{s_I} \frac{1}{\pi_{I_i}} \left( \frac{1}{\pi_{I_i}} - 1 \right) \hat{V}_i \right] &= -E_{p_I(\cdot)} \left[ E_{p_i(\cdot|s_I)} \left( \sum_{s_I} \frac{1}{\pi_{I_i}} \left( \frac{1}{\pi_{I_i}} - 1 \right) \hat{V}_i \middle| s_I \right) \right] = \\
&= -E_{p_I(\cdot)} \left[ \sum_{s_I} \frac{1}{\pi_{I_i}} \left( \frac{1}{\pi_{I_i}} - 1 \right) E_{p_i(\cdot|s_I)}(\hat{V}_i | s_I) \right] = -E_{p_I(\cdot)} \left[ \sum_{s_I} \frac{1}{\pi_{I_i}} \left( \frac{1}{\pi_{I_i}} - 1 \right) V_i \right] = \\
&= - \sum_{U_I} \left( \frac{1}{\pi_{I_i}} - 1 \right) V_i
\end{aligned}$$

Por lo tanto:

$$\begin{aligned}
E(\hat{V}_{PSU}) &= V_{PSU} + \sum_{s_I} \left( \frac{1}{\pi_{I_i}} - 1 \right) V_i - \sum_{U_I} \left( \frac{1}{\pi_{I_i}} - 1 \right) V_i = V_{PSU} \\
\star E_{2ST}(\hat{V}_{SSU}) &= E_{2ST} \left( \sum_{s_I} \frac{\hat{V}_i}{\pi_{I_i}^2} \right) = E_{p_I(\cdot)} \left[ E_{p_i(\cdot|s_I)} \left( \sum_{s_I} \frac{\hat{V}_i}{\pi_{I_i}^2} \middle| s_I \right) \right] = \\
&= E_{p_I(\cdot)} \left[ \sum_{s_I} \frac{1}{\pi_{I_i}^2} E_{p_i(\cdot|s_I)}(\hat{V}_i | s_I) \right] = E_{p_I(\cdot)} \left[ \sum_{s_I} \frac{1}{\pi_{I_i}^2} V_i \right] = \sum_{U_I} \frac{V_i}{\pi_{I_i}} = V_{SSU}
\end{aligned}$$

En conclusión:

$$E_{2ST}(\hat{V}_{2ST}(\hat{t}_\pi)) = E_{2ST}(\hat{V}_{PSU} + \hat{V}_{SSU}) = E_{2ST}(\hat{V}_{PSU}) + E_{2ST}(\hat{V}_{SSU}) = V_{PSU} + V_{SSU} = V_{2ST}(\hat{t}_\pi)$$

Un estimador computacionalmente más sencillo viene dado por:

$$\begin{aligned} \star \hat{V}^* &= \sum \sum_{s_I} \Delta_{I_{ij}}^\check \frac{\hat{t}_{\pi_i}}{\pi_{I_i}} \frac{\hat{t}_{\pi_j}}{\pi_{I_j}} \\ \star E_{2ST}(\hat{V}^*) &= V_{PSU} + \sum_{U_I} \left( \frac{1}{\pi_{I_i}} - 1 \right) V_i = V_{PSU} + \underbrace{\sum_{U_I} \frac{V_i}{\pi_{I_i}}}_{V_{SSU}} - \sum_{U_I} V_i = \\ &= \underbrace{V_{PSU} + V_{SSU}}_{V_{2ST}(\hat{t}_\pi)} - \sum_{U_I} V_i = V_{2ST}(\hat{t}_\pi) - \sum_{U_I} V_i \end{aligned}$$

Por lo tanto, el sesgo de  $\hat{V}^* = - \sum_{U_I} V_i$ , por lo que su sesgo relativo está dado por:

$$\star \frac{B(\hat{V}^*)}{V_{2ST}(\hat{t}_\pi)} = - \frac{- \sum_{U_I} V_i}{\sum \sum_{U_I} \Delta_{I_{ij}}^\check t_{y_i}^\check t_{y_j}^\check + \sum_{U_I} \frac{V_i}{\pi_{I_i}}}$$