

Muestreo por clusters

Daniel Czarniewicz

2017

En el muestreo por clusters el problema principal es que no se tiene ni se puede construir un marco muestral que permita realizar muestreo directo de elementos. La población $U = \{1; \dots; k; \dots; N\}$ se particiona en N_I subpoblaciones (clusters), $\{U_1; \dots; U_i; \dots; U_{N_I}\}$, donde $U = \bigcup_{i \in U_I} U_i$. El set de clusters lo representamos como

$$U_I = \{1; \dots; i; \dots; N_I\}.$$

Llamamos N_i a la cantidad de elementos poblacionales en el i -ésimo cluster. De esta forma,

$$N = \sum_{i \in U_I} N_i = \sum_{U_I} N_i$$

Estrategia de selección

Una muestra s_I se selecciona de U_I de acuerdo con el diseño $p_I(\cdot)$. Es decir, se seleccionan clusters. Luego todos los elementos dentro de los clusters seleccionados son relevados. Es decir, se censa los clusters. La muestra queda entonces conformada por $s = \bigcup_{i \in s_I} U_i$. El tamaño de s_I se denota como: n_I si la muestra es de tamaño fijo, o n_{s_I} si la muestra es de tamaño aleatorio. Téngase presente que si N_i varía, entonces n_s variará. Dado que los clusters seleccionados son censados, $n_s = \sum_{s_I} N_i$.

Probabilidades de inclusión

Probabilidades de seleccionar los distintos clusters:

$$\star \pi_{I_i} = P(\text{"seleccionar el cluster } i") = P(i \in s_I) = \sum_{s_I \ni i} p_I(s_I)$$

$$\star \pi_{I_{ij}} = P(\text{"seleccionar los clusters } i \text{ y } j") = \begin{cases} P(i; j \in s_I) = \sum_{s_I \ni i; j} p_I(s_I) & \text{si } i \neq j \\ \pi_{I_{ii}} = \pi_{I_i} & \text{si } i = j \end{cases}$$

$$\star \Delta_{I_{ij}} = \pi_{I_{ij}} - \pi_{I_i} \pi_{I_j} \quad \star \Delta_{I_{ij}}^\vee = \frac{\Delta_{I_{ij}}}{\pi_{I_{ij}}}$$

Probabilidades de inclusión de las unidades poblacionales:

$$\star \pi_k = P(k \in s) = P(i \in s_I) = \pi_{I_i}$$

$$\star \pi_{kl} = P(k; l \in s) = \begin{cases} P(i \in s_I) & = \pi_{I_i} & \text{si } k; l \in U_i \\ P(i; j \in s_I) & = \pi_{I_{ij}} & \text{si } k \in U_i; l \in U_j \end{cases}$$

$$\star \pi_{kk} = \pi_k$$

El estimador \hat{t}_π

$$\begin{aligned}
 \star \quad t_y &= \sum_U y_k = \sum_{U_I} t_{y_i} \\
 \star \quad \hat{t}_\pi &= \sum_{s_I} t_{y_i}^\vee = \sum_{s_I} \frac{t_{y_i}}{\pi_{I_i}} \\
 \star \quad E(\hat{t}_\pi) &= E\left(\sum_{s_I} t_{y_i}^\vee\right) = \sum_{U_I} E(I_k) \frac{t_{y_i}}{\pi_{I_i}} = \sum_{U_I} \pi_k \frac{t_{y_i}}{\pi_{I_i}} = \sum_{U_I} \pi_{I_i} \frac{t_{y_i}}{\pi_{I_i}} = \sum_{U_I} t_{y_i} = t_y \\
 \star \quad V_{p_I(s_I)}(\hat{t}_\pi) &= \sum \sum_{U_I} \Delta_{I_{ij}} t_{y_i}^\vee t_{y_j}^\vee \\
 \star \quad \hat{V}_{p_I(s_I)}(\hat{t}_\pi) &= \sum \sum_{s_I} \Delta_{I_{ij}}^\vee t_{y_i}^\vee t_{y_j}^\vee
 \end{aligned}$$

Si $p_I(\cdot)$ es de tamaño fijo, entonces se cumple que:

$$\begin{aligned}
 \star \quad V_{p_I(s_I)}(\hat{t}_\pi) &= -\frac{1}{2} \sum \sum_{U_I} \Delta_{I_{ij}} \left(t_{y_i}^\vee - t_{y_j}^\vee\right)^2 \\
 \star \quad \hat{V}_{p_I(s_I)}(\hat{t}_\pi) &= -\frac{1}{2} \sum \sum_{s_I} \Delta_{I_{ij}}^\vee \left(t_{y_i}^\vee - t_{y_j}^\vee\right)^2
 \end{aligned}$$

Si $t_{y_i}^\vee = \frac{t_{y_i}}{\pi_{I_i}}$ es constante para todos los i clusters, entonces $V_{p_I(s_I)}(\hat{t}_\pi) = 0$. Por lo tanto, si podemos elegir $\pi_{I_i} \propto t_{y_i}$ el muestreo por clusters será eficiente. Si lo N_i son conocidos, entonces $\pi_{I_i} \propto N_i$. Dado que $t_{y_i} = N_i \bar{y}_{U_I} = \sum_{U_I} y_k$, esta será una buena elección si hay poca variación entre los \bar{y}_{U_I} . Si todos los \bar{y}_{U_I} son iguales, entonces $V_{p_I(s_I)}(\hat{t}_\pi) = 0$.

Tomar π_{I_i} constante para todos los clusters es una elección pobre si los N_i varían mucho.