

# Nociones básicas sobre muestreo de poblaciones finitas

Daniel Czarniewicz

2017

## Muestreo sin reposición

Sea la población  $U = \{1; \dots; k; \dots; N\}$  donde  $N$  es conocido, pero los valores de  $y_k$  son desconocidos. Se busca estimar el total de la variable  $y$ ,  $t_y = \sum_U y_k$ , o de la media poblacional  $\bar{y}_U = \frac{t_y}{N} = \frac{1}{N} \sum_U y_k$ . Se selecciona una *muestra* de la población, la cual se utiliza para estimar el total y la media.

Llamamos *diseño muestral* a la función  $p(\cdot)$  tal que:

$$p(s) = \Pr(S = s) = \Pr(\text{"seleccionar la muestra } s" | \text{"estrategia de selección"})$$

$p(s)$  es entonces la función de distribución de probabilidad de una variable aleatoria  $S$  con recorrido  $\mathcal{S} = \{s_1; s_2; \dots\}$ .  $\mathcal{S}$  tiene  $2^N$  elementos, contando  $\emptyset$  y  $U$ . Dado que es una función de probabilidad,  $p(s) \geq 0 \forall s \in \mathcal{S}$ , y  $\sum_{s \in \mathcal{S}} p(s) = 1$ .

La inclusión de un elemento  $k$  en la muestra puede indicarse mediante la indicadora:

$$I_k = \begin{cases} 1 & \text{si } k \in s \\ 0 & \text{si } k \notin s \end{cases}$$

Definimos la *probabilidad de inclusión de primer orden* como:

$$\star \pi_k = \Pr(k \in s) = \Pr(I_k = 1) = \sum_{s \ni k} p(s)$$

Existen  $N$  cantidades  $\pi_1; \dots; \pi_k; \dots; \pi_N$  asociadas a un diseño  $p(\cdot)$ . Si  $\pi_k \geq 0 \forall k \in U$  decimos que el muestreo es un *muestreo probabilístico*.

Definimos la *probabilidad de inclusión de segundo orden* como:

$$\star \pi_{kl} = \Pr(k; l \in s) = \Pr(I_k I_l = 1) = \sum_{\substack{s \ni k \\ s \ni l}} p(s)$$

Existen  $\frac{N(N-1)}{2}$  cantidades  $\pi_{12}; \pi_{13}; \dots; \pi_{kl}; \dots; \pi_{N-1, N}$  asociadas a un diseño  $p(\cdot)$ , donde  $\pi_{kl} = \pi_{lk}$  y  $\pi_{kk} = \pi_k$ . Si  $\pi_{kl} \geq 0 \forall k, l \in U$  decimos que el diseño es *medible*. Solo en el caso de diseños medibles es posible obtener estimadores insesgados de la varianza.

Para todo diseño se cumple que:

- $\mathbf{E}_{p(s)}(I_k) = \pi_k \quad \forall k \in U$
- $\mathbf{Var}_{p(s)}(I_k) = \pi_k(1 - \pi_k) \quad \forall k \in U$
- $\Delta_{kl} = \mathbf{Cov}_{p(s)}(I_k; I_l) = \pi_{kl} - \pi_k \pi_l \quad \forall k \neq l \in U$

Llamamos  $n_S$  al tamaño muestral (es decir, al cardinal del conjunto  $s$ ). Para todo diseño se cumple que:

- $n_S = \sum_U I_k$
- $\mathbf{E}_{p(s)}(n_S) = \sum_U \pi_k$
- $\mathbf{Var}_{p(s)}(n_S) = \sum_U \pi_k(1 - \pi_k) + \sum_{k \neq l} \sum_U (\pi_{kl} - \pi_k \pi_l) = \sum_U \pi_k - \left( \sum_U \pi_k \right)^2 + \sum_{k \neq l} \sum_U \pi_{kl}$
- Si  $p(s)$  es de tamaño fijo  $n$ , entonces:

$$\begin{aligned}
 \star E_{p(s)}(n_S) &= \sum_U \pi_k = n & \star \sum_{k \neq l} \sum_U \pi_{kl} &= n(n-1) \\
 \star \sum_{\substack{l \in U \\ l \neq k}} \pi_{kl} &= \sum_{\substack{l \in U \\ l \neq k}} \mathbf{E}_{p(s)}(I_k I_l) = \mathbf{E}_{p(s)} \left[ I_k \left( \sum_U I_l - I_k \right) \right] = \mathbf{E}_{p(s)} \left( I_k \underbrace{\sum_U I_l}_{=n} \right) - \mathbf{E}_{p(s)}(I_k^2) = \\
 &= \mathbf{E}_{p(s)}(I_k n) - \mathbf{E}_{p(s)}(I_k) = n \mathbf{E}_{p(s)}(I_k) - \mathbf{E}_{p(s)}(I_k) = n \pi_k - \pi_k = (n-1) \pi_k
 \end{aligned}$$

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

El principio de  $\pi$ -expansión implica que el elemento muestral  $k$  representa  $1/\pi_k$  elementos en la población. Sea el siguiente estimador de  $t_y$ :

$$\begin{aligned}
 \star \hat{t}_\pi &= \sum_s y_k^\vee = \sum_s \frac{y_k}{\pi_k} \\
 \star \mathbf{E}_{p(s)}(\hat{t}_\pi) &= \mathbf{E}_{p(s)} \left( \sum_s \frac{y_k}{\pi_k} \right) = \sum_U \mathbf{E}_{p(s)}(I_k) \frac{y_k}{\pi_k} = \sum_U \pi_k \frac{y_k}{\pi_k} = \sum_U y_k = t_y \\
 \star \mathbf{Var}_{p(s)}(\hat{t}_\pi) &= \mathbf{Var}_{p(s)} \left( \sum_s y_k^\vee \right) = \sum_U \mathbf{Var}_{p(s)}(I_k) y_k^{\vee 2} + \sum_{k \neq l} \sum_U \mathbf{Cov}_{p(s)}(I_k; I_l) y_k^\vee y_l^\vee = \\
 &= \sum_U \Delta_{kk} y_k^\vee y_k^\vee + \sum_{k \neq l} \sum_U \Delta_{kl} y_k^\vee y_l^\vee = \sum \sum_U \Delta_{kl} y_k^\vee y_l^\vee \\
 \star \hat{\mathbf{Var}}_{p(s)}(\hat{t}_\pi) &= \sum \sum_s \Delta_{kl}^\vee y_k^\vee y_l^\vee \\
 \star \mathbf{E}_{p(s)}(\hat{\mathbf{Var}}_{p(s)}(\hat{t}_\pi)) &= \mathbf{E}_{p(s)} \left( \sum \sum_s \Delta_{kl}^\vee y_k^\vee y_l^\vee \right) = \sum \sum_U \mathbf{E}_{p(s)}(I_k; I_l) \Delta_{kl}^\vee y_k^\vee y_l^\vee = \\
 &= \sum \sum_U \pi_{kl} \frac{\Delta_{kl}}{\pi_{kl}} y_k^\vee y_l^\vee = \sum \sum_U \Delta_{kl} y_k^\vee y_l^\vee = \mathbf{Var}_{p(s)}(\hat{t}_\pi)
 \end{aligned}$$

Si el diseño es de tamaño fijo, son válidas las siguientes expresiones:

$$\begin{aligned}
 \star \mathbf{Var}_{p(s)}(\hat{t}_\pi) &= -\frac{1}{2} \sum \sum_U \Delta_{kl} (y_k^\vee - y_l^\vee)^2 \\
 \star \hat{\mathbf{Var}}_{p(s)}(\hat{t}_\pi) &= -\frac{1}{2} \sum \sum_s \Delta_{kl}^\vee (y_k^\vee - y_l^\vee)^2 \\
 \star \mathbf{E}_{p(s)}(\hat{\mathbf{Var}}_{p(s)}(\hat{t}_\pi)) &= \mathbf{E}_{p(s)} \left( -\frac{1}{2} \sum \sum_s \Delta_{kl}^\vee (y_k^\vee - y_l^\vee)^2 \right) = \\
 &= -\frac{1}{2} \sum \sum_U \mathbf{E}_{p(s)}(I_k; I_l) \Delta_{kl}^\vee (y_k^\vee - y_l^\vee)^2 = -\frac{1}{2} \sum \sum_U \mathbf{E}_{p(s)}(I_k; I_l) \Delta_{kl}^\vee (y_k^\vee - y_l^\vee)^2 = \\
 &= -\frac{1}{2} \sum \sum_U \Delta_{kl} (y_k^\vee - y_l^\vee)^2 = \mathbf{Var}_{p(s)}(\hat{t}_\pi)
 \end{aligned}$$

*Demostración:*

$$\begin{aligned}
\mathbf{Var}_{p(s)}(\hat{t}_\pi) &= -\frac{1}{2} \sum \sum_U \Delta_{kl} (y_k^\vee - y_l^\vee)^2 = \\
&= -\frac{1}{2} \sum \sum_U \Delta_{kl} (y_k^{\vee^2} - 2 y_k^\vee y_l^\vee + y_l^{\vee^2}) = \\
&= -\frac{1}{2} \sum \sum_U \Delta_{kl} y_k^{\vee^2} + \sum \sum_U \Delta_{kl} y_k^\vee y_l^\vee - \frac{1}{2} \sum \sum_U \Delta_{kl} y_l^{\vee^2} = \\
&= \sum \sum_U \Delta_{kl} y_k^\vee y_l^\vee - \sum \sum_U \Delta_{kl} y_k^{\vee^2} = \sum \sum_U \Delta_{kl} y_k^\vee y_l^\vee - \sum_{k \in U} y_k^{\vee^2} \sum_{l \in U} \Delta_{kl} = \\
&= \sum \sum_U \Delta_{kl} y_k^\vee y_l^\vee - \sum_{k \in U} y_k^{\vee^2} \sum_{l \in U} (\pi_{kl} - \pi_k \pi_l) = \\
&= \sum \sum_U \Delta_{kl} y_k^\vee y_l^\vee - \sum_{k \in U} y_k^{\vee^2} \left[ \sum_{l \in U} \pi_{kl} - \sum_{l \in U} \pi_k \pi_l \right] = \\
&= \sum \sum_U \Delta_{kl} y_k^\vee y_l^\vee - \sum_{k \in U} y_k^{\vee^2} \left[ n \pi_k - \pi_k \sum_{l \in U} \pi_l \right] = \\
&= \sum \sum_U \Delta_{kl} y_k^\vee y_l^\vee - \sum_{k \in U} y_k^{\vee^2} (n \pi_k - \pi_k n) = \sum \sum_U \Delta_{kl} y_k^\vee y_l^\vee
\end{aligned}$$

## Muestreo con reposición

Llamamos muestreo (no diseño) con reposición a los esquemas en los que los elementos son repuestos en la población luego de ser seleccionados. Por tanto, dos (o más) extracciones podrían producir el mismo elemento. Llamamos  $k_i$  al elemento seleccionado en la  $i$ -ésima extracción con  $i = 1; \dots; m$ . Al vector que contiene todos los elementos seleccionados lo llamamos *muestra ordenada*:  $os = \{k_1; \dots; k_m\}$ . Llamamos *multiplicidad* a la cantidad de veces que un mismo elemento fue seleccionado. Toda muestra ordenada,  $os$ , induce una muestra,  $s$ , la cual contiene a los elementos sorteados una única vez (se pierde el orden).

$$s = \{k : k = k_i \text{ para alguna extracción } i = 1; \dots; m\}$$

Sean  $p_1; \dots; p_k; \dots; p_N$  números positivos tales que  $\sum_U p_k = 1$ . Dado que los elementos se reponen en la población una vez seleccionado:

$$p_k = \Pr(\text{"seleccionar el elemento } k \text{ en la } i\text{-ésima extracción"}) \quad \forall k \in U$$

De esta forma, la probabilidad de obtener una determinada muestra ordenada será:

$$p(os) = \Pr(os = \{k_1; \dots; k_m\}) = p_{k_1} \times \dots \times p_{k_m} = \prod_{i=1}^m p_{k_i}$$

Sea la variable aleatoria  $r_k$ , la cual mide la cantidad de veces que el elemento  $k$  es extraído en las  $m$  extracciones. Por lo tanto,  $r_k \sim \text{Bin}(m; p_k)$ . Si  $N$  es lo suficientemente grande,  $r_k \stackrel{a}{\sim} \text{Poisson}(m p_k)$ .

$$\Pr(\text{"extraer } r_0 \text{ veces el elemento } k") = \Pr(r_k = r_0) = \binom{m}{r} (p_k)^r (1 - p_k)^{m-r}$$

$$\star \mathbf{E}(r_k) = m p_k \quad \star \mathbf{Var}(r_k) = m p_k (1 - p_k) \doteq m p_k$$

La probabilidad de que el elemento  $k$  nunca sea seleccionado en las  $m$  extracciones, y la probabilidad de que el elemento  $k$  sea seleccionado al menos una vez en las  $m$  extracciones son:<sup>1</sup>

$$\Pr(\text{"no seleccionar } k \text{ en ninguna de las } m \text{ extracciones"}) = \Pr(r_k = 0) = (1 - p_k)^m$$

$$\Pr(\text{"el elemento } k \text{ sea extraído"}) = \Pr(r_k \geq 1) = 1 - \Pr(r_k < 1) = 1 - \Pr(r_k = 0) = 1 - (1 - p_k)^m$$

Por lo tanto, las probabilidades de inclusión de primer y segundo orden serán:

$$\star \pi_k = \Pr(k \in S) = 1 - (1 - p_k)^m$$

$$\begin{aligned} \star \pi_{kl} &= \Pr(k, l \in S) = \Pr(k \in S) \Pr(l \in S) = [1 - (1 - p_k)^m] [1 - (1 - p_l)^m] = \\ &= 1 - (1 - p_l)^m - (1 - p_k)^m + (1 - p_k)^m (1 - p_l)^m = \\ &= 1 - (1 - p_k)^m - (1 - p_l)^m + [(1 - p_k)(1 - p_l)]^m = \\ &= 1 - [(1 - p_k)^m + (1 - p_l)^m - (1 - p_k - p_l + p_k p_l)^m] \quad \forall k \neq l \in U \end{aligned}$$

Si  $p_k = p_l \quad \forall k, l \in U$

$$\begin{aligned} \star \pi_{kl} &= 1 - [2(1 - p_k)^m - (1 - 2p_k + p_k^2)^m] = \\ &= 1 - [2(1 - p_k)^m - (1 - p_k)^{2m}] \quad \forall k \neq l \in U \end{aligned}$$

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

En muestreos con reposición  $t_y$  se estima utilizando un estimador  $p$ -expandido (en lugar del estimador  $\pi$ -expandido):

$$\star \hat{t}_{pwr} = \frac{1}{m} \sum_{i=1}^m \frac{y_k}{p_k}$$

Para hallar las propiedades estadísticas del estimador  $\hat{t}_{pwr}$  se definen las variables  $Z_i = \frac{y_{k_i}}{p_{k_i}}$ . De esta forma,

$$\hat{t}_{pwr} = \frac{1}{m} \sum_{i=1}^m Z_i = \bar{Z}. \text{ Dado que las } Z_i \text{ son iid:}$$

$$\begin{aligned} \star \Pr\left(Z_k = \frac{y_k}{p_k}\right) &= p_k \quad \forall i = 1; \dots; m \\ \star \mathbf{E}(Z_i) &= \sum_U Z_i \Pr\left(Z_k = \frac{y_k}{p_k}\right) = \sum_U \frac{y_k}{p_k} \Pr\left(Z_k = \frac{y_k}{p_k}\right) = \sum_U \frac{y_k}{p_k} p_k = \sum_U y_k = t_y \\ \star \mathbf{Var}(Z_i) &= \mathbf{E}\left[(Z_i - t_y)^2\right] = \mathbf{E}(Z_i^2) - t_y^2 = \sum_U \frac{y_k^2}{p_k} - t_y^2 = \sum_U \frac{y_k^2}{p_k} - 2t_y^2 + t_y^2 = \\ &= \sum_U \frac{y_k^2}{p_k} - 2t_y \sum_U y_k + t_y^2 \sum_U p_k = \sum_U \left[\frac{y_k^2}{p_k} - 2t_y y_k + t_y^2 p_k\right] = \\ &= \sum_U \left[\left(\frac{y_k}{p_k}\right)^2 - 2t_y \left(\frac{y_k}{p_k}\right) + t_y^2\right] p_k = \sum_U \left[\left(\frac{y_k}{p_k} - t_y\right)^2 p_k\right] = V_i \\ \star \mathbf{E}_{p(os)}(\hat{t}_{pwr}) &= \mathbf{E}(\bar{Z}) = \mathbf{E}\left(\frac{1}{m} \sum_{i=1}^m z_i\right) = \frac{1}{m} \sum_{i=1}^m \mathbf{E}(Z_i) = \frac{1}{m}(m, t_y) = t_y \end{aligned}$$

---

<sup>1</sup>Si  $m = 1$ , entonces  $\pi_k = p_k$ .

$$\begin{aligned}
\star \text{Var}_{p(os)}(\hat{t}_{pwr}) &= \text{Var}(\bar{Z}) = \frac{1}{m} \text{Var}(Z_i) = \frac{1}{m} \sum_U \left[ \left( \frac{y_k}{p_k} - t_y \right)^2 p_k \right] = \frac{V_i}{m} \\
\star \hat{V}_i &= \frac{1}{m-1} \sum_{i=1}^m \left( \frac{y_k}{p_k} - \hat{t}_{pwr} \right)^2 \\
\star \hat{\text{Var}}_{p(os)}(\hat{t}_{pwr}) &= \frac{\hat{V}_i}{m} = \frac{1}{m(m-1)} \sum_{i=1}^m \left( \frac{y_k}{p_k} - \hat{t}_{pwr} \right)^2 \\
\star \mathbf{E}_{p(os)}(\hat{\text{Var}}_{p(os)}(\hat{t}_{pwr})) &= \mathbf{E}_{p(os)}(\text{Var}(\bar{Z})) = \frac{1}{m} \mathbf{E}_{p(os)}(\hat{\text{Var}}(Z_i)) = \frac{1}{m} \text{Var}(Z_i) = \frac{V_i}{m} \Leftrightarrow \\
\Leftrightarrow \hat{V}_i &= \frac{1}{m-1} \sum_{i=1}^m (Z_i - \bar{Z})^2 = \frac{1}{m-1} \sum_{i=1}^m \left( \frac{y_k}{p_k} - \hat{t}_{pwr} \right)^2 \\
\therefore \hat{\text{Var}}_{p(os)}(\hat{t}_{pwr}) &\text{ es insesgada para } \text{Var}_{p(os)}(\hat{t}_{pwr})
\end{aligned}$$

$$\begin{aligned}
\text{Si } y_k = c p_k \quad \forall k \in U &\Rightarrow t_y = \sum_U y_k = \sum_U c p_k = c \underbrace{\sum_U p_k}_{=1} = c \Rightarrow \\
\Rightarrow \text{Var}_{p(os)}(\hat{t}_{pwr}) &= \frac{1}{m} \sum_U \left( \frac{y_k}{p_k} - t_y \right)^2 p_k = \frac{1}{m} \sum_U \left( \frac{c p_k}{p_k} - c \right)^2 p_k = \frac{1}{m} \sum_U (c - c)^2 p_k = 0
\end{aligned}$$

En la práctica no es posible establecer  $y_k = c p_k$ . Pero sí es posible establecer  $p_k = \frac{x_k}{\sum_U x_k} \quad \forall k \in U$ , donde  $x_k$  es una variable auxiliar (es decir,  $x_k$  y  $y_k$  están altamente correlacionadas).

## El estimador $\hat{t}_\pi$ en el diseño ordenado

$$\star \hat{t}_\pi = \sum_s y_k^\checkmark = \sum_s \frac{y_k}{\pi_k} = \sum_s \frac{y_k}{1 - (1 - p_k)^m}$$

El cual es insesgado para  $t_y$ , y la expresión para su varianza y el estimador de su varianza son:

$$\begin{aligned}
\star \text{Var}_{p(s)}(\hat{t}_\pi) &= \sum \sum_U \Delta_{kl} y_k^\checkmark y_l^\checkmark \\
\star \hat{\text{Var}}_{p(s)}(\hat{t}_\pi) &= \sum \sum_s \Delta_{kl}^\checkmark y_k^\checkmark y_l^\checkmark
\end{aligned}$$

No se pueden comparar los estimadores  $\hat{t}_{pwr}$  y  $\hat{t}_\pi$ . Ambos son insesgados, pero no se puede concluir sobre sus varianzas. La varianza del  $\hat{t}_{pwr}$  depende de los valores de  $y_k$ , lo cuales son desconocidos.

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

$$\star \hat{t}_{alt} = N \bar{y}_S = \frac{N}{n_S} \sum_s y_k$$

Este estimador tiene como inconveniente que el tamaño muestral  $n_S$  es aleatorio. En general,  $\hat{t}_{alt}$  tiene menor varianza que  $\hat{t}_{pwr}$  o  $\hat{t}_\pi$

## Descomposición de la varianza

Supongamos que la población \$ U = \{ 1; \dots; k; \dots; N \} \$ se encuentra particionada en \$ a \$ grupos de tamaño \$ n \$, \$ \{ S\_1; \dots; S\_r, \dots; S\_a \} \$. Luego entonces:

$$\begin{aligned}
 \sum_U (y_k - \bar{y}_U)^2 &= \sum_{r=1}^a \sum_{S_r} (y_k - \bar{y}_U)^2 = \sum_{r=1}^a \left[ \sum_{S_r} \left( (y_k - \bar{y}_{S_r}) + (\bar{y}_{S_r} - \bar{y}_U) \right)^2 \right] = \\
 &= \sum_{r=1}^a \left[ \sum_{S_r} (y_k - \bar{y}_{S_r})^2 + 2 \sum_{S_r} (y_k - \bar{y}_{S_r})(\bar{y}_{S_r} - \bar{y}_U) + \sum_{S_r} (\bar{y}_{S_r} - \bar{y}_U)^2 \right] = \\
 &= \sum_{r=1}^a \left[ \sum_{S_r} (y_k - \bar{y}_{S_r})^2 + 2 (\bar{y}_{S_r} - \bar{y}_U) \underbrace{\sum_{S_r} (y_k - \bar{y}_{S_r})}_{=0} + \sum_{S_r} (\bar{y}_{S_r} - \bar{y}_U)^2 \right] = \\
 &= \sum_{r=1}^a \left[ \sum_{S_r} (y_k - \bar{y}_{S_r})^2 + \sum_{S_r} (\bar{y}_{S_r} - \bar{y}_U)^2 \right] = \\
 &= \sum_{r=1}^a \sum_{S_r} (y_k - \bar{y}_{S_r})^2 + \sum_{r=1}^a \sum_{S_r} (\bar{y}_{S_r} - \bar{y}_U)^2 = \\
 &= \sum_{r=1}^a \sum_{S_r} (y_k - \bar{y}_{S_r})^2 + \sum_{r=1}^a n (\bar{y}_{S_r} - \bar{y}_U)^2
 \end{aligned}$$

Por lo tanto, tenemos que:

$$\underbrace{\sum_U (y_k - \bar{y}_U)^2}_{SST} = \underbrace{\sum_{r=1}^a \sum_{S_r} (y_k - \bar{y}_{S_r})^2}_{SSW} + \underbrace{\sum_{r=1}^a n (\bar{y}_{S_r} - \bar{y}_U)^2}_{SSB} \Rightarrow \boxed{SST = SSW + SSB}$$

## Tamaño muestral

Supongamos que se quiere estimar \$ t\_y \$ usando el estimador \$ \hat{t}\_y \$ y el estimador de su varianza \$ \hat{V}\_{p(s)}(\hat{t}\_y) \$. Supongamos que ambos estimadores son (aprox) insesgados y que es razonable suponer que:

$$\frac{\hat{t}_y - t_y}{\sqrt{\hat{V}_{p(s)}(\hat{t}_y)}} \stackrel{a}{\sim} N(0; 1)$$

Buscamos un \$ \mathbf{E}(n\_S) \$ tal que para una precisión dada, \$ \varepsilon > 0 \$, y un nivel de confianza dado, \$ 0 < \alpha < 1 \$, nos permita plantear:

$$\Pr(|\hat{t}_y - t_y| < \varepsilon) \doteq 1 - \alpha \Rightarrow \Pr\left(\hat{t}_y - z_{1-\alpha/2} \sqrt{\hat{V}_{p(s)}(\hat{t}_y)} < t_y < \hat{t}_y + z_{1-\alpha/2} \sqrt{\hat{V}_{p(s)}(\hat{t}_y)}\right) \doteq 1 - \alpha$$

$$\Pr(|\hat{t}_y - t_y| < \varepsilon) \doteq 1 - \alpha \Rightarrow \Pr\left(\frac{|\hat{t}_y - t_y|}{\sqrt{\hat{V}_{p(s)}(\hat{t}_y)}} < \frac{\varepsilon}{\sqrt{\hat{V}_{p(s)}(\hat{t}_y)}}\right) \doteq 1 - \alpha$$

Si \$ \frac{\varepsilon}{\sqrt{\mathbf{Var}\_{p(s)}(\hat{t}\_y)}} \doteq z\_{1-\alpha/2} \Rightarrow \varepsilon^2 \doteq z\_{1-\alpha/2}^2 \mathbf{Var}\_{p(s)}(\hat{t}\_y) \$, donde \$ \varepsilon \$ y \$ \alpha \$ están fijos y, en general, \$ \mathbf{Var}\_{p(s)}(\hat{t}\_y) \$ depende de \$ n \$ y de \$ S\_{y\_U}^2 \$. Luego si se cuenta con una buena estimación de \$ S\_{y\_U}^2 \$, se puede despejar \$ n \$.

Para un tamaño de muestra fijo, si disminuimos  $\varepsilon$ , reducimos la amplitud del intervalo, con lo que reducimos la confianza, o sea, aumentamos  $\alpha$ .

En la práctica  $S_{y_U}^2$  es desconocida, pero se pueden ensayar alguna de las siguientes estrategias para obtener una aproximación:

1. Se presume algún tipo de distribución para los valores de  $y$  en la población. Luego se busca una cota para  $S_{y_U}^2$ .
  - Si  $y \sim \text{Ber} \Rightarrow 0 \leq S_{y_U}^2 \leq 1/4$
  - Si  $y \stackrel{a}{\sim} N \Rightarrow S_{y_U}^2 \doteq \frac{y_{k(n)} - y_{k(1)}}{6}$  para un  $\alpha \doteq 0,01$
2. Utilizar datos de algún relevamiento reciente o de alguna variable auxiliar para la que se pueda asumir una variabilidad similar. En estos casos se suele utilizar el  $CV_{y_u}$  dado que es más estable que la varianza. Luego se fija una precisión relativa y se determina el valor de  $n$

$$\Pr(|\hat{t}_y - t_y| < t_y \varepsilon) \doteq 1 - \alpha \Rightarrow \frac{t_y \varepsilon}{\sqrt{\text{Var}_{p(s)}(\hat{t}_y)}} = z_{1-\alpha/2} \Rightarrow \varepsilon = z_{1-\alpha/2} \left[ \frac{\text{Var}_{p(s)}(\hat{t}_y)}{t_y} \right]$$

3. Tomar una pequeña “muestra de iluminación” y calcular  $S_{y_s}^2$  en dicha muestra. Luego se utilizar esta estimación como estimador de  $S_{y_U}^2$  para calcular  $n$ .

## Desarrollos de Taylor para aproximaciones en muestreo de poblaciones finitas

Supongamos que queremos estimar:  $\theta = f(t_1; \dots; t_q) = f(\mathbf{t})$  donde  $t_j = \sum_U y_{jk}$   $j = 1; \dots; q$ , son los totales de las  $q$  variables poblacionales relevadas. Un estimador podría ser  $\hat{\theta} = f(\hat{t}_{1\pi}; \dots; \hat{t}_{q\pi}) = f(\hat{\mathbf{t}}_\pi)$  donde  $\hat{t}_j = \sum_s y_{jk}^\vee$   $j = 1; \dots; q$ , son los totales de las  $q$  variables poblacionales relevadas, estimados en la muestra  $s$ .

Si  $f(\hat{\mathbf{t}}_\pi)$  es lineal tenemos que  $\theta = a_0 + \sum_{j=1}^q a_j t_j = a_0 + \mathbf{a}'\mathbf{t}$ . Luego entonces  $\hat{\theta} = a_0 + \sum_{j=1}^q a_j \hat{t}_{j\pi} = a_0 + \mathbf{a}'\hat{\mathbf{t}}_\pi$  estima  $\theta$  de forma tal que:

- $\mathbf{E}(\hat{\theta}) = \mathbf{E}(a_0 + \mathbf{a}'\hat{\mathbf{t}}_\pi) = \mathbf{E}(a_0) + \mathbf{E}(\mathbf{a}'\hat{\mathbf{t}}_\pi) = a_0 + \mathbf{a}'\mathbf{t} \Rightarrow \hat{\theta}$  es insesgado para  $\theta$ .
- $\text{Var}(\hat{\theta}) = \text{Var}(a_0 + \mathbf{a}'\hat{\mathbf{t}}_\pi) = \sum_{j=1}^q \sum_{j'=1}^q a_{jj} \text{Cov}(\hat{t}_{j\pi}; \hat{t}_{j'\pi}) = \mathbf{a}' V(\hat{\mathbf{t}}_\pi) \mathbf{a}$ , donde

$$\text{Cov}(\hat{t}_{j\pi}; \hat{t}_{j'\pi}) = \sum \sum_U \Delta_{kl} y_{jk}^\vee y_{j'k}^\vee$$

Podemos reescribir  $\hat{\theta}$  de la siguiente forma:

$$\star \hat{\theta} = a_0 + \sum_{j=1}^q a_j \hat{t}_{j\pi} = a_0 + \mathbf{a}'\hat{\mathbf{t}}_\pi = a_0 + \sum_{j=1}^q a_j \sum_s y_{jk}^\vee = a_0 + \sum_{j=1}^q \sum_s a_j y_{jk}^\vee = a_0 + \sum_s u_k^\vee$$

$$\text{con } u_k = \sum_{j=1}^q a_j y_{jk} \text{ y } u_k^\vee = \frac{u_k}{\pi_k}$$

$$\star \text{Var}(\hat{\theta}) = \sum \sum_U \Delta_{kl} u_k^\vee u_l^\vee$$

$$\star \mathbf{\hat{Var}}(\hat{\theta}) = \sum \sum_s \Delta_{kl}^{\check{}} u_k^{\check{}} u_l^{\check{}}$$

Si  $f(\hat{\mathbf{t}}_\pi)$  no es lineal,  $\hat{\theta}$  debe aproximarse linealmente, y luego se podrán calcular  $\mathbf{Var}(\hat{\theta})$  y  $\mathbf{\hat{Var}}(\hat{\theta})$ . La técnica aproxima  $\hat{\theta}$  por un pseudo-estimador,  $\hat{\theta}_0$ , que es lineal en  $\hat{\mathbf{t}}_\pi$ . En general  $\hat{\theta}_0$  dependerá de cantidades desconocidas (de ahí que se le llama pseudo-estimador). La técnica para hallar  $\hat{\theta}_0$  consiste en la aproximación de Taylor de primer orden de la función  $f$ , en el entorno de un punto  $\mathbf{t}$ , y despreciar el término de error.

$$\hat{\theta} \doteq \hat{\theta}_0 = \theta + \sum_{j=1}^q a_j (\hat{t}_{j\pi} - t_j) \quad \text{donde} \quad a_j = \left. \frac{\partial f}{\partial t_{j\pi}} \right|_{\hat{\mathbf{t}}_\pi = \mathbf{t}}$$

En muestras grandes  $\hat{\mathbf{t}}_\pi \approx \mathbf{t} \Rightarrow \hat{\theta}_0 = \hat{\theta}$  y  $\mathbf{AVar}(\hat{\theta}) = \mathbf{Var}(\hat{\theta}_0)$ .

$$\begin{aligned} \star \mathbf{AVar}(\hat{\theta}) &\doteq V(\hat{\theta}_0) = V\left(\sum_{j=1}^q a_j \hat{t}_{j\pi}\right) = V\left(\sum_{j=1}^q a_j \sum_s \frac{y_{jk}}{\pi_k}\right) = \\ &= \mathbf{Var}\left(\sum_s u_k^{\check{}}\right) = \sum \sum_U \Delta_{kl} u_k^{\check{}} u_l^{\check{}} \end{aligned}$$

Como  $\mathbf{E}(\hat{\theta}_0) = \theta \Rightarrow MSE(\hat{\theta}) \doteq MSE(\hat{\theta}_0) = \mathbf{Var}(\hat{\theta}_0) = \mathbf{AVar}(\hat{\theta})$

Las cantidades  $u_k$  dependen de  $a_j = \left. \frac{\partial f}{\partial t_{j\pi}} \right|_{\hat{\mathbf{t}}_\pi = \mathbf{t}}$  que es desconocida ya que  $t_j$  es desconocido  $\forall j$ . De todas

formas, la estimación puntual será  $\hat{\theta}_0 = f(\hat{\mathbf{t}}_\pi)$ . Para estimar la varianza se reemplaza  $a_j$  por  $\hat{a}_j = \left. \frac{\partial f}{\partial t_{j\pi}} \right|_{\hat{\mathbf{t}}_\pi = \hat{\mathbf{t}}_0}$ ,

siendo  $\hat{\mathbf{t}}_0$  el total observado en la muestra. Luego entonces  $\hat{u}_k = \sum_{j=1}^q \hat{a}_j y_{jk}$ , con lo que puede calcularse

$\mathbf{\hat{Var}}(\hat{\theta}) = \sum \sum_s \Delta_{kl}^{\check{}} \hat{u}_k^{\check{}} \hat{u}_l^{\check{}}$ . Esto es válido ya que  $\hat{u}_k$  es consistente para estimar  $u_k$ .



## Estimador de una razón

El problema consiste en estimar un cociente entre totales poblacionales:

$$\star R = \frac{t_y}{t_z} = \frac{\bar{y}_U}{\bar{z}_U}$$

Sea el estimador:

$$\star \hat{R} = f(\hat{t}_{y\pi}; \hat{t}_{z\pi}) = \frac{\hat{t}_{y\pi}}{\hat{t}_{z\pi}} = \frac{\bar{y}_s}{\bar{z}_s}$$

Utilizando la linealización de Taylor:

$$\hat{R} = \hat{R}_0 = R + a_1(\hat{t}_{y\pi} - t_y) + a_2(\hat{t}_{z\pi} - t_z)$$

donde  $a_1 = \left. \frac{\partial f}{\partial \hat{t}_{y\pi}} \right|_{\substack{\hat{t}_{y\pi}=t_y \\ \hat{t}_{z\pi}=t_z}} = \frac{1}{t_z}$  y  $a_2 = \left. \frac{\partial f}{\partial \hat{t}_{z\pi}} \right|_{\substack{\hat{t}_{y\pi}=t_y \\ \hat{t}_{z\pi}=t_z}} = -\frac{t_y}{t_z^2} = -\frac{R}{t_z}$

Luego entonces:

$$\begin{aligned} \hat{R} \doteq \hat{R}_0 &= R + \frac{1}{t_z}(\hat{t}_{y\pi} - t_y) - \frac{R}{t_z}(\hat{t}_{z\pi} - t_z) = R + \frac{\hat{t}_{y\pi}}{t_z} - \frac{t_y}{t_z} - \frac{R\hat{t}_{z\pi}}{t_z} + \frac{Rt_z}{t_z} = \\ &= R + \frac{\hat{t}_{y\pi}}{t_z} - R - \frac{R\hat{t}_{z\pi}}{t_z} + R = R + \frac{\hat{t}_{y\pi}}{t_z} - \frac{R\hat{t}_{z\pi}}{t_z} = R + \frac{1}{t_z}(\hat{t}_{y\pi} - R\hat{t}_{z\pi}) = \\ &= R + \frac{1}{t_z} \sum_s \frac{y_k - Rz_k}{\pi_k} = R + \sum_s \frac{u_k}{\pi_k} \quad \text{donde } u_k = \frac{1}{t_z}(y_k - Rz_k) \end{aligned}$$

En conclusión:

$$\boxed{\hat{R} \doteq \hat{R}_0 = R + \sum_s \frac{u_k}{\pi_k}}$$

$\hat{R} \doteq \hat{R}_0$  es aproximadamente insesgado para  $R$ .

$$\begin{aligned} \star \mathbf{E}(\hat{R}) &\doteq \mathbf{E}(\hat{R}_0) = E\left(R + \sum_s \frac{u_k}{\pi_k}\right) = R + \sum_U u_k = R + \sum_U \frac{y_k}{t_z} - R \sum_U \frac{z_k}{t_z} = \\ &= R + \frac{1}{t_z} \sum_U y_k - \frac{R}{t_z} \sum_U z_k = R + \frac{t_y}{t_z} - \frac{R}{t_z} t_z = \frac{t_y}{t_z} = R \\ \star \mathbf{AVar}(\hat{R}) &= \mathbf{Var}(\hat{R}_0) = \sum \sum_U \Delta_{kl} u_k^\vee u_l^\vee = \frac{1}{t_z^2} \sum \sum_U \Delta_{kl} E_k^\vee E_l^\vee \\ &\quad \text{donde } E_k = y_k - Rz_k \\ \star \hat{\mathbf{Var}}(\hat{R}_0) &= \sum \sum_s \Delta_{kl}^\vee \hat{u}_k^\vee \hat{u}_l^\vee = \frac{1}{\hat{t}_z^2} \sum \sum_s \Delta_{kl}^\vee e_k^\vee e_l^\vee \\ &\quad \text{donde } e_k = y_k - \hat{R} z_k \\ \star \mathbf{E}(\hat{\mathbf{Var}}(\hat{R}_0)) &= \mathbf{E}\left(\sum \sum_s \Delta_{kl}^\vee \hat{u}_k^\vee \hat{u}_l^\vee\right) = \sum \sum_U \Delta_{kl} \hat{u}_k^\vee \hat{u}_l^\vee \end{aligned}$$

Por lo tanto,  $\hat{\mathbf{Var}}(\hat{R}_0)$  es aproximadamente insesgado para estimar  $\mathbf{Var}(\hat{R})$ .

Las expresiones anteriores para  $\mathbf{AVar}(\hat{R})$  y  $\hat{\mathbf{Var}}(\hat{R})$  son equivalentes a:

$$\begin{aligned} \star \mathbf{AVar}(\hat{R}) &= \frac{1}{t_z^2} \left[ \mathbf{Var}(\hat{t}_{y\pi}) + R^2 V(\hat{t}_{z\pi}) - 2R \mathbf{Cov}(\hat{t}_{y\pi}; \hat{t}_{z\pi}) \right] \\ \star \hat{\mathbf{Var}}(\hat{R}) &= \frac{1}{\hat{t}_z^2} \left[ \hat{\mathbf{Var}}(\hat{t}_{y\pi}) + \hat{R}^2 \hat{\mathbf{Var}}(\hat{t}_{z\pi}) - 2\hat{R} \hat{\mathbf{Cov}}(\hat{t}_{y\pi}; \hat{t}_{z\pi}) \right] \end{aligned}$$

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

El objetivo es estimar  $t_y$ , y se cuenta con una variable auxiliar  $z$  conocida  $\forall k \in U$ . Sea el “estimador de razón”:

$$\star \hat{t}_{yra} = \frac{\hat{t}_y \pi}{\hat{t}_z \pi} t_z = \hat{R} t_z = \frac{t_z}{\hat{t}_z \pi} \hat{t}_y \pi$$

$$\star \mathbf{AVar}(\hat{t}_{yra}) = \mathbf{AVar}(\hat{R} t_z) = t_z^2 \mathbf{AVar}(\hat{R}) = \sum \sum_U \Delta_{kl} \left( \frac{y_k - R z_k}{\pi_k} \right) \left( \frac{y_l - R z_l}{\pi_l} \right)$$

$$\star \hat{\mathbf{Var}}(\hat{t}_{yra}) = \sum \sum_s \Delta_{kl}^{\checkmark} \left( \frac{y_k - \hat{R} z_k}{\pi_k} \right) \left( \frac{y_l - \hat{R} z_l}{\pi_l} \right)$$

La lógica detrás de este estimador es la misma que en el  $\hat{t}_{alt} = \frac{N}{\hat{N}} \hat{t}_y \pi$  donde  $z_k = 1 \ \forall k \in U$ .