

MUEST- T9. 1. MUESTREO de CONGLOMERADOS SIN SUBMUESTREO.

3. COEFICIENTE de CORRELACIÓN INTRACONGLOMERADO y su INTERPRETACIÓN.
 2. ESTIMADORES, VARIANZAS y sus ^{ESTIMACIONES} ~~INTERPRETACIONES~~.
 - EFFECTO de DISEÑO.
 4. UTILIZACIÓN de ESTIMADORES de RAZÓN
-

1 - MUESTREO de CONGLOMERADOS

El muestreo de conglomerados considera una población finita con un total de M unidades elementales o átomos, agrupadas en N conglomerados ~~o~~ o unidades primarias, de forma que constituyen una partición de la población. La unidad de muestreo es el conglomerado y de la población se extrae una muestra aleatoria de n conglomerados, a partir de la cual se estiman los parámetros poblacionales.

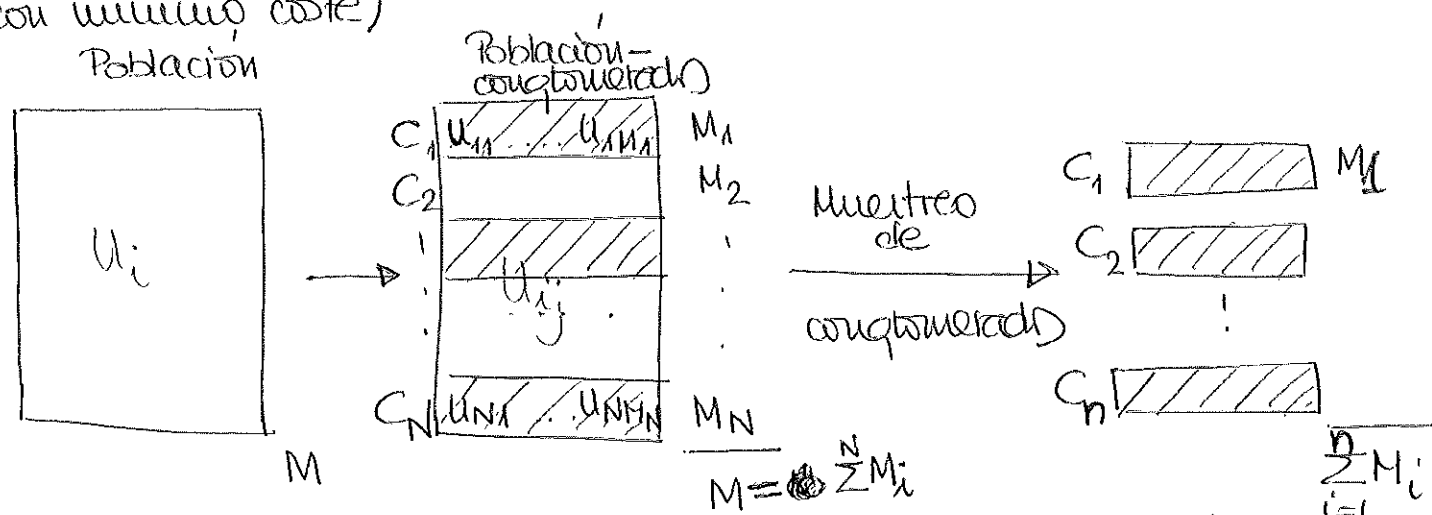
- Si de cada conglomerado muestral se observan todas las unidades elementales \rightarrow muestreo de conglomerados monetápico o sin submuestreo.
- Si de cada conglomerado se selecciona una muestra \rightarrow muestreo polietápico.

El n.º de unidades elementales de cada conglomerado se denomina tamaño del conglomerado $\left\{ \begin{array}{l} \text{tamaño igual} \\ \text{tamaño } \neq \end{array} \right.$

~~Los conglomerados deben ser~~

El muestreo por conglomerados debería ser de modo que los conglomerados fuesen lo más homogéneos entre sí (a la hora de constituir la muestra no se pierde información) y dentro de cada conglomerado se recoja la mayor heterogeneidad posible (similar a la de la población, población a escala).

La situación ideal es que un único conglomerado pudiera representar fielmente a la población (muestra de tamaño 1 con mínimo coste)



Unidad elemental, u_{ij} $\left\{ \begin{array}{l} i=1 \dots N \rightarrow n^{\circ} \text{ conglomerados} \\ j=1 \dots M_i \rightarrow \text{tamaño muestra conglomerado} \end{array} \right.$

Unidad elemental, u_{ij} $\left\{ \begin{array}{l} i=1 \dots n \rightarrow n^{\circ} \text{ conglomer. muestra} \\ j=1 \dots M_i \rightarrow \text{en muestreo monoetápico se observan todas las unid. del conglomerado} \end{array} \right.$

El muestreo de conglomerados sin submuestreo consiste en seleccionar aleatoriamente n conglomerados (se pueden utilizar \neq tipos de muestreo), y observar dentro de cada conglomerado todas sus unidades elementales, es decir $\sum_{i=1}^n M_i$ unid. elementales.

Como cada conglomerado se observa de manera exhaustiva, sin efectuar submuestreo dentro de él, también se denomina muestreo de conglomerados monoetápico, MCM(n).

Antes de continuar, quizá sería conveniente explicar detalladamente las diferencias entre muestreo estratificado y muestreo por conglomerados.

~~El~~ El muestreo estratificado ~~se~~ trata de obtener una partición de la población en estratos heterógenos entre sí y homogéneos dentro. El muestreo por conglomerado divide a la población en conglomerados homogéneos entre sí y heterógenos dentro.

La situación ideal del muestreo estratificado es obtener una muestra de tamaño 1 de cada estrato \rightarrow 1 unid. elem.
La situación ideal del muestreo por conglomerado es ~~haber~~ tener información de un único conglomerado ~~y observar~~ y en el caso monoetápico, observar de él todas sus unidades elem.
 \rightarrow Mⁱ unid. elementales.

Es muy frecuente que los conglomerados estén definidos como áreas geográficas (división territorial de la población), por lo que al muestreo por conglomerados también se le conoce como muestreo por áreas. Se utiliza por razones de economía en coste, en tiempo, en recursos, etc., y algunas veces porque se disminuye el sesgo al ser más fácil la supervisión.

Nótese que para efectuar m.a.s. hace falta una lista detallada con todos los ~~elementos~~ de la población (marco), y en muestreo estratificado una lista detallada por estratos y code uno de los estratos. (mucho coste y mucho tiempo)

En la práctica, no se dispone de tales listas. Es preferible dividir la población en áreas y elaborar las listas solamente de las áreas seleccionadas.

Ventajas del muestreo por conglomerados:

- + Marco + fácil: No se necesita un marco muy específico, como en el caso de mas j de me.
- + Marco + barato: Al seleccionar previamente las áreas sobre las que elaborar el listado de unidad, elem., el marco de conglomerados es mas fácil de conseguir, en términos de tiempo, dinero y de efectivos.
- + Marco ya existe: Se pueden utilizar archivos ya existentes por necesidades administrativas como censos.
- + Proceso mas rápido y barato: la concentración de unidades disminuye el tiempo de ~~desp~~ recogida de datos, reduce el coste por desplazamiento y permite utilizar recursos disponibles por zona (personal y materiales).

Inconvenientes:

- Menor precisión en las estimaciones: Aunque lo ideal es que la heterogeneidad dentro de cada conglomerado sea máxima, siempre va a existir un grado de homogeneidad inevitable dentro de los conglomerados.
- La eficiencia disminuye al aumentar el tamaño de los conglomerados: cuando en la práctica este tipo de muestreo es muy útil en poblaciones muy numerosas en las que se puedan construir conglomerados grandes.

2 - ESTIMADORES, VARIANZAS Y ESTIMACIONES

El muestreo por conglomerados distingue las siguientes situaciones:

- Conglomerados de igual tamaño (AQUÍ)
- Conglomerados de \neq tamaño $\left\{ \begin{array}{l} \text{parecido} \\ \text{muy } \neq \end{array} \right\}$ En 4. RAZÓN.

Además, el método de muestreo elegido para seleccionar los conglomerados muestrales lleva a estimadores \neq s, con distintas varianzas. Distinguimos:

- Muestreo sin reposición $\left\{ \begin{array}{l} \text{probab.} = (*) \text{ + utilizado} \\ \text{probab.} \neq \end{array} \right.$
- Muestreo con reposición $\left\{ \begin{array}{l} \text{probab.} = \\ \text{probab.} \neq \end{array} \right.$

En este epígrafe, consideramos el caso más sencillo: todos los conglomerados tienen el mismo tamaño $\overline{M} = \frac{M}{N}$ y las probab. son iguales

Sea:

$N = n^{\circ}$ conglomerados de la población

$n = n^{\circ}$ conglomerados en la muestra

$\overline{M} =$ tamaño del conglomerado ($M_1 = \dots = M_N = \overline{M}$)

$N\overline{M} = u^{\circ}$ total de unid. elementales en la poblac. ($N\overline{M} = M$)

$n\overline{M} = u^{\circ}$ total de unid. elementales en la muestra.

Sea $\Theta = \left[\sum_{i=1}^N Y_i \right] = \sum_{i=1}^N \sum_{j=1}^{\overline{M}} Y_{ij}$, parámetro poblacional a estimar

↓
notación anterior

$N = u^{\circ}$ unid. elementales de la poblac.

$N\overline{M}$ tamaño poblacional.

a) SIN reposición, probab IGUALES : (probab $\neq \Rightarrow \pi_i$) .

El estimador lineal insesgado de θ es el estimador

de Horvitz y Thompson ^{probab. iguales $\pi_i = \frac{n}{N}$}

$$\hat{\theta}_{HT} = \sum_{i=1}^n \sum_{j=1}^{\bar{M}} \frac{y_{ij}}{\pi_i} = \sum_{i=1}^n \sum_{j=1}^{\bar{M}} \frac{y_{ij}}{n/N} = \frac{N}{n} \sum_{i=1}^n \sum_{j=1}^{\bar{M}} y_{ij}.$$

Es insesgado de θ :

$$e_i = \begin{cases} 1 & \text{si cong. } i \text{ e muestra con probab. } \pi_i \\ 0 & \text{no} \end{cases} \quad \left| \quad e_i \rightarrow \text{Bern}(\pi_i) \right.$$

$$\begin{aligned} E[\hat{\theta}_{HT}] &= E\left[\sum_{i=1}^n \sum_{j=1}^{\bar{M}} \frac{y_{ij}}{\pi_i}\right] = E\left[\sum_{i=1}^n \sum_{j=1}^{\bar{M}} \frac{y_{ij}}{\pi_i} e_i\right] = \sum_{i=1}^n \sum_{j=1}^{\bar{M}} \frac{y_{ij}}{\pi_i} E[e_i] = \\ &= \sum_{i=1}^n \sum_{j=1}^{\bar{M}} \frac{y_{ij}}{\pi_i} \pi_i = \sum_{i=1}^n \sum_{j=1}^{\bar{M}} y_{ij} = \theta. \end{aligned}$$

La expresión del estimador $\hat{\theta}_{HT}$ en los parámetros más utilizados es:

$$\begin{aligned} \theta = X = \sum_{i=1}^n \sum_{j=1}^{\bar{M}} X_{ij} &\Rightarrow \hat{X} = \frac{N}{n} \sum_{i=1}^n \sum_{j=1}^{\bar{M}} X_{ij} = \frac{N\bar{M}}{n} \sum_{i=1}^n \frac{1}{\bar{M}} \sum_{j=1}^{\bar{M}} X_{ij} = \\ &= NM \cdot \frac{1}{n} \sum_{i=1}^n \bar{X}_i = NM \cdot \bar{\bar{X}} \quad \bar{X}_i \end{aligned}$$

\hookrightarrow el estimador del total es el u° total de unidades de la población multiplicado por el estimador de la media (estimador de expansión)

$$\theta = \bar{X} = \sum_{i=1}^n \sum_{j=1}^{\bar{M}} \frac{X_{ij}}{NM} \Rightarrow \hat{\bar{X}} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{\bar{M}} X_{ij} = \frac{1}{n} \sum_{i=1}^n \frac{1}{\bar{M}} \sum_{j=1}^{\bar{M}} X_{ij} = \bar{\bar{X}}$$

\hookrightarrow el estimador insesgado de la media poblacional es la media muestral de la media ^{pobla} de los conglomerados ^{pobl.}

$$\theta = P = \sum_{i=1}^n \sum_{j=1}^{\bar{M}} \frac{A_{ij}}{NM} \Rightarrow \hat{P} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^{\bar{M}} A_{ij} = \frac{1}{n} \sum_{i=1}^n \frac{1}{\bar{M}} \sum_{j=1}^{\bar{M}} A_{ij} = \frac{1}{n} \sum_{i=1}^n \bar{P}_i$$

$$\theta = A = \sum_{i=1}^n \sum_{j=1}^{\bar{M}} A_{ij} \Rightarrow \hat{A} = NM \cdot \hat{P} = NM \cdot \frac{1}{n} \sum_{i=1}^n \bar{P}_i$$

VARIANZAS :

Para la media:

$$V(\bar{x}) = V\left(\frac{1}{n} \sum_{i=1}^n \bar{x}_i\right) \stackrel{\text{SR}}{=} (1-f) \cdot \frac{\sum_{i=1}^N (\bar{x}_i - \bar{x})^2}{N-1} \stackrel{\text{MS}_x^2 = S_b^2}{=} \frac{(1-f)}{n \bar{M}} \cdot \frac{\sum_{i=1}^N \bar{M} (\bar{x}_i - \bar{x})^2}{N-1} =$$

$$\stackrel{\text{SR}}{=} \frac{(1-f)}{n \bar{M}} \cdot S_b^2 = \frac{(1-f)}{n} \cdot S_b^2 \quad \text{doble } S_x^2 \quad \text{poblc} \quad \text{entre conglomerados}$$

$$\text{doble } S_b^2 = \frac{\sum_{i=1}^N \bar{M} (\bar{x}_i - \bar{x})^2}{N-1} \quad \text{between}$$

de la cual se pueden deducir las expresiones de la variancia de los otros estimadores:

Para el total poblacional:

$$\hat{X} = N \bar{M} \bar{x} \Rightarrow V(\hat{X}) = N^2 \bar{M}^2 V(\bar{x}) = N^2 \bar{M} \cdot \frac{S_b^2}{n} (1-f)$$

$$\hat{P} = \frac{1}{n} \sum_{i=1}^n P_i \Rightarrow V(\hat{P}) = \frac{(1-f)}{n \bar{M}} \cdot S_b^2 = \frac{(1-f)}{n} \cdot \frac{\sum_{i=1}^N (P_i - \bar{P})^2}{N-1}$$

$$\text{doble } S_b^2 = \frac{\sum_{i=1}^N (P_i - \bar{P})^2}{N-1}$$

$$\hat{A} = N \bar{M} \hat{P} \Rightarrow V(\hat{A}) = N^2 \bar{M}^2 V(\hat{P}) = N^2 \bar{M} \cdot \frac{(1-f)}{n} S_b^2$$

ESTIMACIONES de las VARIANZAS

Como la cuasivariancia muestral entre conglomerados \hat{S}_b^2 es un estimador insesgado de la variancia poblacional S_b^2 :

$$E[\hat{S}_b^2] = E\left[\frac{\bar{M}}{n-1} \sum_{i=1}^n (\bar{x}_i - \bar{x})^2\right] = \bar{M} E\left[\frac{\sum_{i=1}^n (\bar{x}_i - \bar{x})^2}{n-1}\right] = \bar{M} E[\hat{S}_x^2] =$$

$$= \bar{M} \cdot S_x^2 = S_b^2$$

$$\hat{V}(\bar{x}) = \frac{(1-f)}{n \bar{M}} \cdot S_b^2 \quad \text{con } \hat{S}_b^2 = \frac{\bar{M}}{n-1} \sum_{i=1}^n (\bar{x}_i - \bar{x})^2$$

$$\hat{V}(\hat{X}) = N^2 \bar{M} \frac{(1-f)}{n} \hat{S}_b^2$$

$$\hat{V}(\hat{P}) = \frac{(1-f)}{n \bar{M}} \hat{S}_b^2 \quad \text{con } \hat{S}_b^2 = \frac{\bar{M}}{n-1} \sum_{i=1}^n (P_i - \bar{P})^2$$

$$\hat{V}(\hat{A}) = N^2 \bar{M} \frac{(1-f)}{n} \hat{S}_b^2$$

MIRAR

⊗ DESCOMPOSICIÓN de la VARIANZA

$$SCT = SCD + SCE$$

$$\sum_{i=1}^N \sum_{j=1}^{\bar{M}} (X_{ij} - \bar{X})^2 = \sum_{i=1}^N \sum_{j=1}^{\bar{M}} (X_{ij} - \bar{X}_i + \bar{X}_i - \bar{X})^2 = \sum \sum (X_{ij} - \bar{X}_i)^2 + \sum \sum (\bar{X}_i - \bar{X})^2 + 2 \sum \sum (X_{ij} - \bar{X}_i)(\bar{X}_i - \bar{X})$$

$$\sum_{i=1}^N \sum_{j=1}^{\bar{M}} (X_{ij} - \bar{X})^2 \stackrel{=0}{=} \sum_{i=1}^N \sum_{j=1}^{\bar{M}} (X_{ij} - \bar{X}_i)^2 + \sum_{i=1}^N \sum_{j=1}^{\bar{M}} (\bar{X}_i - \bar{X})^2$$

$$\begin{aligned} \downarrow & \quad \downarrow & \quad \downarrow \\ (N\bar{M}-1) S_T^2 &= N(\bar{M}-1) S_W^2 + (N-1) S_b^2 \\ \underbrace{NM \sigma_T^2}_{\sigma^2} &= \underbrace{NM \sigma_W^2}_{\sigma_W^2} + \underbrace{(N-1) \sigma_b^2}_{\sigma_b^2/n} \end{aligned}$$

En el caso de b muestra:

$$\sum_{i=1}^n \sum_{j=1}^{\bar{M}} (X_{ij} - \bar{X})^2 = \sum \sum (X_{ij} - \bar{X}_i + \bar{X}_i - \bar{X})^2 = \sum_{j=1}^{\bar{M}} \sum (X_{ij} - \bar{X}_i)^2 + \sum \sum (\bar{X}_i - \bar{X})^2 + 2 \sum \sum (X_{ij} - \bar{X}_i)(\bar{X}_i - \bar{X})$$

$$\downarrow$$

$$(n\bar{M}-1) \hat{S}_T^2 =$$

$$\downarrow \quad \downarrow$$

$$n(\bar{M}-1) \hat{S}_W^2 + (n-1) \hat{S}_b^2$$

\hat{S}_b^2 es un estimador insesgado de S_b^2

\hat{S}_W^2 es un estimador insesgado de S_W^2

MIRAR
ATRÁS →

Un estimador insesgado de S^2 es:

$$S_o^2 = \frac{N(\bar{M}-1)}{N\bar{M}-1} \hat{S}_W^2 + \frac{N-1}{N\bar{M}-1} \hat{S}_b^2$$

aunque por muestras de más de 50 conglomerados se puede considerar como estimador insesgado de S^2 a:

$$\hat{S}_T^2 = \frac{1}{n\bar{M}-1} \sum \sum (X_{ij} - \bar{X})^2$$

b) CON reposición, probabilidades IGUALES:

El estimador lineal insesgado de θ es el estimador de Hansen y Hurwitz:

$$\hat{\theta}_{HH} = \sum_{i=1}^n \sum_{j=1}^M \frac{y_{ij}}{n p_i} \quad \text{probab. } p_i = 1/N$$

$$= \sum_{i=1}^n \sum_{j=1}^M \frac{y_{ij}}{n/N} = \frac{N}{n} \sum_{i=1}^n \sum_{j=1}^M y_{ij}$$

cuya expresión en probabilidades iguales coincide con el estimador de Horvitz y Thompson.

Varianza:

Para la media poblacional:

$$V(\bar{x}) = V\left(\frac{1}{n} \sum_{i=1}^n \bar{x}_i\right) = \frac{\sigma_x^2}{n} = \frac{1}{n} \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N} = \frac{1}{nN} \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{M} = \frac{\sigma_b^2}{nM}$$

donde $\sigma_b^2 = \frac{1}{N} \sum_{i=1}^N M (x_i - \bar{x})^2 \equiv \text{varianza entre conglomerados}$

(expresión similar al m.a.s. substituyendo σ^2 por σ_b^2 y n por nM , n.º totl de unidades elementales en la muestra).

Para el total poblacional:

$$V(\hat{X}) = V(NM \cdot \bar{x}) = N^2 M^2 V(\bar{x}) = N^2 M \cdot \frac{\sigma_b^2}{n}$$

Para la proporción:

$$V(\hat{P}) = \frac{\sigma_p^2}{nM} = \frac{\frac{1}{N} \sum_{i=1}^N (P_i - P)^2}{nM} = \frac{\sum_{i=1}^N (P_i - P)^2}{nN}$$

Para el total de clase:

$$V(\hat{A}) = V(NM \hat{P}) = N^2 M^2 V(\hat{P}) = N^2 M^2 \frac{\sum_{i=1}^N (P_i - P)^2}{nN}$$

Estimaciones

$$\hat{V}(\bar{x}) = \frac{\hat{\sigma}_b^2}{nM}$$

$$\hat{V}(\hat{X}) = N^2 M^2 \hat{V}(\bar{x})$$

$$\hat{V}(\hat{P}) = \frac{\hat{\sigma}_p^2}{nM}$$

$$V(\hat{A}) = N^2 M^2 \hat{V}(\hat{P})$$

3. COEF. de CORRELACIÓN INTRACONGLOMERADOS Y SU INTERPRETACIÓN. EFECTO de DISEÑO

El coeficiente de correlación intraconglomerados es una medida de la homogeneidad dentro de los conglomerados. Interesa que tenga un valor muy pequeño, pues en realidad por conglomerados lo ideal es la heterogeneidad dentro de los conglomerados.

Se define como el coef. de correlación lineal entre todos los pares de unidades elementales pertenecientes a cada conglomerado, para todos los conglomerados.

Para cada i, j, k

$$\delta = \frac{\text{cov}(X_{ij}, X_{ik})}{\sigma(X_{ij}) \sigma(X_{ik})} = \frac{E[(X_{ij} - E[X_{ij}])(X_{ik} - E[X_{ik}])]}{\sigma^2}$$

$\sqrt{\frac{n^2 \text{ pares} = N \cdot M \cdot (M-1)}{M}}$

$$\rightarrow \frac{\sum_{i=1}^N \sum_{j \neq k}^M (X_{ij} - \bar{X})(X_{ik} - \bar{X})}{N \bar{M} (\bar{M} - 1) \sigma^2} = \frac{\sum_{i=1}^N \sum_{j \neq k}^M (X_{ij} - \bar{X})(X_{ik} - \bar{X})}{(\bar{M} - 1)(N \bar{M} - 1) S^2}$$

$\sigma^2 = \frac{N \bar{M} - 1}{N \bar{M}} S^2$

Utilizando el coef. de correlación intraconglomerados se puede expresar la varianza de los estimadores:

(SR) $\rightarrow V(\bar{X}) = (1-f) \cdot \frac{S_b^2}{n \bar{M}} = \frac{(1-f)}{n \bar{M}} \cdot \frac{M}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2 =$

$S_b^2 \approx S^2 [1 + (\bar{M}-1)\delta]$

$$= \frac{1-f}{n(N-1)} \sum_{i=1}^N \left[\frac{1}{\bar{M}} \sum_{j=1}^M (X_{ij} - \bar{X}) \right]^2 = \frac{1-f}{n(N-1) \bar{M}^2} \sum_{i=1}^N \left(\sum_{j=1}^M (X_{ij} - \bar{X}) \right)^2$$

$$= \frac{1-f}{n(N-1) \bar{M}^2} \left[\sum_i \sum_j (X_{ij} - \bar{X})^2 + \sum_i \sum_{j \neq k} (X_{ij} - \bar{X})(X_{ik} - \bar{X}) \right] =$$

$$= \frac{1-f}{n(N-1) \bar{M}^2} \left[(N \bar{M} - 1) S^2 + (N \bar{M} - 1)(\bar{M} - 1) S^2 \delta \right] =$$

$$= \frac{(1-f) S^2 (N \bar{M} - 1)}{n(N-1) \bar{M}^2} [1 + (\bar{M} - 1) \delta] \xrightarrow{N \rightarrow \infty} \frac{(1-f) S^2}{n \bar{M}} [1 + (\bar{M} - 1) \delta]$$

$$(SR) \rightarrow V(\bar{x}) = (1-f) \cdot \underbrace{\frac{S^2}{n\bar{M}}}_{S_{mas}^2} [1 + (\bar{M}-1)\delta]$$

que se puede expresar en función de la varianza de la media muestral obtenida con u.a.s.

$$V(\bar{x}) = V_{mas}(\bar{x}) [1 + (\bar{M}-1)\delta]$$

por lo que podemos hacer una comparación de varianzas a partir de la interpretación de δ :

$\delta > 0 \Rightarrow$ Var. conglomerados es mayor que Var. mas. para muestras del mismo tamaño $n\bar{M}$.

Esta diferencia será máxima para $\delta = 1$, caso más desfavorable. El caso más favorable, varianzas mínimas será cuando $\delta = -\frac{1}{\bar{M}-1}$, en el que la var. de conglomerados será 0.

Para $\delta = 0$, ambos métodos son igual de precisos.

El término $(\bar{M}-1)$ se interpreta como el aumento de varianzas debido a seleccionar n conglomerados de tamaño \bar{M} , en vez de hacer un u.a.s. de $n\bar{M}$ unidades elementales.

$\delta < 0 \Rightarrow$ Var. conglomerados es menor que la varianzas obtenida con u.a.s.

En la práctica, suele ocurrir que los elementos de cada conglomerado guarden cierto parecido entre sí, con lo que la correlación es positiva y el muestreo por conglomerados es menos preciso que el u.a.s.

$$V_{mc}(\bar{x}) = V_{mas}(\bar{x}) [1 + (\bar{M}-1)\delta] \Rightarrow \begin{cases} \delta > 0 \Rightarrow \text{Congl. peor que mas} \\ \delta = 0 \Rightarrow \text{Congl. igual que mas} \\ \delta < 0 \Rightarrow \text{Congl. mejor que mas} \end{cases}$$

Si llamamos n_a y n_c al tamaño de una muestra, expresado en unidades elementales para obtener una precisión dada, en el caso en el que los dos tipos de muestras tengan la misma precisión:

$$V_{nc}(\bar{x}) = V_{mas}(\bar{x}) [1 + (\bar{M} - 1)\delta] = V_{mas}(\bar{x})$$

~~$$(1-f) \frac{S_b^2}{n_c} = (1-f) \frac{S^2}{n_a}$$~~

por lo que:

$$(1-f) \frac{S_b^2}{n_c} [1 + (\bar{M} - 1)\delta] = (1-f) \frac{S^2}{n_a}$$

$$\Rightarrow n_c = n_a [1 + (\bar{M} - 1)\delta],$$

es decir, la cantidad $[1 + (\bar{M} - 1)\delta]$ es aquella por la que hay que multiplicar el tamaño muestral de mas para obtener el tamaño muestral de conglomerados. Se llama efecto de diseño.

Podemos interpretar los valores de δ utilizando toda la información anterior:

Si $V_{nc}(\bar{x}) \simeq V_{mas}(\bar{x}) [1 + (\bar{M} - 1)\delta]$, se puede escribir como

$$(1-f) \frac{S_b^2}{n\bar{M}} \simeq (1-f) \frac{S^2}{n\bar{M}} [1 + (\bar{M} - 1)\delta], \text{ de donde:}$$

$$S_b^2 \simeq S^2 [1 + (\bar{M} - 1)\delta] \longrightarrow S_b^2 - S^2 \simeq (\bar{M} - 1) S^2 \cdot \delta$$

$$\delta \simeq \frac{S_b^2 - S^2}{(\bar{M} - 1) S^2} \in \left[-\frac{1}{\bar{M} - 1}, 1 \right]$$

$$\bullet \delta = -\frac{1}{\bar{M} - 1} \Rightarrow S_b^2 = 0$$

$$\Rightarrow V(\bar{x}) = 0$$

$$\delta \simeq \frac{\frac{S_b^2}{S^2} - 1}{(\bar{M} - 1)}$$

L

\Rightarrow toda la variabilidad procede de dentro de los conglomerados, que son homogéneos entre sí
 \Rightarrow 1 solo conglomerado proporciona toda la información \rightarrow caso ideal.

$$\bullet -\frac{1}{M-1} < \delta < 0 \Rightarrow V_{MC}(\bar{x}) < V_{MAS}(\bar{x}) \quad \text{ó} \quad n_c < n_a$$

\Rightarrow El muestreo por conglomer. es + preciso que el mas.

\rightarrow Caso excepcional, conglomerados heterogéneos dentro

$$\bullet \delta = 0 \Rightarrow S_b^2 = S^2, \quad V_{MC}(\bar{x}) = V_{MAS}(\bar{x}) \quad \text{ó} \quad n_c = n_a.$$

\rightarrow Los dos tipos de muestreo son igual de precisos.

\rightarrow La variabilidad entre conglomerados coincide con la variabilidad entre las unidades elementales.

$$\bullet 0 < \delta < 1 \Rightarrow [1 + (\bar{M}-1)\delta] > 1 \Rightarrow V_{MC}(\bar{x}) > V_{MAS}(\bar{x}) \quad \text{ó} \quad n_c > n_a$$

\rightarrow Existe homogeneidad dentro de los conglomerados (caso + habitual), y $S_b^2 > S^2$.

\rightarrow Aunque el m.c. es menor preciso, se suele utilizar por motivos de coste.

$$\bullet \delta = 1 \Rightarrow S_b^2 \simeq \bar{M} S^2, \quad S_w^2 = 0 \quad \text{y} \quad n_c = \bar{M} n_a$$

\rightarrow Caso más desfavorable, no hay variabilidad dentro de los conglomerados \Rightarrow se recomienda observar 1 sola observación de cada conglomerado.

$$V(\bar{x}) \simeq (1-f) \cdot \frac{S^2}{n\bar{M}} [1 + (\bar{M}-1) \cdot 1] = (1-f) \underbrace{\frac{S^2}{n}}_{\text{m.a.s. con unidades elementales}}$$

muy superior a la var_{MAS} obtenida con $n\bar{M}$ unidades elementales

Ⓢ En la EPA, $[1 + (\bar{M}-1)\delta] \simeq 2$, pero reduce coste

$$(CR) \rightarrow \delta = \frac{\sum_{i=1}^N \sum_{j \neq z}^{\bar{M}} (X_{ij} - \bar{X})(X_{iz} - \bar{X})}{N\bar{M}(\bar{M}-1) \sigma^2}$$

$$\text{donde } \sigma^2 = \frac{1}{N\bar{M}} \sum_{i=1}^N \sum_{j=1}^{\bar{M}} (X_{ij} - \bar{X})^2$$

por lo que podemos expresar la variancia del estimador,

$$\begin{aligned} V(\bar{\bar{X}}) &= \frac{\frac{1}{N} \sum_{i=1}^N (\bar{X}_i - \bar{X})^2}{n} = \frac{1}{nN} \sum_{i=1}^N \left[\frac{1}{\bar{M}} \sum_{j=1}^{\bar{M}} X_{ij} - \frac{\bar{M} \cdot \bar{X}}{\bar{M}} \right]^2 = \\ &= \frac{1}{nN\bar{M}^2} \sum_{i=1}^N \left(\sum_{j=1}^{\bar{M}} X_{ij} - \sum_{j=1}^{\bar{M}} \bar{X} \right)^2 = \frac{1}{nN\bar{M}^2} \sum_{i=1}^N \left(\sum_{j=1}^{\bar{M}} (X_{ij} - \bar{X}) \right)^2 = \\ &= \frac{1}{nN\bar{M}^2} \left[\sum_{i=1}^N \sum_{j=1}^{\bar{M}} (X_{ij} - \bar{X})^2 + \sum_{i=1}^N \sum_{j \neq z}^{\bar{M}} (X_{ij} - \bar{X})(X_{iz} - \bar{X}) \right] = \\ &= \frac{1}{nN\bar{M}^2} \left[N\bar{M}\sigma^2 + N\bar{M}(\bar{M}-1)\sigma^2 \cdot \delta \right] = \\ &= \frac{N\bar{M}\sigma^2}{nN\bar{M}^2} \left[1 + (\bar{M}-1)\delta \right] = \frac{\sigma^2}{n\bar{M}} \left[1 + (\bar{M}-1)\delta \right] \end{aligned}$$

$$(CR) \rightarrow V(\bar{\bar{X}}) = \frac{\sigma^2}{n\bar{M}} \left[1 + (\bar{M}-1)\delta \right]$$

de expresión de la que se deducen las mismas conclusiones que en el caso de SIN reposición:

$$\bullet \delta \in \left[-\frac{1}{\bar{M}-1}, 1 \right]$$

$$\bullet V_{MC}(\bar{\bar{X}}) = V_{MAS}(\bar{X}) \left[1 + (\bar{M}-1)\delta \right] \Rightarrow \begin{cases} \delta < 0 \rightarrow \text{compl. Mejor mas} \\ \delta = 0 \rightarrow \text{igual} \\ \delta > 0 \rightarrow \text{compl. Peor mas} \end{cases}$$

y cuanto más se acerque a los límites, más se acerque la variancia / pérdida en precisión del MC respecto al MAS

4. UTILIZACIÓN de ESTIMADORES de RAZÓN

Hasta ahora hemos considerado el caso más sencillo, todos los conglomerados tienen el mismo tamaño ($M_i = \bar{M}$). Pero lo habitual es que los conglomerados tengan \neq tamaño.

a) Si M_i son similares,

podemos considerar $\bar{M} = \frac{1}{N} \sum_{i=1}^N M_i$ como la media de todos los tamaños de los conglomerados y utilizar la fórmula estudiada hasta ahora.

En el caso de muestreo sin reposición con probab. iguales $\bar{X} = \frac{1}{n} \sum_{i=1}^n \bar{X}_i = \frac{1}{n\bar{M}} \sum_{i=1}^n X_i$, estimador insesgado de \bar{X}

$$V(\bar{X}) = \frac{(1-f)}{n\bar{M}} S_b^2 = \frac{1-f}{n\bar{M}} \sum_{i=1}^N \frac{(X_i - \bar{X})^2}{(N-1)\bar{M}} = \frac{1-f}{n\bar{M}^2} \sum_{i=1}^N \frac{(X_i - \bar{X})^2}{N-1}$$

$$\hat{V}(\bar{X}) = \frac{(1-f)}{n\bar{M}} \hat{S}_b^2 = \dots = \frac{1-f}{n\bar{M}^2} \sum_{i=1}^N \frac{(X_i - \bar{X})^2}{n-1}$$

b) Si M_i son muy distintos,

la varianza del estimador tomando la media muestral de los tamaños puede ser muy grande si los tamaños de los conglomerados difieren mucho entre sí, porque también habrá una variabilidad alta entre los totales $X_i = \sum_{j=1}^{M_i} X_{ij}$.

En este caso, la precisión puede mejorarse utilizando el estimador de razón \hat{X}_R , que aunque es sesgado, puede ser más acurado.

$$\hat{\bar{X}}_R = \bar{\bar{X}}_R = \hat{R} = \frac{\sum_{i=1}^n X_i}{\sum_{i=1}^n M_i}, \text{ estimador de } R = \frac{\sum_{i=1}^N X_i}{\sum_{i=1}^N M_i} = \bar{\bar{X}}$$

Por ser un estimador de razón, su varianza aproximada (para n m.f.c. grande) es:

$$(SR) \rightarrow V(\hat{\bar{X}}) = V(\hat{R}) = V(\hat{R} - R) = V\left(\frac{\sum_{i=1}^n X_i}{\sum_{i=1}^n M_i} - R\right) =$$

$$= V\left(\frac{\hat{X}}{\hat{M}} - R\right) = V\left(\frac{\hat{X} - R\hat{M}}{\hat{M}}\right) = \frac{V(\hat{X} - R\hat{M})}{\hat{M}^2} = \frac{V(N\hat{X} + R\hat{M})}{M^2 \hat{M}}$$

$$\approx \frac{N^2}{M^2} \cdot \frac{1-f}{n} \cdot \frac{1}{N-1} \sum_{i=1}^N (X_i - RM_i)^2$$

Tuiz
C. Pérez

$$V(\hat{R}) = \frac{N^2(1-f)}{n M^2} \cdot \frac{1}{N-1} \sum_{i=1}^N M_i^2 (\bar{X}_i - \bar{\bar{X}})^2$$

$$\hat{V}(\hat{R}) = \frac{N^2}{M^2} \cdot \frac{1-f}{n} \cdot \frac{1}{N-1} \sum_{i=1}^N M_i^2 (\bar{X}_i - \bar{\bar{X}})^2$$

Para el total: $\hat{X} = M\hat{R} \Rightarrow V(\hat{X}) = M^2 V(\hat{R})$

Para la proporción: $V(\hat{P}) = \frac{N^2}{M^2} \cdot \frac{1-f}{n} \cdot \frac{1}{N-1} \sum_{i=1}^N M_i^2 (P_i - \bar{P})^2$

$$\hat{V}(\hat{P}) = \frac{N^2}{M^2} \cdot \frac{1-f}{n} \cdot \frac{1}{N-1} \sum_{i=1}^N M_i^2 (P_i - \bar{P})^2$$

Para el total de clase: $\hat{A} = M\hat{P} \Rightarrow V(\hat{A}) = M^2 V(\hat{P})$

$$(CR) \rightarrow V(\hat{\bar{X}}) = V(\hat{R}) = \frac{N^2}{M^2} \cdot \frac{1}{n} \cdot \frac{1}{N} \sum_{i=1}^N M_i^2 (\bar{X}_i - \bar{\bar{X}})^2$$

$$\hat{V}(\hat{\bar{X}}) = \frac{N^2}{M^2} \cdot \frac{1}{n} \cdot \frac{1}{N-1} \sum_{i=1}^N M_i^2 (\bar{X}_i - \bar{\bar{X}})^2$$

Para el total: $\hat{X} = M\hat{\bar{X}} \Rightarrow V(\hat{X}) = M^2 V(\hat{\bar{X}})$

Para la proporción: $V(\hat{P}) = \frac{N^2}{M^2} \cdot \frac{1}{n} \cdot \frac{1}{N} \sum_{i=1}^N M_i^2 (P_i - \bar{P})^2$

$$\hat{V}(\hat{P}) = \frac{N^2}{M^2} \cdot \frac{1}{n} \cdot \frac{1}{N-1} \sum_{i=1}^N M_i^2 (P_i - \bar{P})^2$$

Para el total de clase: $\hat{A} = M\hat{P} \Rightarrow V(\hat{A}) = M^2 V(\hat{P})$

Comentari en ①

En el caso de probab. desiguales, habría que acudir a la fórmula general de los estimadores lineales inses-
gados de Horvitz y Thompson ($\pi_i = P(\text{comp. i.º muestro})$) y
Hansen y Horvitz ($P_i = P(\text{comp. i.º muestro})$).

Los métodos más interesantes eran aquellos con probab. proporcionales a los tamaños M_i , en cuyo caso

$$(SR) \rightarrow \hat{\bar{X}} = \frac{\hat{X}_{HH}}{M} = \frac{M \bar{\bar{X}}}{M} = \bar{\bar{X}}$$

\hookrightarrow la expresión del estimador coincide con la de prob.

Los valores de las varianzas y sus estimadores dependen del valor de π_{ij} en cada mt. de selección.

$$(CR) \rightarrow \hat{\bar{X}} = \frac{\hat{X}_{HH}}{M} = \frac{M \bar{\bar{X}}}{M} = \bar{\bar{X}}$$

\hookrightarrow la expresión es igual.

$$V(\hat{\bar{X}}_{HH}) = \frac{1}{nM} \sum_{i=1}^N M_i (\bar{X}_i - \bar{\bar{X}})^2$$

$$\hat{V}(\hat{\bar{X}}_{HH}) = \frac{1}{n(n-1)} \sum (\bar{X}_i - \bar{\bar{X}})^2$$

L