

# PROBAB. DESIGUALES

(1)

**PROBAB**  $e_i = 0^o$  veces que aparece  $u_i$  en  $k$  muestra.

SR  $\rightarrow e_i = \{0, 1\}$ ,  $e_i \rightarrow B(1, \pi_i)$   $\left\{ \begin{array}{l} E[e_i] = \pi_i \\ V[e_i] = \pi_i(1-\pi_i) \\ \text{cov}[e_i, e_j] = \pi_{ij} - \pi_i \pi_j \end{array} \right.$

Propiedades: P1 -  $\sum_{i=1}^N \pi_i = n$

P2 -  $\sum_{i \neq j}^N \pi_i = n - \pi_j$

P3 -  $\sum_{i \neq j}^N \pi_{ij} = (n-1)\pi_j$

P4 -  $\sum_{i \neq j}^N (\pi_{ij} - \pi_i \pi_j) = -\pi_j(1-\pi_j)$

para  $u_j$   
fijo

CR  $\rightarrow e_i = \{0, 1, \dots, n\}$ ,  $e_i \rightarrow B(n, \pi_i)$   $\left\{ \begin{array}{l} E[e_i] = n\pi_i \\ V[e_i] = n\pi_i(1-\pi_i) \\ E[e_i e_j] = n(n-1)\pi_i \pi_j \\ \text{cov}[e_i, e_j] = -n\pi_i \pi_j \end{array} \right.$

**ESTIMADOR**  $\theta = \sum_{i=1}^N y_i \rightarrow \hat{\theta} = \sum_{i=1}^N y_i w_i$

SR  $\rightarrow \hat{\theta}_{HT} = \sum_{i=1}^N \frac{y_i}{\pi_i}$

CR  $\rightarrow \hat{\theta}_{HH} = \sum_{i=1}^N \frac{y_i}{n\pi_i}$

$\left[ \begin{array}{l} \hat{X}_{HT} = \sum_{i=1}^N \frac{x_i}{\pi_i} \rightarrow \hat{X}_{HT} = \frac{1}{N} \hat{K}_{HT} \\ \hat{A}_{HT} = \sum_{i=1}^N \frac{A_i}{\pi_i} \rightarrow \hat{P}_{HT} = \frac{1}{N} \hat{A}_{HT} \\ \hat{X}_{HH} = \sum_{i=1}^N \frac{x_i}{n\pi_i} \rightarrow \hat{X}_{HH} = \frac{1}{N} \hat{X}_{HH} \\ \hat{A}_{HH} = \sum_{i=1}^N \frac{A_i}{n\pi_i} \rightarrow \hat{P}_{HH} = \frac{1}{N} \hat{A}_{HH} \end{array} \right.$

**VARIANZA**  $\hat{\theta}$  insesgado  $\Rightarrow$  Precisión  $\sim$  Acurciedad

SR  $\rightarrow V(\hat{\theta}_{HT}) = \sum_{i=1}^N \frac{y_i^2}{\pi_i} (1-\pi_i) + 2 \sum_{i=1}^N \sum_{j>i}^N \frac{y_i y_j}{\pi_i \pi_j} (\pi_{ij} - \pi_i \pi_j)$

CR  $\rightarrow V(\hat{\theta}_{HH}) = \frac{1}{n} \left[ \sum_{i=1}^N \frac{y_i^2}{\pi_i} - \theta^2 \right] \sim \frac{1}{n} \sum_{i=1}^N \left( \frac{y_i}{\pi_i} - \theta \right)^2 \pi_i$

**ESTIMADOR VARIANZA** Intersección

SR  $\xrightarrow{\text{Dividido por } \pi} \hat{V}(\hat{\theta}_{HT}) = \sum_{i=1}^N \frac{y_i^2}{\pi_i^2} (1-\pi_i) + 2 \sum_{i=1}^N \sum_{j>i}^N \frac{y_i y_j}{\pi_i \pi_j} \frac{\pi_{ij} - \pi_i \pi_j}{\pi_j}$

Waller y Gaudy: YG  $\rightarrow \hat{V}(\hat{\theta}_{HT}) = \sum_{i=1}^N \sum_{j>i}^N \left( \frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right)^2 \frac{\pi_{ij} - \pi_i \pi_j}{\pi_j}$

CR  $\rightarrow \hat{V}(\hat{\theta}_{HH}) = \frac{1}{n(n-1)} \left[ \sum_{i=1}^N \left( \frac{y_i}{\pi_i} \right)^2 - n \hat{\theta}_{HH}^2 \right] \sim \frac{1}{n(n-1)} \sum_{i=1}^N \left( \frac{y_i}{\pi_i} - \hat{\theta}_{HH} \right)^2$   
dividido por  $(n-1)\pi_i$

# PROBABILIDADES OPTIMAS

SR  $\rightarrow \pi_i = K M_i \Rightarrow \pi_i = \frac{n}{M} \cdot M_i$  (prop. a  $M_i$ )

CR  $\rightarrow P_i = K \pi_i \Rightarrow P_i = \frac{1}{M} \cdot M_i$  (prop. a  $\pi_i$ )

## MTS. ESPECIALES SIN REPOSIC Y PPT

Una Generalizaci3n  $\rightarrow U_i \sim \pi_i$  bolas color  $i$ . 1<sup>a</sup> extr. sale color  $i \rightarrow U_i$  muestra  $\rightarrow$  se sacan  $\pi_i$  bolas

$n=1 \rightarrow \pi_i = \frac{M_i}{M}$   
 $n=2 \rightarrow \pi_i = \frac{M_i}{M} \left( 1 + \sum_{j \neq i} \frac{M_j}{M - \pi_j} \right)$

Brewer (1968)  $\rightarrow$  Par  $n=2$ ;

$n=1 \rightarrow P(U_i \in \text{muestra}) = \pi_i = \frac{K_i}{K} = \frac{P_i(1-P_i)/(1-2P_i)}{\sum_j K_j}$   $P_i = \frac{\pi_i}{\pi} / P_i < 1/2$

$n=2 \rightarrow \pi_i = 2P_i$  (sin reposic.)

$\rightarrow \pi_{ij} \equiv$  proporc. a  $P_i$  y a  $P_j$

$\rightarrow \hat{\theta}_{HT}$

Durbin (1967)  $\rightarrow$  Par  $n=2$

$n=1$ , sin repos.  $P_i = \frac{M_i}{M}$

$n=2$ ,  $\pi_i$  prop. a  $K_j = P_i \left( \frac{1}{1-2P_i} + \frac{1}{1-2P_j} \right) / P_j < \frac{1}{2}$ .

$\rightarrow \pi_i, \pi_{ij}, \hat{\theta}_{HT}$  coinciden con Brewer

Esquema mixto SCG  $\rightarrow U_i \sim \pi_i$  bolas color  $i$ .

$j$ -extr. sale bola color  $i \Rightarrow$  s3lo se saca esa bola (pueden  $\pi_i - 1$ )  
 $\Rightarrow U_i$  puede elegirse  $0, 1, \dots, \min\{n, M_i\}$  veces  
 $\Rightarrow$  probab. gradualmente variables

Esquema mixto  $\rightarrow$  SR para bola  
 $\rightarrow$  CR para unidades

$e_i \rightarrow H(M, n, P_i)$

$\hat{\theta}_{SCG} = \sum \frac{Y_i}{n P_i} = \hat{\theta}_{HH}$ , pero  $V(\hat{\theta}_{SCG}) \neq V(\hat{\theta}_{HH})$  y  
 $\hat{V}(\hat{\theta}_{SCG}) = \frac{M-n}{M} \hat{V}(\hat{\theta}_{HH}) \Rightarrow$  + preciso

# MUEST\_T3. MUESTREO CON PROBAB. DESIGUALES.

ESTIMADORES LINEALES.

VARIANZAS de los estimadores y sus ESTIMACIONES.

~~COMPARACIÓN entre el muestreo~~

PROBABILIDADES ÓPTIMAS de SELECCIÓN.

ESTIMADORES ESPECIALES de SELECCIÓN

SIN reposición y PROBAB. PROPORCIONALES al tamaño

## 1. MUESTREO CON PROBAB. DESIGUALES

Muestreo probabilístico en el que la probabilidad de que una unidad poblacional,  $u_i$ ,  $i=1 \dots N$ , forme parte de la muestra es distinta para cada unidad (en cada extracción)

$$P(u_i \in M) \neq P(u_j \in M)$$

Muestreo idóneo para unidades compuestas, en el que la importancia de cada unidad es proporcional al n° de unid. elementales que la componen (probab. <sup>(PPT)</sup> proporcionales al tamaño) ó bien con unidades elementales con distinta importancia. ( $\neq w_i$ )

Estas probabilidades de pertenecer a la muestra dependen del tipo de muestreo, con o sin reposición.

Yo creo que el n° de muestras posibles será igual que en el caso de probab. iguales, aunque la probab. de la muestra sea  $\neq$  para cada muestra.

#S	NO import orden	SÍ import orden
SIN repos	$C_{N,n} = \binom{N}{n}$	$V_{N,n} = \binom{N}{n} n!$
CON repos	$CR_{N,n} = \binom{N+n-1}{n}$	$VR_{N,n} = N^n$

## 2 - MUESTREO SIN REPOSICIÓN y PROBAB. DESIGUALES

Procedimiento de selección de la muestra:

- con probabilidades desiguales,
- se seleccionan las unidades de la muestra 1 a 1 sin reposición,
- no importa el orden de aparición.

Por lo que las muestras obtenidas mediante este procedimiento se caracterizan por:

- muestras con elementos repetidos son imposibles,
- muestras no equiprobables (pq las unid.  $u_i$  tienen  $\neq$  probab. de pertenecer a la muestra).

### PROBABILIDADES

Dada la población  $\{u_1, \dots, u_N\}$  de la que seleccionamos sin reposición una muestra de tamaño  $n$ , para cada elemento poblacional  $u_i$  definiremos:

$$e_i = \begin{cases} 1 & \text{si } u_i \in \text{Muestra, con probab. } \pi_i \\ 0 & \text{si } u_i \notin \text{Muestra, con probab. } 1 - \pi_i \end{cases}$$

$e_i$	$P(e_i)$
1	$\pi_i$
0	$1 - \pi_i$

$e_i$  es una v.a. que mide el nº de veces que cada unidad poblacional aparece en la muestra, definida a través de su función de probabilidad.

$e_i \rightarrow B(1, \pi_i)$  por lo que:

$$E[e_i] = \pi_i = 1 \cdot \pi_i + 0 \cdot (1 - \pi_i)$$

$$E[e_i^2] = \pi_i = 1^2 \cdot \pi_i + 0^2 \cdot (1 - \pi_i) = \pi_i$$

$$V[e_i] = pq = \pi_i(1 - \pi_i) = E[e_i^2] - E[e_i]^2 = \pi_i - \pi_i^2$$

Para cada  $i \neq j$  consideramos la ocurrencia conjunta a través de la v.a. auxiliar

$$e_i \cdot e_j = \begin{cases} 1 & \text{si } (u_i, u_j) \in \text{Muestra con probab } \pi_{ij} \\ 0 & \text{si } (u_i, u_j) \notin \text{Muestra con probab } 1 - \pi_{ij} \end{cases}$$

$$E[e_i \cdot e_j] = 1 \cdot \pi_{ij} + 0 \cdot (1 - \pi_{ij}) = \pi_{ij}$$

$$\text{cov}[e_i, e_j] = E[e_i \cdot e_j] - E[e_i]E[e_j] = \pi_{ij} - \pi_i \pi_j$$

ΔTRAS  
→

ESTIMADOR lineal insesgado

En el caso de que la característica poblacional que se quiere estudiar sobre los elementos poblacionales  $u_i, i=1, \dots, N$ , sea el parámetro  $\theta = \sum_{i=1}^N y_i$

un buen estimador lineal es

$$\hat{\theta} = \sum_{i=1}^N w_i y_i$$

Para que este estimador sea insesgado  $E[\hat{\theta}] = \theta$ , en el caso de m.a.p.d.s.r., la ponderación  $w_i = \frac{1}{\pi_i}$ , s.r.  $E[\hat{\theta}] = E\left[\sum_{i=1}^N w_i y_i\right] = E\left[\sum_{i=1}^N w_i y_i e_i\right] = \sum_{i=1}^N w_i y_i E[e_i] = \sum_{i=1}^N w_i y_i \pi_i$

$$E[\hat{\theta}] = \theta \Leftrightarrow w_i \pi_i = 1 \Leftrightarrow w_i = \frac{1}{\pi_i}, \forall i.$$

Luego el estimador lineal insesgado para  $\theta$  en el caso de muestreo sin reposición es el estimador de Horvitz y Thompson (1952), cuya expresión es:

$$\hat{\theta}_{HT} = \sum_{i=1}^N \frac{y_i}{\pi_i}$$

Particularizando para los parámetros poblacionales más comunes:

Parámetro	Estimador Horvitz y Thompson
Total poblacional	$\theta = X = \sum_{i=1}^N x_i \quad y_i = x_i \quad \hat{\theta}_{HT} = \sum_{i=1}^N \frac{x_i}{\pi_i}$
Media poblacional	$\theta = \bar{X} = \sum_{i=1}^N \frac{x_i}{N} \quad y_i = \frac{x_i}{N} \quad \hat{\theta}_{HT} = \frac{1}{N} \sum_{i=1}^N \frac{x_i}{\pi_i} = \frac{1}{N} \cdot \hat{X}_{HT}$
Total de clase	$\theta = A = \sum_{i=1}^N A_i \quad y_i = A_i \quad \hat{\theta}_{HT} = \sum_{i=1}^N \frac{A_i}{\pi_i}$
Proporción de clase	$\theta = P = \sum_{i=1}^N \frac{A_i}{N} \quad y_i = \frac{A_i}{N} \quad \hat{\theta}_{HT} = \frac{1}{N} \sum_{i=1}^N \frac{A_i}{\pi_i} = \frac{1}{N} \cdot \hat{A}_{HT}$

## VARIANZA del estimador de Horvitz y Thompson

La precisión de un estimador se analiza a partir de tres conceptos:

- error de muestreo,  $G(\hat{\theta}) = \sqrt{V(\hat{\theta})}$
- acuracidad,  $ECM(\hat{\theta}) = E[\hat{\theta} - \theta]^2 = V(\hat{\theta}) + b(\hat{\theta})^2$
- sesgo,  $b(\hat{\theta}) = E[\hat{\theta}] - \theta$

En el caso de que el estimador sea insesgado,  $E[\hat{\theta}] = \theta$  resulta que  $b(\hat{\theta}) = 0$  y por tanto  $ECM(\hat{\theta}) = V(\hat{\theta})$ , por lo que bastará con medir la varianza del estimador para estudiar su precisión y su acuracidad.

En el caso de muestreo sin reposición, la varianza del estimador de Horvitz y Thompson es:

$$\begin{aligned}
 V(\hat{\theta}_{HT}) &= V\left(\sum_{i=1}^N \frac{y_i}{\pi_i}\right) = V\left(\sum_{i=1}^N \frac{y_i}{\pi_i} e_i\right) = \sum_{i=1}^N V\left(\frac{y_i}{\pi_i} e_i\right) + \\
 &\quad + 2 \sum_{i=1}^N \sum_{j>i}^N \text{cov}\left(\frac{y_i}{\pi_i} e_i, \frac{y_j}{\pi_j} e_j\right) = \\
 &\quad = \sum_{i=1}^N \left(\frac{y_i}{\pi_i}\right)^2 V(e_i) + 2 \sum_{i=1}^N \sum_{j>i}^N \frac{y_i}{\pi_i} \cdot \frac{y_j}{\pi_j} \text{cov}(e_i, e_j) = \\
 &\quad = \sum_{i=1}^N \frac{y_i^2}{\pi_i^2} \cdot \pi_i(1-\pi_i) + 2 \sum_{i=1}^N \sum_{j>i}^N \frac{y_i}{\pi_i} \cdot \frac{y_j}{\pi_j} (\pi_{ij} - \pi_i \pi_j) \\
 V(e_i) &= \pi_i(1-\pi_i) \\
 \text{cov}(e_i, e_j) &= \pi_{ij} - \pi_i \pi_j
 \end{aligned}$$

$$V(\hat{\theta}_{HT}) = \sum_{i=1}^N \frac{y_i^2}{\pi_i} \frac{1-\pi_i}{\pi_i} + 2 \sum_{i=1}^N \sum_{j>i}^N \frac{y_i}{\pi_i} \cdot \frac{y_j}{\pi_j} (\pi_{ij} - \pi_i \pi_j)$$

## ESTIMACIÓN de la VARIANZA del estimador de Horvitz y Thompson

La expresión de  $V(\hat{\theta}_{HT})$  depende de toda la unid. poblaciones,  $i=1 \dots N$ , y solamente tenemos información de las unidades muestrales,  $i=1 \dots n$ . Por ello, se hace necesario estimar la varianza del estimador en función de la información muestral.

un estimador insesgado de  $V(\hat{\theta}_{HT})$  es:

$$\hat{V}(\hat{\theta}_{HT}) = \sum_{i=1}^n \frac{y_i^2}{\pi_i^2} (1-\pi_i) + 2 \sum_{i=1}^n \sum_{j>i}^n \frac{y_i}{\pi_i} \cdot \frac{y_j}{\pi_j} \cdot \frac{\pi_{ij} - \pi_i \pi_j}{\pi_{ij}}$$

Efectivamente, es insesgado,  $E[\hat{V}(\hat{\theta}_{HT})] = V(\hat{\theta}_{HT})$ .

$$\begin{aligned} E[\hat{V}(\hat{\theta}_{HT})] &= E\left[\sum_{i=1}^n \frac{y_i^2}{\pi_i^2} (1-\pi_i) + 2 \sum_{i=1}^n \sum_{j>i}^n \frac{y_i}{\pi_i} \cdot \frac{y_j}{\pi_j} \cdot \frac{\pi_{ij} - \pi_i \pi_j}{\pi_{ij}}\right] = \\ &= E\left[\sum_{i=1}^n \frac{y_i^2}{\pi_i^2} (1-\pi_i) e_i\right] + 2 E\left[\sum_{i=1}^n \sum_{j>i}^n \frac{y_i}{\pi_i} \cdot \frac{y_j}{\pi_j} \cdot \frac{\pi_{ij} - \pi_i \pi_j}{\pi_{ij}} e_i e_j\right] = \\ &= \sum_{i=1}^n \frac{y_i^2}{\pi_i^2} (1-\pi_i) \underbrace{E[e_i]}_{\pi_i} + 2 \sum_{i=1}^n \sum_{j>i}^n \frac{y_i}{\pi_i} \cdot \frac{y_j}{\pi_j} \cdot \frac{\pi_{ij} - \pi_i \pi_j}{\pi_{ij}} \cdot \underbrace{E[e_i e_j]}_{\pi_{ij}} = \\ &= V(\hat{\theta}_{HT}). \end{aligned}$$

Este estimador, aunque es el más utilizado en la práctica, es algo inestable (puede proporcionar valores negativos). Otro estimador más estable se debe a

Yates y Grundy:

$$\hat{V}(\hat{\theta}_{HT}) = \sum_{i=1}^n \sum_{j>i}^n \left( \frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right)^2 \cdot \frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}} \quad \leftarrow Y-G$$

También es insesgado

$$\begin{aligned} \text{Dem: } E\left[\sum_{i=1}^n \sum_{j>i}^n \left( \frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right)^2 \cdot \frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}}\right] &= E\left[\sum_{i=1}^n \sum_{j>i}^n \left( \frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right)^2 \cdot e_i e_j \cdot \frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}}\right] = \\ &= \sum_{i=1}^n \sum_{j>i}^n \left( \frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right)^2 \cdot \frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}} \cdot \pi_{ij} = \\ &= \sum_{i=1}^n \sum_{j>i}^n \left( \frac{y_i^2}{\pi_i^2} + \frac{y_j^2}{\pi_j^2} \right) (\pi_i \pi_j - \pi_{ij}) - 2 \sum_{i=1}^n \sum_{j>i}^n \frac{y_i y_j}{\pi_i \pi_j} (\pi_i \pi_j - \pi_{ij}) = \end{aligned}$$

Desarrollando el cuadrado:

$$\begin{aligned} &= \sum_{i=1}^n \sum_{j>i}^n \left( \frac{y_i^2}{\pi_i^2} + \frac{y_j^2}{\pi_j^2} \right) (\pi_i \pi_j - \pi_{ij}) - 2 \sum_{i=1}^n \sum_{j>i}^n \frac{y_i y_j}{\pi_i \pi_j} (\pi_i \pi_j - \pi_{ij}) = \\ &= \sum_{i=1}^n \frac{y_i^2}{\pi_i^2} \sum_{j \neq i}^n (\pi_i \pi_j - \pi_{ij}) + 2 \sum_{i=1}^n \sum_{j>i}^n \frac{y_i y_j}{\pi_i \pi_j} (\pi_{ij} - \pi_i \pi_j) = \\ \text{Como } \sum_{j \neq i}^n (\pi_i \pi_j - \pi_{ij}) &= \pi_i (1 - \pi_i) \\ &= \sum_{i=1}^n \frac{y_i^2}{\pi_i^2} \pi_i (1 - \pi_i) + 2 \sum_{i=1}^n \sum_{j>i}^n \frac{y_i y_j}{\pi_i \pi_j} (\pi_{ij} - \pi_i \pi_j) = V(\hat{\theta}_{HT}) \end{aligned}$$

### 3. MUESTREO CON REPOSICIÓN y PROBAB. DESIGUALES

Procedimiento de selección de la muestra

- con probab. desiguales
- las muestras con elementos repetidos son posibles,
- cada elemento poblacional puede estar 0, 1, ..., n veces en la muestra.

#### PROBABILIDADES

Para cada unidad poblacional  $u_i, i=1, \dots, N$ , definimos la v.a. auxiliar  $e_i \equiv n^\circ$  veces que aparece  $u_i$  en la muestra.

$e_i \rightarrow B(n, P_i)$ , donde  $P_i = P(u_i \text{ en muestra en cada extrac.})$

$$E[e_i] = nP_i$$

$$V[e_i] = nP_i(1-P_i)$$

La probab. de una muestra cualquiera seguirá el modelo multinomial, donde cada elemento  $u_i$  va a estar en la muestra  $e_i$  veces,  $\sum_{i=1}^N e_i = n$ .

$$(e_1, e_2, \dots, e_N) \rightarrow \text{Multinomial}(n, P_1, \dots, P_N) \quad \sum_{i=1}^N P_i = 1$$

$$P(e_1, \dots, e_N) = \frac{n!}{e_1! \dots e_N!} P_1^{e_1} \dots P_N^{e_N}, \quad e_i = 0, 1, \dots, N \quad \sum_{i=1}^N e_i = n$$

Para la ocurrencia conjunta, definimos la v.a. auxiliar  $e_i \cdot e_j = "n^\circ \text{ veces que aparecen } (u_i, u_j) \text{ en la muestra}"$ .

$$e_i \cdot e_j = 0, 1, \dots, n^2$$

A partir de la f. generatriz de momentos de la dist. multinomial se deducen la esperanza y la covarianza:

$$E[e_i \cdot e_j] = n(n-1) P_i P_j$$

$$\begin{aligned} \text{cov}[e_i, e_j] &= E[e_i \cdot e_j] - E[e_i] E[e_j] = n(n-1) P_i P_j - nP_i \cdot nP_j \\ &= n^2 P_i P_j - nP_i P_j - n^2 P_i P_j = -nP_i P_j \end{aligned}$$

Ya tenemos definido el vector esperanza y la matriz de var-cov:  $E[e_1, \dots, e_N] = (nP_1, \dots, nP_N)$

$$\Sigma(e_1, \dots, e_N) = \begin{pmatrix} nP_1(1-P_1) & -nP_1P_2 & \dots & -nP_1P_N \\ -nP_2P_1 & nP_2(1-P_2) & & \\ \vdots & & \ddots & \\ -nP_NP_1 & & & nP_N(1-P_N) \end{pmatrix}$$



ESTIMADOR lineal insesgado

Para el parámetro poblacional  $\Theta = \sum_{i=1}^N Y_i$ , utilizamos como estimador lineal  $\hat{\Theta} = \sum_{i=1}^N Y_i w_i$

Para que sea insesgado,  $E[\hat{\Theta}] = \Theta$

$$E[\hat{\Theta}] = E\left[\sum_{i=1}^N Y_i w_i\right] = E\left[\sum_{i=1}^N Y_i w_i e_i\right] = \sum_{i=1}^N Y_i w_i E[e_i]$$

en el caso de muestreo con reposicionamiento  $E[e_i] = nP_i$ , por lo que

$$E[\hat{\Theta}] = \sum_{i=1}^N Y_i w_i nP_i = \left(\sum_{i=1}^N Y_i\right) \stackrel{= \Theta}{=} \Leftrightarrow w_i = \frac{1}{nP_i}, \forall i = 1, \dots, N$$

Luego el estimador lineal insesgado para  $\Theta$  en el caso de muestreo con reposición es el estimador de Hausen y Hurwitz (1943):

$$\hat{\Theta}_{HH} = \sum_{i=1}^n \frac{Y_i}{nP_i}$$

cuya expresión se puede particularizar para los parámetros poblacionales más utilizados:

Parámetro	Estimador Hausen y Hurwitz, $\hat{\Theta}_{HH}$		
Total poblacional	$\Theta = \sum_{i=1}^N X_i = X$	$Y_i = X_i \rightarrow$	$\hat{X}_{HH} = \sum_{i=1}^n \frac{X_i}{nP_i}$
Media poblacional	$\Theta = \sum_{i=1}^N \frac{X_i}{N} = \bar{X}$	$Y_i = \frac{X_i}{N} \rightarrow$	$\hat{\bar{X}}_{HH} = \frac{1}{N} \sum_{i=1}^n \frac{X_i}{nP_i} = \frac{1}{N} \hat{X}$
Total de clase	$\Theta = \sum_{i=1}^N A_i = A$	$Y_i = A_i \rightarrow$	$\hat{A}_{HH} = \sum_{i=1}^n \frac{A_i}{nP_i}$
Proporción poblacional	$\Theta = \sum_{i=1}^N \frac{A_i}{N} = P$	$Y_i = \frac{A_i}{N} \rightarrow$	$\hat{P}_{HH} = \frac{1}{N} \sum_{i=1}^n \frac{A_i}{nP_i} = \frac{1}{N} \hat{A}_{HH}$

VARIANZA del estimador de Hansen y Hurwitz

$$\begin{aligned}
 V(\hat{\theta}_{HH}) &= V\left(\sum_{i=1}^N \frac{y_i}{nP_i}\right) = V\left(\sum_{i=1}^N \frac{y_i}{nP_i} e_i\right) = \sum_{i=1}^N \frac{y_i^2}{n^2 P_i^2} \underbrace{V(e_i)}_{nP_i(1-P_i)} + \\
 &+ \sum_{i \neq j} \frac{y_i}{nP_i} \cdot \frac{y_j}{nP_j} \underbrace{\text{cov}(e_i, e_j)}_{-nP_i P_j} = \\
 &= \sum_{i=1}^N \frac{y_i^2}{n^2 P_i^2} \cdot nP_i(1-P_i) + \sum_{i \neq j} \frac{y_i}{nP_i} \cdot \frac{y_j}{nP_j} (-nP_i P_j) = \\
 &= \sum_{i=1}^N \frac{y_i^2 (1-P_i)}{nP_i} - \sum_{i \neq j} \frac{y_i y_j}{n} = \sum_{i=1}^N \frac{y_i^2}{nP_i} - \frac{1}{n} \sum_{i=1}^N y_i^2 - \sum_{i \neq j} \frac{y_i y_j}{n} = \\
 &\text{Como } \left(\sum_{i=1}^N y_i\right)^2 = \sum_{i=1}^N y_i^2 + \sum_{i \neq j} y_i y_j \\
 &= \sum_{i=1}^N \frac{y_i^2}{nP_i} - \frac{1}{n} \left[ \sum_{i=1}^N y_i^2 + \sum_{i \neq j} y_i y_j \right] = \sum_{i=1}^N \frac{y_i^2}{nP_i} - \frac{1}{n} \underbrace{\left(\sum_{i=1}^N y_i\right)^2}_{\theta^2} = \\
 &= \sum_{i=1}^N \frac{y_i^2}{nP_i} - \frac{1}{n} \theta^2 = \frac{1}{n} \left( \sum_{i=1}^N \frac{y_i^2}{P_i} - \theta^2 \right) =
 \end{aligned}$$

$$V(\hat{\theta}_{HH}) = \frac{1}{n} \left[ \sum_{i=1}^N \left( \frac{y_i}{P_i} \right)^2 P_i - \theta^2 \right]$$

Podemos llegar a una formulación equivalente, ya que:

$$\begin{aligned}
 \sum_{i=1}^N \left( \frac{y_i}{P_i} - \theta \right)^2 P_i &= \sum_{i=1}^N \left( \frac{y_i}{P_i} \right)^2 P_i + \theta^2 \underbrace{\sum_{i=1}^N P_i}_1 - 2\theta \underbrace{\sum_{i=1}^N \frac{y_i}{P_i} P_i}_{\theta} = \\
 &= \sum_{i=1}^N \frac{y_i^2}{P_i} - \theta^2
 \end{aligned}$$

Otra expresión equivalente de la varianza del estimador de Hansen y Hurwitz es:

$$V(\hat{\theta}_{HH}) = \frac{1}{n} \sum_{i=1}^N \left( \frac{y_i}{P_i} - \theta \right)^2 P_i$$

Otra manera equivalente sería:

$$V(\hat{\theta}_{HH}) = \frac{1}{n} \sum_{i=1}^N \sum_{j>i}^N \left( \frac{y_i}{P_i} - \frac{y_j}{P_j} \right)^2 P_i P_j$$

# ESTIMACIÓN de la Varianza del estimador de Hausen y Hurwitz

Un estimador insesgado para  $V(\hat{\theta}_{HH})$  viene dado por:

$$\hat{V}(\hat{\theta}_{HH}) = \frac{1}{n(n-1)} \left[ \sum_{i=1}^n \left( \frac{y_i}{p_i} \right)^2 - n \hat{\theta}_{HH}^2 \right]$$

Es un estimador insesgado de  $V(\hat{\theta}_{HH})$ , porque  $E[\hat{V}(\hat{\theta}_{HH})] = V(\hat{\theta}_{HH})$

$$E[\hat{V}(\hat{\theta}_{HH})] = E \left[ \sum_{i=1}^n \left( \frac{y_i}{p_i} \right)^2 - n \hat{\theta}_{HH}^2 \right] =$$

$$= \frac{1}{n(n-1)} \left[ \sum_{i=1}^n \left( \frac{y_i^2}{p_i^2} \right) \cdot \underbrace{E[e_i]}_{n \hat{\theta}_{HH}} - n E[\hat{\theta}_{HH}^2] \right] =$$

$$= \frac{1}{n(n-1)} \left[ \sum_{i=1}^n \frac{y_i^2}{p_i^2} - \left[ V(\hat{\theta}_{HH}) + E[\hat{\theta}_{HH}]^2 \right] \right] =$$

$$= \frac{1}{n-1} \left[ \sum_{i=1}^n \frac{y_i^2}{p_i^2} - V(\hat{\theta}_{HH}) - \theta^2 \right] = \frac{1}{n-1} \left[ \underbrace{\left( \sum_{i=1}^n \frac{y_i^2}{p_i^2} - \theta^2 \right)}_{n \cdot V(\hat{\theta}_{HH})} - V(\hat{\theta}_{HH}) \right] =$$

$$= \frac{n-1}{n-1} V(\hat{\theta}_{HH}) = V(\hat{\theta}_{HH})$$

Una expresión equivalente es:

$$\hat{V}(\hat{\theta}_{HH}) = \frac{1}{n(n-1)} \sum_{i=1}^n \left( \frac{y_i}{p_i} - \hat{\theta}_{HH} \right)^2$$

ya que

$$\begin{aligned} \sum_{i=1}^n \left( \frac{y_i}{p_i} - \hat{\theta}_{HH} \right)^2 &= \sum_{i=1}^n \left( \frac{y_i^2}{p_i^2} + \hat{\theta}_{HH}^2 - 2 \hat{\theta}_{HH} \cdot \frac{y_i}{p_i} \right) = \\ &= \sum_{i=1}^n \left( \frac{y_i^2}{p_i^2} \right) + n \hat{\theta}_{HH}^2 - 2 \hat{\theta}_{HH} \underbrace{\sum_{i=1}^n \frac{y_i}{p_i}}_{n \cdot \hat{\theta}_{HH}} = \\ &= \sum_{i=1}^n \left( \frac{y_i^2}{p_i^2} \right) + n \hat{\theta}_{HH}^2 - 2 n \hat{\theta}_{HH}^2 = \\ &= \sum_{i=1}^n \left( \frac{y_i^2}{p_i^2} \right) - n \hat{\theta}_{HH}^2 \end{aligned}$$

MUEST\_T3

10



#### 4. PROBABILIDADES ÓPTIMAS de SELECCIÓN

En muchas ocasiones, es conveniente asignar a las unidades probabilidades de pertenecer a la muestra teniendo en cuenta el tamaño de la unidad. Es el caso de unidades de muestreo compuestas de unidades elementales.

Para cada unidad poblacional  $u_i$ ,  $i = 1 \dots N$ , defino  $M_i \equiv$  tamaño de la unidad  $\rightarrow$  n.º de unidades elementales que componen  $u_i$ .  $\sum_{i=1}^N M_i = M$

Muestreo SIN reposición;

$$\theta = \sum_{i=1}^N Y_i \rightarrow \hat{\theta}_{HT} = \sum_{i=1}^N \frac{y_i}{\pi_i} \quad \text{tg } \pi_i = P(u_i \in \text{muestra})$$

Si fuera  $\pi_i = n \cdot \frac{y_i}{\theta}$ , el valor del estimador sería

$$\hat{\theta}_{HT} = \sum_{i=1}^N y_i \cdot \frac{\theta}{n y_i} = \frac{n\theta}{n} = \theta \Rightarrow V(\hat{\theta}_{HT}) = 0$$

En este caso ideal ( $\theta$  descon.) el estimador proporcionaría siempre el valor del parámetro. Este resultado supone utilizar probabilidades  $\pi_i$  proporcionales a una variable conocida, en general de tamaño  $M_i$ , que se supone correlacionada con  $Y_i$ .

$$\pi_i = K \cdot M_i \Rightarrow n = \sum_{i=1}^N \pi_i = K \sum_{i=1}^N M_i = KM \Rightarrow K = \frac{n}{M}$$

$$\pi_i = \frac{n M_i}{M}$$

En la práctica, las unidades de muestreo suelen ser conglomerados de tamaño  $M_i$ , aunque a veces este modelo se utiliza con unidades simples, donde  $M_i$  con la distribución ponderaciones fue dada mayor o menor importancia a las unidades poblacionales.

Muestreo Con reposición:

$$\theta = \sum_{i=1}^N Y_i \longrightarrow \hat{\theta}_{HH} = \sum_{i=1}^n \frac{Y_i}{n P_i} \quad \text{tg} \quad P_i = P(u_i \in \text{muestra})$$

En el hipotético caso de que

$P_i = \frac{Y_i}{\theta}$ , el valor del estimador sería

$$\hat{\theta}_{HH} = \sum_{i=1}^n \frac{Y_i}{n \cdot \frac{Y_i}{\theta}} = \frac{n\theta}{n} = \theta \Rightarrow V(\hat{\theta}_{HH}) = 0.$$

Luego en este caso también es razonable asignar probabilidades de selección  $P_i$  proporcionales al tamaño de las unidades poblacionales  $M_i$ .

$$P_i = K \cdot M_i \Rightarrow 1 = \sum_{i=1}^N P_i = K \sum_{i=1}^N M_i = K M \Rightarrow K = \frac{1}{M}$$

## 5. ESTIMADORES ESPECIALES DE SELECCIÓN

$$P_i = \frac{M_i}{M}$$

REPOSICIÓN Y PROBAB. PROPORCIONALES AL TAMAÑO (pp1)

De una población  $\{u_1 \dots u_N\}$  de unidades complejas (o no equiprobables) con tamaño (o ponderaciones)

$$\{M_1 \dots M_N\} \quad \text{tg} \quad \sum_{i=1}^N M_i = M.$$

En el caso del muestreo SIN reposición tiene sentido que las probabilidades  $\pi_i$  de pertenecer a la muestra sean proporcionales a los tamaños  $M_i$ .

- Modelo de una generalidad
- Método de Brewer
- Método de Durbin
- Esquema mixto de Sánchez-Crespo y Gabeira

a) Esquema de urna generalizado

Cada unidad  $u_i$  aparece representada por  $M_i$  bolas del mismo color en la urna. Se realiza la extracción, sale bola de color  $i \Rightarrow u_i \in \text{muestra}$

$\Rightarrow$  se retiran todas las bolas de color  $i$

$\Rightarrow$  Quedan  $\pi - \pi_i$  bolas

$\hookrightarrow$  siguiente extracción.

$$n=1 \rightarrow \pi_i = P(u_i \in M) = \frac{M_i}{M}$$

$$\begin{aligned} n=2 \rightarrow \pi_i &= P(u_i \in \pi) = P(u_i \in \pi^{\text{extr}}) + P(u_i \notin \pi^{\text{extr}} \cap u_i \in \pi^{\text{ext}}) \\ &= P(u_i \in \pi^{\text{extr}}) + \sum_{j \neq i} P(u_j \in \pi^{\text{extr}}) \cdot P(u_i \in \pi^{\text{ext}} / u_j \in \pi^{\text{extr}}) \\ &= \frac{M_i}{M} + \sum_{j \neq i} \frac{M_j}{M} \cdot \frac{M_i}{M - M_j} = \frac{M_i}{M} \left( 1 + \underbrace{\sum_{j \neq i} \frac{M_j}{M - M_j}}_K \right) \end{aligned}$$

etc.  
 $\pi_{ij} = P(u_i \in \pi^{\text{ext}} \cap u_j \in \pi^{\text{ext}}) + P(u_j \in \pi^{\text{ext}} \cap u_i \in \pi^{\text{ext}})$  proporcional a  $M_i$

b) Esquema de Brewer (1963)

Propósito  
 Para  $n=2$ , la primera unidad se extrae sin reposición y con probab. proporcional a  $k_i = P_i \frac{(1-P_i)}{(1-2P_i)}$ ,

siendo  $P_i = \frac{M_i}{M} < 1/2$ . (por que  $k_i > 0$ )

La segunda extracción también es sin reposición y con probabilidad proporcional a  $\frac{P_i}{1-2P_i}$ .

$$\begin{aligned} \text{Sea } K &= \sum_{i=1}^N k_i = \sum_{i=1}^N P_i \frac{(1-P_i)}{(1-2P_i)} = \frac{1}{2} \left( \sum_{i=1}^N \frac{P_i}{1-2P_i} + 1 \right) \\ \pi_i &= P(u_i \in \text{muestra}) = P(u_i \in \pi^{\text{ext}}) + P(u_i \notin \pi^{\text{ext}} \cap u_i \in \pi^{\text{ext}}) = \\ &= P(u_i \in \pi^{\text{ext}}) + \sum_{j \neq i} P(u_j \in \pi^{\text{ext}}) \cdot P(u_i \in \pi^{\text{ext}} / u_j \in \pi^{\text{ext}}) = \frac{k_i}{K} + \sum_{j \neq i} \frac{k_j}{K} \cdot \frac{P_i}{1-P_j} = \\ &= \frac{k_i}{K} + \frac{P_i}{K} \sum_{j \neq i} \frac{P_j (1-P_j)}{(1-P_j)} = \frac{P_i}{K} \left[ \frac{P_i (1-P_i)}{1-2P_i} + \sum_{j \neq i} \frac{P_j}{1-2P_j} \right] = \\ &= \frac{P_i}{K} \left[ 1 + \frac{P_i}{1-2P_i} + \sum_{j \neq i} \frac{P_j}{1-2P_j} \right] = \frac{P_i}{K} \left[ 1 + \underbrace{\sum_{i=1}^N \frac{P_i}{1-2P_i}}_{2K} \right] = \frac{P_i}{K} \cdot 2K = 2P_i \end{aligned}$$

la probabilidad  $\pi_{ij}$  de que las unidades  $u_i$  y  $u_j$  pertenezcan a la muestra será:

$$\begin{aligned}\pi_{ij} &= P((u_i, u_j) \in \text{Muestra}) = P(u_i \in 1^s \cap u_j \in 2^s) + P(u_j \in 1^s \cap u_i \in 2^s) \\ &= P(u_i \in 1^s) \cdot P(u_j \in 2^s / u_i \in 1^s) + P(u_j \in 1^s) \cdot P(u_i \in 2^s / u_j \in 1^s) = \\ &= \frac{k_i}{K} \cdot \frac{P_j}{1-P_i} + \frac{k_j}{K} \cdot \frac{P_i}{1-P_j} = \frac{1}{K} \left[ \frac{P_i(1-P_i)}{1-2P_i} \cdot \frac{P_j}{1-P_i} + \frac{P_j(1-P_j)}{1-2P_j} \cdot \frac{P_i}{1-P_j} \right] \\ &= \frac{1}{K} \cdot P_i P_j \left[ \frac{1}{1-2P_i} + \frac{1}{1-2P_j} \right] \rightarrow \text{proporcional a } P_i. \\ &\quad \Rightarrow \text{proporcional a } \pi_i.\end{aligned}$$

Como es muestreo SIN reposicionamiento se utilizó el estimador de Horvitz y Thompson que para  $n=2$  tiene como expresión:

$$\hat{\theta}_{HT} = \frac{y_i}{\pi_i} + \frac{y_j}{\pi_j} = \frac{y_i}{2P_i} + \frac{y_j}{2P_j} = \frac{1}{2} \left( \frac{y_i}{P_i} + \frac{y_j}{P_j} \right)$$

$V(\hat{\theta}_{HT})$  se obtiene substituyendo  $\pi_i$  y  $\pi_{ij}$  en la fórmula general.

### Método de Durbin (1967)

Para  $n=2$ , la primera unidad se extrae sin reposición con probabilidad  $P_i = \frac{M_i}{M}$  y la segunda  $j$  con probabilidad proporcional a:

$$k_j = P_j \left( \frac{1}{1-2P_i} + \frac{1}{1-2P_j} \right) \quad / \quad P_j < 1/2 \quad (P_j = \frac{M_j}{M}).$$

Resulta inmediato comprobar que los valores  $\pi_i$  y  $\pi_{ij}$  coinciden con los obtenidos por el mt. de Brewer, luego tb. coinciden el estimador y su varianza.



Esquema mixto de Sánchez-Grespo y Gabeirán (1987)

Es un esquema de urnas con  $M_i$  bolas del color  $i$  representando a la unidad  $U_i$ .

En la 1ª extracción sale bola de color  $i$  ~~se retira solo~~ se retira solamente esa bola (quedan  $M_i - 1$  bolas color  $i$ )

$\Rightarrow$  las unidades tienen probabilidades gradualmente variables

$\Rightarrow U_i$  puede ser elegida  $0, 1, \dots, \min\{n, M_i\}$  veces

Es un método mixto porque  $\left\{ \begin{array}{l} \text{sin reposición (bola no se repone)} \\ \text{con reposición (} U_i \text{ puede)} \end{array} \right.$

Consideramos la v.a. auxiliar  $e_i \equiv u_i^2$  veces que aparece  $U_i$  en la muestra,

$$e_i \rightarrow H(M, n, P_i) \left\{ \begin{array}{l} E[e_i] = n P_i \\ V[e_i] = \frac{M-n}{M-1} n P_i (1-P_i) \\ \text{cov}(e_i, e_j) = - \frac{M-n}{M-1} n P_i P_j \end{array} \right.$$

$$P(e_1 \dots e_N) = \frac{\binom{M_1}{e_1} \dots \binom{M_N}{e_N}}{\binom{M}{n}} \quad / \quad \sum e_i = n$$

El estimador insesgado de  $\theta$  coincide con el de Hansen y Hurwitz, pero su varianza no:  $(P(U_i \in \text{Muestra}) = e_i \text{ cada extrac.})$

$$\textcircled{*} \hat{\theta}_{scg} = \sum \frac{y_i}{n P_i} = \hat{\theta}_{HH}$$

$$V(\hat{\theta}_{scg}) = \frac{M-n}{M-1} V(\hat{\theta}_{HH}) \quad , \quad \hat{V}(\hat{\theta}_{scg}) = \frac{M-n}{M} \hat{V}(\hat{\theta}_{HH})$$

$\Rightarrow \hat{\theta}_{scg}$  es más preciso que  $\hat{\theta}_{HH}$ .

Existe una versión mejorada de Gabeirán, en la que se retira de la urna  $b$  bolas en cada extracción,  $b = \min\{n, M_i\}$  de modo que todas las unidades siguen estando representadas en la siguiente extracción

$$\begin{aligned} P(U_i \in 1^a) &= \frac{M_i}{M} = P_i \\ P(U_i \in 2^a) &= P(U_i \notin 1^a \cap U_i \in 2^a) = \sum_{j \neq i} P(U_j \in 1^a) \cdot P(U_i \in 1^a / U_j \in 2^a) = \\ &= \frac{M_i}{M} \cdot \frac{M_i-1}{M-1} + \sum_{j \neq i} \frac{M_j}{M} \cdot \frac{M_i}{M-1} = \frac{M_i}{M} \left[ \frac{M_i-1}{M-1} + \sum_{j \neq i} \frac{M_j}{M-1} \right] = \frac{M_i}{M} \end{aligned}$$