

# MODELOS DE SUPERPOBLACION EN EL MUESTREO POR CUOTAS: ESTIMACION DE LA MEDIA POBLACIONAL

*Marta Guijarro Garv*

**RESUMEN.**— Se utiliza un modelo de analisis de la varianza como modelo de superpoblacion en la estimacion de la media poblacional mediante muestreo por cuotas, comprobandose que la expresion matricial del estimador basado en el modelo es la misma para los casos de uno y dos factores. En cada uno de estos contextos se estudia la insesgadez del estimador y se obtiene la expresion de su varianza.

## 1. INTRODUCCION

Basicamente, el muestreo por cuotas consiste en clasificar los  $N$  elementos de la poblacion finita objeto de estudio segun las distintas modalidades de uno o varios factores que representaremos por  $i, j, \dots, h$  ( $i = 1, \dots, I$ ;  $j = 1, \dots, J; \dots$ ;  $h = 1, \dots, H$ ). Ası,  $N_{ij\dots h}$  expresa el numero de unidades poblacionales pertenecientes a la casilla  $(i, j, \dots, h)$  de una tabla de contingencia multiples, siendo  $n_{ij\dots h}$  el numero de ellas que forman parte de la muestra de tamano  $n$ .

En este trabajo supondremos que el muestreo es por cuotas marginales, es decir, tomaremos una unica fraccion de muestreo,  $f$ , para cada conjunto de cuotas:

$$n_{i\dots} = fN_{i\dots}; n_{j\dots} = fN_{j\dots}; \dots; n_{\dots h} = fN_{\dots h}$$

indicando por la suma en todas las modalidades de la caracterıstica correspondiente.

La falta de aleatoriedad en la selección de la muestra y la consecuente imposibilidad de realizar inferencias nos llevan a la necesidad de utilizar enfoques basados en modelos, entre los que se encuentra el de los modelos de superpoblaciones que adoptaremos en este trabajo.

Admitiremos, por tanto, que la población finita ha sido generada como muestra aleatoria de una superpoblación infinita, esto es, que el valor de la variable de interés en cada individuo de la población es una realización de una variable aleatoria, formulando la distribución conjunta del total de variables mediante un modelo de análisis de la varianza.

En la sección 2 obtendremos, tanto la expresión matricial del estimador de la media poblacional basado en un modelo con un único factor, como la de su varianza. En la sección 3 realizaremos idéntico análisis para el caso de dos factores<sup>1</sup>.

## 2. MODELOS DE ANALISIS DE LA VARIANZA CON UN FACTOR

La población finita formada por  $N$  unidades está dividida en  $I$  cuotas, determinadas por las correspondientes modalidades de un único factor o categoría. Denotaremos por  $N_i$ , el número de elementos de  $i$ -ésima cuota ( $i = 1, \dots, I$ ).

El valor de la característica de interés en la unidad  $k$ -ésima de la  $i$ -ésima cuota, vendrá dado por  $y_{ik}$  con  $i = 1, \dots, I, k = 1, \dots, N_i$ .

Siguiendo el enfoque proporcionado por los modelos de superpoblación, el vector de valores poblacionales,  $y = (y_{11}, \dots, y_{1N_1}, \dots, y_{I1}, \dots, y_{IN_I})'$  es una realización del vector de variables aleatorias  $Y = (Y_{11}, \dots, Y_{1N_1}, \dots, Y_{I1}, \dots, Y_{IN_I})'$ , cuya distribución conjunta viene formulada a través de un modelo de análisis de la varianza,  $\xi$ , dado por:

$$Y = A\beta + \varepsilon \quad [1]$$

donde  $A$  es una matriz de dimensión  $N \times I$ ,

$$A = \begin{pmatrix} I_{N_1} & 0 & \dots & 0 \\ 0 & I_{N_2} & \dots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \dots & I_{N_I} \end{pmatrix}$$

con  $I_{N_i} = (1, \dots, 1)'$ ,  $N_i$ -vector unidad ( $i = 1, \dots, I$ ),  $\beta = (\beta_1, \dots, \beta_I)'$ , vector paramétrico de dimensión  $I \times 1$  y  $\varepsilon = (\varepsilon_{11}, \dots, \varepsilon_{1N_1}, \dots, \varepsilon_{I1}, \dots, \varepsilon_{IN_I})'$  vector

<sup>1</sup> Deville (1991) considera modelos de análisis de la varianza aunque con hipótesis ligeramente distintas a las formuladas en este trabajo. Sin embargo, dicho autor no utiliza un enfoque matricial, más adecuado al contexto del muestreo en superpoblaciones.

aleatorio de dimensión  $N \times I$  verificando:

$$\begin{aligned} E_{\xi}(\epsilon) &= 0 \\ E_{\xi}(\epsilon\epsilon') &= \Sigma \end{aligned}$$

donde  $E_{\xi}(\cdot)$  denota la esperanza respecto al modelo y  $\Sigma$  es una matriz de dimensión  $N \times N$  con la siguiente estructura:

$$\Sigma = \begin{pmatrix} \Sigma_1 & 0 & \dots & 0 \\ 0 & \Sigma_2 & \dots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \dots & \Sigma_I \end{pmatrix}$$

siendo  $\Sigma_i = \text{diag}(\sigma_i^2, \dots, \sigma_i^2) N_i \times N_i$ -matriz diagonal ( $i = 1, \dots, I$ ).

Sin pérdida de generalidad, listaremos las unidades de la muestra en primer lugar, teniendo en cuenta que ésta incluye  $n_i$  individuos de la categoría  $i$  ( $i = 1, \dots, I$ ). Tendremos, así, las siguientes particiones de  $Y, A, \epsilon$  y  $I$ , vector de unos de dimensión  $N \times I$ :

$$\begin{aligned} Y &= (Y'_s, Y'_r)' \\ A &= \begin{pmatrix} A_s \\ A_r \end{pmatrix} \\ \epsilon &= (\epsilon'_s, \epsilon'_r)' \\ I &= (I'_s, I'_r)' \end{aligned}$$

donde

$$\begin{aligned} Y_s &= (Y_{11}, \dots, Y_{1n_1}, \dots, Y_{In_1})' \\ Y_r &= (Y_{11}, \dots, Y_{1N_1-n_1}, \dots, Y_{I1}, \dots, Y_{IN_I-n_I})' \\ A_s &= \begin{pmatrix} 1_{n_1} & 0 & \dots & 0 \\ 0 & 1_{n_2} & \dots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \dots & 1_{n_I} \end{pmatrix} \\ A_r &= \begin{pmatrix} 1_{N_1-n_1} & 0 & \dots & 0 \\ 0 & 1_{N_2-n_2} & \dots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \dots & 1_{N_I-n_I} \end{pmatrix} \\ \epsilon_s &= (\epsilon_{11}, \dots, \epsilon_{1n_1}, \dots, \epsilon_{I1}, \dots, \epsilon_{In_1})' \\ I_s &= (I'_{n_1}, \dots, I'_{n_I})' \end{aligned}$$

Siguiendo a Royall (1970), el estimador de la media poblacional basado en el modelo, resulta de dividir por  $N$  la suma de los valores observados y la predicción de los no observados. En este caso:

$$e_R = N^{-1} (I'_s Y_s + I'_r \hat{Y}_r) = N^{-1} (I'_s Y_s + I'_r A_r \hat{\beta})$$

Dado que el estimador de Gauss-Markov de  $\beta_i$  es la media de los valores observados en la categoría  $i$ -ésima, es decir,

$$\hat{\beta} = \bar{y}_i = \left( \sum_{k=1}^{n_i} y_{ik} \right) / n_i \quad i = 1, \dots, \dots, I$$

es estimador del vector paramétrico es  $\hat{\beta} = (\bar{y}_1, \dots, \bar{y}_I)'$ .

Esto nos lleva a

$$I'_s Y_s = I'_s A_s \hat{\beta}$$

con lo que, en este caso, el estimador predictivo de la media poblacional es

$$e_R = N^{-1} (I'_s A_s \hat{\beta} + I'_r A_r \hat{\beta}) = N^{-1} (I'_s A_s + I'_r A_r) \hat{\beta} = N^{-1} I' A \hat{\beta}$$

## 2.1. ESPERANZA Y VARIANZA DEL ESTIMULADOR PREDICTIVO

Veamos que el sesgo del estimador es cero. En efecto:

$$N^{-1} E_{\xi} (I' A \hat{\beta} - I' Y) = N^{-1} E_{\xi} (I' A \hat{\beta} - I' A \hat{\beta} - I' \epsilon) = N^{-1} [I' A E_{\xi} (\hat{\beta} - \beta) - I' E_{\xi} (\epsilon)] = 0$$

puesto que

$$E_{\xi} (\hat{\beta}_i - \beta_i) = E_{\xi} (\bar{e}_i) = 0 \quad i = 1, \dots, N$$

con

$$\bar{e}_i = \frac{1}{n_i} \sum_{k=1}^{n_i} e_{ik}$$

es decir,  $\hat{\beta}$  es un estimador insesgado de  $\beta$ .

Por tanto, y considerando que

$$Var_{\xi} \left( \frac{1}{N} I' A \beta \right) = \frac{1}{N^2} I' A V ar_{\xi}(\beta) A' I$$

bastará calcular  $Var_{\xi}(\hat{\beta})$  para obtener la varianza del estimador:

$$Var_{\xi}(\hat{\beta}) = E_{\xi}[(\hat{\beta} - \beta)(\hat{\beta} - \beta)'] = \begin{pmatrix} \sigma_1^2/n_1 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & \sigma_I^2/n_I \end{pmatrix}$$

ya que

$$E_{\xi}(\hat{\beta}_i - \beta_i)^2 = E_{\xi}(\bar{\epsilon}_i)^2 = \sigma_i^2/n_i \quad i = 1, \dots, I$$

Si denotamos por

$$V = \begin{pmatrix} \sigma_1^2/n_1 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & \sigma_I^2/n_I \end{pmatrix}$$

tendremos:

$$Var_{\xi}(e_R) = N^{-2} I' A V A' I$$

Esta varianza puede ser estimada sustituyendo  $\sigma_i^2$  por

$$\hat{\sigma}_i^2 = \frac{1}{n_i - 1} \sum_{k=1}^{n_i} (Y_{ik} - \bar{y}_i)^2$$

### 3. MODELO DE ANALISIS DE LA VARIANZA CON DOS FACTORES

Supongamos, ahora, la influencia de dos factores sobre los valores poblacionales de la variable de interés. Denotaremos por  $N_{ij}$  ( $i = 1, \dots, I$ ;  $j = 1, \dots, J$ ), el número de individuos de la población que poseen las modalidades  $i$  y  $j$ , es decir, que pertenecen a la casilla,  $(i, j)$  de la correspondiente tabla de contingencia múltiple. Así, el valor de la característica objeto de estudio, en la unidad  $k$ -ésima de la casilla  $(i, j)$  se representará por  $y_{ijk}$  ( $k = 1, \dots, N_{ij}$ ).

El vector de valores poblacionales  $y = (y_{111}, \dots, y_{11N_{11}}, \dots, y_{1J1}, \dots, y_{1JN_{1J}})'$ , es, en este caso, una realización del vector de variables aleatorias  $Y = (Y_{111}, \dots, Y_{11N_{11}}, \dots, Y_{1J1}, \dots, Y_{1JN_{1J}})'$  siendo el modelo de superpoblación el dado por la expresión [1], donde  $A$  es ahora de dimensión  $N \times (I + J + 1)$ :

$$A = \begin{pmatrix} I_{N_{11}} & I_{N_{11}} & \dots & 0 & I_{N_{11}} & \dots & 0 \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ I_{N_{1J}} & I_{N_{1J}} & \dots & 0 & 0 & \dots & I_{N_{1J}} \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ I_{N_{I1}} & 0 & \dots & I_{N_{I1}} & I_{N_{I1}} & \dots & 0 \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ I_{N_{IJ}} & 0 & \dots & I_{N_{IJ}} & 0 & \dots & I_{N_{IJ}} \end{pmatrix}$$

con  $IN_{ij} = (1, \dots, 1)'$ ,  $N_{ij}$ -vector unidad ( $i = 1, \dots, I$ ;  $j = 1, \dots, J$ ),  $\beta = (\mu, \alpha_1, \dots, \alpha_J, \gamma_1, \dots, \gamma_J)$ , vector paramétrico de dimensión  $(I + J + 1) \times 1$  y  $\varepsilon = (\varepsilon_{111}, \dots, \varepsilon_{11n_{11}}, \dots, \varepsilon_{IJ}, \dots, \varepsilon_{IJn_{IJ}})'$ , vector aleatorio de dimensión  $N \times 1$ , cuyas componentes son centradas, independientes y, además,

$$\text{Var}_{\xi}(\varepsilon_{ijk}) = \tau_i^2 + w_j^2 \quad i = 1, \dots, I, \quad j = 1, \dots, J, \quad k = 1, \dots, N_{ij}$$

Consideraremos las mismas particiones de  $Y$ ,  $A$ ,  $\varepsilon$  y  $I$  empleadas en el caso de un único factor. En esta nueva situación tendremos:

$$Y_s = (Y_{111}, \dots, Y_{11n_{11}}, \dots, Y_{IJ}, \dots, Y_{IJn_{IJ}})'$$

$$Y_r = (Y_{111}, \dots, Y_{11N_{11}-n_{11}}, \dots, Y_{IJ}, \dots, Y_{IJN_{IJ}-n_{IJ}})'$$

$$A_s = \begin{pmatrix} I_{n_{11}} & I_{n_{11}} & \dots & 0 & I_{n_{11}} & \dots & 0 \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ I_{n_{1J}} & I_{n_{1J}} & \dots & 0 & 0 & \dots & I_{n_{1J}} \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ I_{n_{I1}} & 0 & \dots & I_{n_{I1}} & I_{n_{I1}} & \dots & 0 \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ I_{n_{IJ}} & 0 & \dots & I_{n_{IJ}} & 0 & \dots & I_{n_{IJ}} \end{pmatrix}$$

$$A_r = \begin{pmatrix} I_{N_{11}-n_{11}} & I_{N_{11}-n_{11}} & \dots & 0 & I_{N_{11}-n_{11}} & \dots & 0 \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ I_{N_{1J}-n_{1J}} & I_{N_{1J}-n_{1J}} & \dots & 0 & 0 & \dots & I_{N_{1J}-n_{1J}} \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ I_{N_{I1}-n_{I1}} & 0 & \dots & I_{N_{I1}-n_{I1}} & I_{N_{I1}-n_{I1}} & \dots & 0 \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ I_{N_{IJ}-n_{IJ}} & 0 & \dots & I_{N_{IJ}-n_{IJ}} & 0 & \dots & I_{N_{IJ}-n_{IJ}} \end{pmatrix}$$

$$\varepsilon_s = (\varepsilon_{111}, \dots, \varepsilon_{11n_{11}}, \dots, \varepsilon_{IJ}, \dots, \varepsilon_{IJn_{IJ}})'$$

$$I_s = (I'_{n_{11}}, \dots, I'_{n_{IJ}})'$$

con  $I_{n_{ij}}$  y  $I_{N_{ij}-n_{ij}}$  vectores unidad de dimensiones  $n_{ij}$  y  $N_{ij} - n_{ij}$ , respectivamente.

Admitiremos que se cumplen las siguientes condiciones de identificabilidad del modelo:

$$\sum_j n_{ij} \alpha_i = 0 \text{ y } \sum_j n_{ij} \beta_j = 0$$

Utilizaremos el muestreo por cuotas marginales, con muestras de tamaño fijo  $n = fN$ , incluyendo  $n_i = fN_i$  individuos por cada modalidad  $i$  del primer factor y  $n_j = fN_j$  por cada  $j$  del segundo.

Al igual que en la situación presentada en la sección anterior, nuestro objetivo es calcular el estimador de Royall, basándonos en la predicción de los valores observados. Es decir,

$$e_R = \frac{1}{N} (I'_s Y_s + I'_r \hat{Y}_r) = \frac{1}{N} (I'_s Y_s + I'_r A_r \hat{\beta})$$

La estimación del vector paramétrico por mínimos cuadrados ordinarios nos lleva a:

$$\hat{\beta} = (\bar{y}, \bar{y}_1 - \bar{y}, \dots, \bar{y}_I - \bar{y}, \bar{y}_1 - \bar{y}, \dots, \bar{y}_J - \bar{y})'$$

donde  $\bar{y}$  es la media de los valores muestrales;  $\bar{y}_i$  e  $\bar{y}_j$  son las medias de las observaciones muestrales de las modalidades  $i$ -ésima y  $j$ -ésima del primer y segundo factor, respectivamente ( $i = 1, \dots, I; j = 1, \dots, J$ ).

Podemos obtener una expresión más sencilla del estimador predictivo teniendo en cuenta que,

$$I'_s A_s \hat{\beta} = n\bar{y} + \sum_i \sum_j n_{ij} (\bar{y}_i - \bar{y}) + \sum_i \sum_j (\bar{y}_j - \bar{y}) = n\bar{y} = I'_s Y_s$$

Por tanto,

$$e_R = \frac{1}{N} (I'_s A_s \hat{\beta} + I'_r A_r \hat{\beta}) = \frac{1}{N} I'_r A_r \hat{\beta} \quad 2$$

2 Operando obtenemos

$$I'_r A_r \hat{\beta} = (N - n) \bar{y}$$

con lo que el estimador de Royall coincide, en este caso, con  $\bar{y}$ .

### 3.1. ESPERANZA Y VARIANZA DEL ESTIMADOR PREDICTIVO

Siguiendo el mismo razonamiento que en el apartado 2.1, el estimador  $e_R$  es insesgado de la media poblacional, puesto que la independencia entre ambos factores nos lleva a que  $\hat{\beta}$  es un estimador insesgado de  $\beta$ .

Con el fin de hallar la varianza del estimador, denotemos por  $m$  el vector de dimensión  $N \times I$  de componentes  $m_{ijk} = E_{\xi}(Y_{ijk})$ . A partir de él obtenemos, de manera habitual, el correspondiente vector  $m_s$  para los valores de la muestra, pudiendo, así, escribir:

$$\begin{aligned} \frac{1}{N^2} E_{\xi}(I'A\hat{\beta} - I'Y)^2 &= \frac{1}{N^2} E_{\xi} \left[ \frac{N}{n} 1'_s(Y_s - m_s) - I'(Y - m) \right]^2 = \\ \frac{1}{N^2} E_{\xi} \left( \frac{N}{n} I'_s \epsilon_s - I'\epsilon \right)^2 &= \frac{1}{N^2} \left[ E_{\xi} \left( \frac{N}{n} I'_s \epsilon_s \right)^2 + E_{\xi}(I'\epsilon)^2 - 2E_{\xi} \left( \frac{N}{n} I'_s \epsilon_s I'\epsilon \right) \right] = \\ \frac{1}{n^2} \sum_i \sum_j n_{ij} (\tau_i^2 + \omega_j^2) &+ \frac{1}{N^2} \sum_i \sum_j N_{ij} (\tau_i^2 + \omega_j^2) - \frac{2}{Nn} \sum_i \sum_j n_{ij} (\tau_i^2 + \omega_j^2) \end{aligned}$$

Ahora bien, aplicando las condiciones del muestreo por cuotas marginales, deducimos que:

$$\sum_i \sum_j n_{ij} (\tau_i^2 + \omega_j^2) = \sum_i n_i \tau_i^2 + \sum_j n_j \omega_j^2 = \frac{n}{N} \left( \sum_i N_i \tau_i^2 + \sum_j N_j \omega_j^2 \right)$$

con lo que la varianza del estimador predictivo puede expresarse como:

$$\frac{1}{N^2} E_{\xi}(I'A\hat{\beta} - I'Y)^2 = \frac{1}{Nn} (1-f) \left( \sum_i N_i \tau_i^2 + \sum_j N_j \omega_j^2 \right)$$

Obsérvese que esta varianza depende sólo de las cuotas y no de la muestra; este hecho justifica, en cierta medida, el empleo de cuotas marginales.

### BIBLIOGRAFIA

- Deville, J. (1991): A theory of quota surveys. *Survey Methodology* 17 n. 2, 163-181.  
 Gourieroux, C. (1981): Theory des sondages. *Economica*.  
 Herson, J. y Royall, R. M. (1973): Robust estimation in finite populations. *Journal of American Statistical Association* 68, 880-893.  
 Kendall, M. G. y Stuard, A. (1967): *The advanced theory of statistics*. New York: Hafner.

- Rohatgi, V. K. (1976): *An introduction to probability theory and mathematical statistics*. New York: John Wiley.
- Royall, R. M. (1970): On finite population sampling theory under certain linear regression models. *Bioetrika* 57, 377-387.
- Tam, S. M. (1988b): Some results son robust estimation in finite population sampling. *Journal of American Statistical Association* 83, 242-248.
- Wolter, K. M. (1985): *Introduction to variance estimation*. New York: Springer-Verlag.