

# MUEST-T16. ERRORES AJENOS AL MUESTREO I:

MARCOS IMPERFECTOS.

EL PROBLEMA de las UNIDADES VACÍAS.

ESTIMACIÓN del TOTAL y de la MEDIA.

CÁLCULO de la VARIANZA y COMPARACIÓN con la VARIANZA del MARCO DEPURADO.

---

## 1. ERRORES AJENOS AL MUESTREO

En las encuestas por muestreo puede definirse el "error" de una determinada estimación como la diferencia entre el valor observado  $\hat{\theta}$  y el valor desconocido de la característica poblacional  $\theta$  que tratamos de estimar.

$$\text{error} = |\hat{\theta} - \theta|$$

Los errores se deben a causas diversas, pudiendo clasificarse en errores de carácter aleatorio y errores de carácter sistemático o sesgos.

Pueden originarse errores en los resultados de una muestra particular debido a los respondientes, a los encuestadores, codificadores, etc, así como a la posible interdependencia entre ellos.

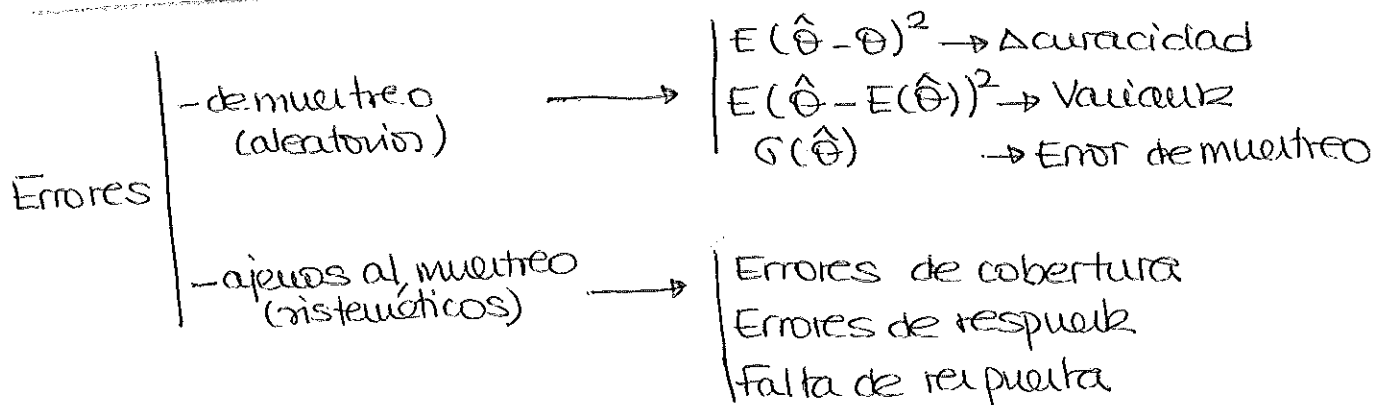
Los errores de carácter aleatorio y los errores de carácter sistemático tienen, en general, distintas fuentes, efectos y métodos de medida. La reducción de los errores aleatorios requiere hacer "más de algo" (aumentar el tamaño muestral), mientras que la reducción de los errores sistemáticos requiere hacer "algo más" (supervisar, controlar...)

Otra clasificación distingue entre errores de muestreo (originados por la variabilidad de los valores obtenidos en el muestreo, de carácter aleatorio) y los errores ajenos al muestreo (que se producen por causas no proba-

bilísticos, con errores no aleatorios)

mientras los errores de muestreo decrecen al aumentar el tamaño de la muestra, los errores ajenos al muestreo suelen crecer al aumentar el tamaño de la investigación, o en casos no decrecen.

los errores de muestreo se estiman con los datos de la muestra, los errores ajenos al muestreo suelen requerir para su estimación datos extramuestrales.



• Errores de cobertura  $\equiv$  marco imperfecto

El marco no ~~es~~ se corresponde con la poblac. objetivo y subestima o sobreestima la característica.

• Errores de respuesta  $\equiv$  debidos al proceso de la encuesta <sup>recopile de la info y tróbul.</sup>  
 distintas respuestas, posibles influencias de entrevistadores, codificadores, etc.

• Falta de respuesta  $\equiv$  imposibilidad de obtener toda la información por falta de unidades respondientes, cuestionarios incompletos...

## ARTICULO 50

1. LAS INFRACCIONES SE CLASIFICAN EN MUY GRAVES, GRAVES Y LEVES.

2. SON INFRACCIONES MUY GRAVES:

A) EL INCUMPLIMIENTO DEL DEBER DEL SECRETO ESTADISTICO.

B) LA UTILIZACION PARA FINALIDADES DISTINTAS DE LAS ESTADISTICAS DE LOS DATOS PERSONALES OBTENIDOS DIRECTAMENTE DE LOS INFORMANTES POR LOS SERVICIOS ESTADISTICOS.

C) EL SUMINISTRO DE DATOS FALSOS A LOS SERVICIOS ESTADISTICOS COMPETENTES.

D) LA RESISTENCIA NOTORIA, HABITUAL O CON ALEGACION DE EXCUSAS FALSAS EN EL ENVIO DE LOS DATOS REQUERIDOS, CUANDO HUBIERE OBLIGACION DE SUMINISTRARLOS.

E) LA COMISION DE UNA INFRACCION GRAVE CUANDO EL INFRACITOR HUBIERE SIDO SANCIONADO POR OTRAS DOS GRAVES DENTRO DEL PERIODO DE UN AÑO.

3 graves en 1 año

3. SON INFRACCIONES GRAVES:

A) LA NO REMISION O EL RETRASO EN EL ENVIO DE LOS DATOS REQUERIDOS CUANDO SE PRODUCIESE GRAVE PERJUICIO PARA EL SERVICIO, Y HUBIERE OBLIGACION DE SUMINISTRARLOS.

B) EL ENVIO DE DATOS INCOMPLETOS O INEXACTOS CUANDO SE PRODUCIESE GRAVE PERJUICIO PARA EL SERVICIO, Y HUBIERE OBLIGACION DE SUMINISTRARLOS.

3 leves en 1 año

C) LA COMISION DE UNA INFRACCION LEVE CUANDO EL INFRACITOR HUBIERA SIDO SANCIONADO POR OTRAS DOS LEVES DENTRO DEL PERIODO DE UN AÑO.

4. SON INFRACCIONES LEVES:

A) LA REMISION O EL RETRASO EN EL ENVIO DE DATOS CUANDO NO HUBIERE CAUSADO PERJUICIO GRAVE PARA EL SERVICIO, Y HUBIERE OBLIGACION DE SUMINISTRARLOS.

B) EL ENVIO DE DATOS INCOMPLETOS O INEXACTOS CUANDO NO HUBIERE CAUSADO PERJUICIO GRAVE PARA EL SERVICIO, Y HUBIERE OBLIGACION DE SUMINISTRARLOS.

ARTICULO 51 SANCIONES: de 10.000 a 5.000.000 Ptas

1. LAS INFRACCIONES MUY GRAVES SERAN SANCIONADAS CON MULTAS DE 500.001 A 5.000.000 DE PESETAS.

2. LAS INFRACCIONES GRAVES SERAN SANCIONADAS CON MULTAS DE 50.001 A 500.000 PESETAS.

10.000 - 50.000  
50.001 - 500.000  
500.001 - 5.000.000

3. LAS INFRACCIONES LEVES SE SANCIONARAN CON MULTAS DE 10.000 A 50.000 PESETAS.

4. LA CUANTIA DE LAS SANCIONES ESTABLECIDAS EN LOS APARTADOS ANTERIORES SE GRADUARA ATENDIENDO, EN CADA CASO, A LA PROPIA GRAVEDAD DE LA INFRACCION, A LA NATURALEZA DE LOS DAÑOS Y PERJUICIOS CAUSADOS Y A LA CONDUCTA ANTERIOR DE LOS INFRACTORES.

ARTICULO 52

## 2 - FALTA de RESPUESTA y SUS EFECTOS

Prácticamente en todos los casos (muestreo exhaustivo) ~~como~~<sup>y</sup> encuestas (por muestreo), en los que se utilizan cuestionarios y entrevistadores para la recogida de datos se produce el problema de la falta de respuesta  $\Rightarrow$  falta información de una parte de la población o de <sup>una parte de</sup> la muestra seleccionada.

Falta de respuesta total  $\rightarrow$  una unidad respondiente no ha contestado ninguna pregunta del cuest.

Falta de respuesta parcial  $\rightarrow$  la unidad ha dejado alguna pregunta sin responder.

Las posibles causas del problema de la no respuesta son:

(X) MIRAR ATRAS

1 - Imposibilidad de identificar la unidad sobre el terreno o de acceder a la misma (zonas difíciles, inconveniencias tiempo...)

2 - Ausencia temporal del entrevistado

3 - Incapacidad para contestar por parte del entrevistado

$\hookrightarrow$  Contribuye el hecho de que el respondiente deba ser una persona específica, al no ser válida la respuesta de otra persona del hogar, p. ej.

4 - Negativa a cooperar en la encuesta por parte del entrevistado

$\hookrightarrow$  Por razones personales, subjetivas

$\hookrightarrow$  Si las preguntas son incómodas (íntimas, económicas)

$\hookrightarrow$  Por temas de cumplimiento y envío en cuest. por correo

$\rightarrow$  Destacar la legislación estadística, que obliga a cumplimiento de cuestionarios del INE. (ATRAS)

5 - Pérdida de información

$\hookrightarrow$  En el proceso de recogida, codificación...

$\hookrightarrow$  Falta de conocimiento del entrevistador

adecuado, motivación...

EFFECTOS:

Población:  $N \begin{cases} N_1 & \text{, que contestan} \\ N_2 & \text{, que no contestan} \end{cases} \rightarrow W = \frac{N_2}{N}$  } *efectos*

$W \equiv$  proporción de los no respondientes

los parámetros poblacionales quedan:

$$\bar{X} = \sum_{i=1}^N X_i = (1-W)\bar{X}_1 + W\bar{X}_2$$

$$\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i = (1-W)\bar{X}_1 + W\bar{X}_2$$

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})^2 = \underbrace{\sum_{h=1}^2 \frac{N_h}{N} \sigma_h^2}_{\text{desv. intr.}} + \underbrace{\sum_{h=1}^2 \frac{N_h}{N} (\bar{X}_h - \bar{X})^2}_{\text{cada entr. con total}}$$

La falta de respuesta produce:

1 - Disminución del tamaño de la muestra  $\Rightarrow$  disminuye la precisión  $N \rightarrow n_1$

2 - Aparece un sesgo independiente del tamaño muestral

El primer problema se puede solucionar eligiendo un tamaño muestral mfc. grande, para que  $n_1$  tb. lo sea.

El segundo problema si fue es importante, porque:

$$\bar{X} \rightarrow \bar{X}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} X_i \quad \text{con } E[\bar{X}_1] = \bar{X}_1 \Rightarrow \text{sesgado para } \bar{X}$$

cuyo sesgo es

$$B(\bar{X}_1) = E(\bar{X}_1) - \bar{X} = \bar{X}_1 - \bar{X} = \bar{X}_1 - (1-W)\bar{X}_1 - W\bar{X}_2 = W(\bar{X}_1 - \bar{X}_2)$$

$\Rightarrow$  el sesgo es proporcional al  $n^2$  de unid. no respondientes.

$$\bar{X} / B(\bar{X}) \left\{ \begin{array}{l} \text{indep. de } n \\ \text{proporcional a } N_2 \end{array} \right.$$

En el caso del total poblacional, si consideramos

$\hat{X}_1 = N\bar{x}_1$ , también es sesgado pues  $E[\hat{X}_1] = N\bar{X}_1$

$$B(\hat{X}_1) = NW(\bar{X}_1 - \bar{X}_2)$$

$\hat{X}_2 = \frac{N}{n} \sum_{i=1}^{n_1} X_i$  tb. es sesgado.

La cuasivariante muestral tb. es sesgado para la variante poblacional.

$$\hat{X}_1 = N\bar{x}_1 = \frac{N}{n_1} \sum_{i=1}^{n_1} X_i \quad \text{con} \quad E[\hat{X}_1] = N \cdot \bar{X}_1$$

$$B(\hat{X}_1) = N(\bar{X}_1 - \bar{X}_2)$$

$$\hat{X}_2 = \frac{N}{n} \sum_{i=1}^{n_1} X_i \quad \text{con} \quad E[\hat{X}_2] = \frac{N}{n} \cdot E\left[\sum_{i=1}^{n_1} X_i\right] = \frac{N}{n} n_1 \bar{X}_1$$

$$B(\hat{X}_2) = E[\hat{X}_2] - X$$

$$s_1^2 = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (X_i - \bar{x}_1)^2 \quad \text{con} \quad E[s_1^2] = \sigma_1^2 \neq \sigma^2$$

$$B(s_1^2) = \sigma_1^2 - \sigma^2 = W(\sigma_1^2 - \sigma_2^2)$$

Parámetro	Estimador	Sesgo
$X = X_1 + X_2$	$\hat{X}_1 = N \cdot \bar{x}_1$	$N \cdot W_2 (\bar{X}_1 - \bar{X}_2)$
	$\hat{X}_2 = N \cdot \frac{\sum X_i}{n}$	$\frac{N}{n} \cdot n_1 \bar{X}_1 - X$
$\bar{X} = W_1 \bar{X}_1 + W_2 \bar{X}_2$	$\bar{x}_1 = \frac{\sum X_i}{n_1}$	$W_2 (\bar{X}_1 - \bar{X}_2)$
$\sigma^2 = \sum_{k=1}^2 \frac{N_k}{N} \sigma_k^2 + \sum_{k=1}^2 \frac{N_k}{N} (\bar{X}_k - \bar{X})^2$	$\hat{s}_1^2 = \frac{1}{n_1 - 1} \sum (X_i - \bar{x}_1)^2$	$W_2 (\sigma_1^2 - \sigma_2^2)$

### Modelos de Deming (1952)

Procedimiento para disminuir la ~~fa~~ no respuesta, especialmente producida por ausencia de la persona específica seleccionada para la encuesta, realizando visitas sucesivas.

Para ello se fija un n° mínimo de "revisitas" que deben hacerse a cada unidad antes de abandonarlo como "contacto imposible".

Este procedimiento encarece el coste de la encuesta y puede retrasar considerablemente los resultados finales  $\rightarrow$  planificar previamente la duración de la recogida de datos.

Deming desarrolló el siguiente modelo:

se divide la población en  $L$  clases de acuerdo con la probabilidad de que el encuestado se encuentre en casa.

$w_{ij}$  = probab. de que un encuestado de la clase  $j$  sea entrevistado en la visita  $i$

$P_j$  = proporción de la población de la clase  $j$ .

$\left. \begin{matrix} \bar{X}_j \\ S_j^2 \end{matrix} \right\}$  media y varianzas pobl. de la clase  $j$ .

$\text{Sp } w_{ij} > 0 \quad \forall j = 1, \dots, L$

Después de  $i$  visitas la composición de la muestra está formada por  $L$  clases, con elementos del estrato  $j$  que han respondido en los  $i$  primeros intentos, más una clase que no ha respondido en estos  $i$  primeros intentos.

Si  $n_0$  es el tamaño inicial de la muestra,  $n_{ij}$  es una multinomial.

Después de  $i$  visitas, el n° de respuestas  $n_i$  es una binomial:

$$n_i \rightarrow B\left(n_0, \sum_{j=1}^L w_{ij} P_j\right)$$

$$E[n_i] = n_0 \cdot \sum_{j=1}^L w_{ij} P_j$$



El total muestral se puede escribir como

$$x = x_1 + x_2$$

donde  $x_2 \equiv$  unidades sin respuesta en la primera vuelta.  
Si llamamos  $f_{21} \equiv$  fracción de muestreo en el estrato 2 en  
segunda vuelta,  $f_{21} = \frac{n_{21}}{n_2}$ , un estimador insesgado de  $x_2$

$$\text{es: } \hat{\hat{x}}_2 = \frac{n_2}{n_{21}} x_{21} = n_2 \cdot \bar{x}_{21} \quad \left( \bar{x}_{21} = \sum \frac{x_{i21}}{n_{21}} \right).$$

$$\hat{X} = \hat{x}_1 + \hat{\hat{x}}_2 = N \left( \frac{1}{n} x_1 + \frac{1}{n} \hat{\hat{x}}_2 \right) = \frac{N}{n} (x_1 + \hat{\hat{x}}_2)$$

$$E[\hat{X}] = E \left[ \underbrace{E_{n_1 n_2}}_{\text{fijos}} (\hat{X}) \right] = E \left[ \frac{N}{n} (x_1 + n_2 \underbrace{E_{n_1 n_2}(\bar{x}_{21})}_{\bar{x}_2}) \right] = E \left( \frac{N}{n} x \right) = E[N\bar{x}] = X$$

luego  $\hat{X}$  es insesgado para  $X$ .

$$V(\hat{X}) = E V_{n_1 n_2}(\hat{X}) + V E_{n_1 n_2}(\hat{X})$$

$$\begin{aligned} \text{donde } V_{n_1 n_2}(\hat{X}) &= V_{n_1 n_2} \left[ \frac{N}{n} (x_1 + \hat{\hat{x}}_2) \right] = \left( \frac{N}{n} \right)^2 V_{n_1 n_2}(\hat{\hat{x}}_2) = \left( \frac{N}{n} \right)^2 n_2^2 V_{n_1 n_2}(\bar{x}_{21}) \\ &= \frac{1}{f^2} \cdot n_2^2 \cdot \frac{n_2 - n_{21}}{n_2} \cdot \frac{S_{21}^2}{n_{21}} = \frac{1 - f_{21}}{f^2 \cdot f_{21}} \cdot n_2 S_{21}^2 \end{aligned}$$

$$\begin{aligned} E V_{n_1 n_2}(\hat{X}) &= \frac{1 - f_{21}}{f^2 \cdot f_{21}} E[n_2 S_{21}^2] = \frac{1 - f_{21}}{f^2 \cdot f_{21}} n S_2^2 E \left[ \underbrace{\frac{n_{21}}{n}}_{N_2/N} \right] = \\ &= \frac{1 - f_{21}}{f \cdot f_{21}} N_2 S_2^2 \end{aligned}$$

$$E_{n_1 n_2}(\hat{X}) = \frac{N}{n} X$$

$$V E_{n_1 n_2}(\hat{X}) = V \left( \frac{N}{n} x \right) = N^2 V(\bar{x}) = N^2 \cdot \frac{1 - f}{n} S^2$$

$$\text{Entonces: } V(\hat{X}) = \frac{1 - f_{21}}{f \cdot f_{21}} N_2 S_2^2 + N^2 \cdot \frac{1 - f}{n} S^2$$

donde  $S^2 \equiv$  cuasiv. poblacional y  $S_2^2 \equiv$  cuasiv. pobl. estrato 2.

que puede expresarse como:  $V(\hat{X})$  es la suma de la varianza del estimador del total si no hubiera falta de respuesta más el incremento de varianza debido al submuestreo de los no respondientes.

Para estimar la media poblacional  $\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$

$$\hat{\bar{X}} = \frac{1}{n} (X_1 + \hat{X}_2)$$

$$V(\hat{\bar{X}}) = \frac{1-f}{n} S^2 + \frac{1-f_2}{nNf_2} N_2 S_2^2$$

Para determinar el tamaño de la muestra, consideramos la función de coste total:

$$C = C_0 n + C_1 n_1 + C_2 n_2$$

Como  $n_2$  es una cantidad aleatoria, utilizamos el coste esperado:

$$C = n \left( C_0 + C_1 \frac{n_1}{n} + C_2 \frac{n_2}{n} \right) = n \left( C_0 + C_1 \left( \frac{n_1}{n} \right) + C_2 \cdot f_2 \frac{n_2}{n} \right)$$

$$\Rightarrow E(C) = n (C_0 + C_1 (1-W) + C_2 f_2 W)$$

Utilizando multiplicadores de Lagrange para optimizar el tamaño muestral para obtener una varianza dada  $V$ :

$$n = \underbrace{\frac{N^2 S^2}{V + N S^2}}_{\downarrow} + \frac{N N_2 \frac{1-f_2}{f_2} S_2^2}{V + N S^2} \quad \left. \vphantom{\frac{N^2 S^2}{V + N S^2}} \right\} \rightarrow \text{Penalización debida a la falta de respuesta}$$

Tamaño muestral necesario para estimar el total con un error  $E = \sqrt{V}$  si no hubiera falta de respuesta

### 3. TRATAMIENTO de la FALTA de RESPUESTA

Una primera idea consistió en seguir seleccionando unidades aleatoriamente hasta conseguir el tamaño muestral requerido. Conseguiríamos disminuir la varianza del estimador, pero no el sesgo, que permanecería constante ya que ~~se~~ todas las unidades de la muestra correspondían al estrato de los que sí responden.

Si se dispone de información directa del estrato de no respondedores que justifique la hipótesis  $\bar{X}_1 \approx \bar{X}_2$ , el estimador insesgado para el total sería:

$$\hat{X}_3 = \frac{N}{n} \sum_{i=1}^n X_i \quad , \text{ sólo insesgado si } \bar{X}_1 \approx \bar{X}_2$$

Si esta hipótesis no es sostenible, estaríamos en la situación de  $\hat{X}_1$ , que es sesgado, aunque su varianza es menor.

Otras técnicas están basadas en modelos, aplicables a situaciones particulares en las que se produce la falta de respuesta.

#### Mét. de Hansen y Hurwitz

Método de aplicación general, que inicialmente fue diseñado para encuestas que se enviaban por correo.

Sea  $n \equiv$  tamaño muestral, con  $\begin{cases} n_1 \text{ unid. que sí responden} \\ n_2 \text{ unid. que NO responden} \end{cases}$  (por m.a.s.)

De las  $n_2$  unid. que no responden, seleccionamos una muestra aleatoria de tamaño  $n_{21} \leq n_2$  y se envían entrevistadores para conseguir su respuesta en 2ª vuelta.

Para un  $n_i$  fijo, el  $u^o$  de entrevistas obtenidas en cada una de las clases es una multinomial y condicionada,  
 $E(\bar{X}_i/n_i) = \bar{X}_i$ , que no depende de  $n_i \Rightarrow$  la media de la muestra después de  $i$  visitas, independientemente del  $u^o$  de ~~visitas~~ respuestas, es también  $\bar{X}_i$ .

$$B(\bar{x}_i) = \bar{x}_i - \bar{X}$$

$$\sqrt{V(\bar{x}_i/n_i)}$$

### Modelo de respuesta aleatorizada. Modelo de Warner (1965)

Es difícil contar con la colaboración veraz de los encuestados frente a preguntas íntimas. Se puede dejar de contestar o bien dar deliberadamente respuestas falsas.

El individuo selecciona una de dos preguntas aleatoriamente, y contesta Sí o NO. El entrevistador desconoce la pregunta seleccionada  $\rightarrow$  Queda garantizado el secreto.

$\pi_A \rightarrow$  proporción desconocida de personas que contestan Sí a una pregunta íntima

$P \rightarrow$  probabilidad conocida de que se elige la pr. íntima

$\pi_y \rightarrow$  proporción conocida de respuestas Sí a una pr. intrascendente.

$n \rightarrow$  tamaño de la muestra

$$P(\text{Sí}) = \lambda = P \cdot \pi_A + (1-P) \pi_y \Rightarrow \pi_A = \frac{\lambda - (1-P) \pi_y}{P}$$

estimando  $\lambda$  con  $\hat{\lambda}$ :

$$\hat{\pi}_A = \frac{\hat{\lambda} - (1-P) \pi_y}{P}$$

$$V(\hat{\pi}_A) = V\left(\frac{\hat{\lambda}}{P}\right) = \frac{\lambda(1-\lambda)}{P^2 \cdot n}$$

#### 4. IMPUTACIÓN. TÉCNICAS DE REPONDERACIÓN

Con el propósito de disminuir el posible sesgo introducido por la falta de respuesta se suelen utilizar varios tipos de ajuste.

- Ajuste sobre el terreno: el entrevistador recibe el listado de las unidades muestrales y un listado de suplentes.

↳ No reduce el sesgo de los que no responden.

##### - Reponderaciones

Si se dispone información suplementaria sobre la proporción de unidades para ciertas clases de la población, se puede estratificar la muestra a posteriori y utilizar las mencionadas proporciones como reponderaciones.

Si no existe información suplementaria, se pueden utilizar como pesos las inversas de las tasas de respuesta.

##### - Imputación

Cuando un cuestionario no está completo, se puede realizar una imputación para los datos que faltan basados en la posible correlación entre el dato omitido y el resto de los datos disponibles.

Fichero caliente:

- 1.- Se establecen una serie de caracteres que se suponen correlacionados con el que queremos imputar.
- 2.- Se introducen en el ordenador unos valores iniciales, fichero frío, obtenidos de encuestas anteriores.
- 3.- Si en la primera ficha de la encuesta falta el dato, se imputa al fichero caliente. Si está el dato, se imputa al fichero frío,

y así sucesivamente

---

Se ha demostrado que el procedimiento "archivo caliente" produce un incremento ~~de~~ la variancia del estimador que no se reflejan en los métodos de estimación disponibles hasta ahora.

