

MUEST-T5 . ESTIMADOR de REGRESIÓN en el m.a.s.
 SESGO, VARIANZA y SUS ESTIMACIONES.
 COMPARACIONES CON EL EST. \hat{X}_R y \hat{X} .
 REFERENCIA a la ESTIM. GENERALIZADA
 por REGRESIÓN.

Var $\begin{cases} b_{cde} \\ b_{vcrmiu} \\ b_{pbliu} \end{cases}$

1. ESTIMADOR de REGRESIÓN en el MUEST. ALÉNT. SIMPLE

Los métodos indirectos aprovechan la información conocida relativa a una var. auxiliar Y correlacionada con la var. de estudio X para obtener estimaciones más precisas para X que las obtenidas solamente con la información de X .

Si entre X e Y hay una relación lineal (o aproximadamente lineal) ~~entonces~~ se puede utilizar un método de estimación basado en la regresión lineal de X sobre Y .
 (Si la recta pasase por el origen, se utilizaría el estimador de razón).

Para la población estudiada se observa el par (X_i, Y_i) de la variable en estudio y la variable auxiliar, con media poblacional conocida, \bar{Y} conocido.

Si los puntos están situados sobre una recta:

$$X_i = a + bY_i \Rightarrow \sum_{i=1}^N X_i = Na + b \sum_{i=1}^N Y_i \Rightarrow \bar{X} = a + b\bar{Y}$$

Teniendo una muestra (n):

$$x_i = a + bY_i \Rightarrow \sum_{i=1}^n x_i = na + b \sum_{i=1}^n Y_i \Rightarrow \bar{x} = a + b\bar{y}$$

luego $\bar{X} - \bar{x} = b(\bar{Y} - \bar{y})$, por lo que:

$$\bar{x}_{reg} = (\hat{\bar{X}}_{reg}) = \bar{x} + b(\bar{Y} - \bar{y}) \rightarrow \text{estimador de regresión de la media poblacional}$$

$$\begin{aligned}\hat{X}_{reg} &= N\bar{X}_{reg} = N\bar{x} + b(N\bar{y} - N\hat{y}) = \\ &= \hat{X} + b(\bar{y} - \hat{y}) \quad \leftarrow \text{estimador de regresión} \\ &\quad \text{del total poblacional.}\end{aligned}$$

En el caso de que $\bar{y} = \hat{y}$, el estimador de regresión de la media poblacional coincide con la media muestral

$$\bar{y} = \hat{y} \Rightarrow \bar{X}_{reg} = \bar{x}$$

Si $\bar{y} \neq \hat{y}$, también se espera que $\bar{X} \neq \bar{x}$, por lo que $b(\bar{y} - \hat{y})$ es el sumando de ajuste para una posible estimación por defecto o por exceso.

CASOS PARTICULARES:

$$\bar{X}_{reg} = \bar{x} + b(\bar{y} - \hat{y})$$

- $b = 0 \Rightarrow \bar{X}_{reg} = \bar{x}$ (estimador simple)
 \hookrightarrow ocurre cuando no hay correlación lineal entre X e Y .

- $b = \frac{\bar{x}}{\bar{y}} \Rightarrow \bar{X}_{reg} = \hat{X}_R$ (estimador de razón)

$$\bar{X}_{reg} = \bar{x} + \frac{\bar{x}}{\bar{y}}(\bar{y} - \hat{y}) = \bar{x} + \frac{\bar{x}}{\bar{y}}\bar{y} - \bar{x} = \frac{\bar{x}}{\bar{y}}\bar{y} = \hat{X}_R.$$

- $b = 1 \Rightarrow \bar{X}_{reg} = \bar{x}_{dif}$ (estimador por diferencia)

$$\bar{X}_{reg} = \bar{x} + (\bar{y} - \hat{y}) = \bar{x} - \hat{y} + \bar{y}$$

Para el total poblacional: $\hat{X}_{reg} = N\bar{X}_{reg}$

* César Pérez da la expresión para los estimadores de regresión del total de clase y de la proporción poblacional, ¿tiene sentido la regresión?

$$\hat{P}_{reg} = \hat{P}_x + b(P_y - \hat{P}_y) \quad \text{y} \quad \hat{A}_{reg} = N\hat{P}_{reg}.$$

2. SESGO, VARIANZA y ESTIMACIONES

SESGO :

El estimador de regresión es, en general, sesgado.

$$E[\bar{x}_{reg}] = E[\bar{x} + b(\bar{y} - \bar{y})] = E[\bar{x}] + \bar{y} E[b] - E[b\bar{y}] = \\ = E[\bar{x}] - (E[b\bar{y}] - E[b]E[\bar{y}]) = \bar{x} - \text{cov}(b, \bar{y})$$

$$B[\bar{x}_{reg}] = E[\bar{x}_{reg}] - \bar{x} = -\text{cov}(b, \bar{y})$$

El estimador de regresión será insesgado si $\text{cov}(b, \bar{y})$ es 0, que puede ocurrir cuando:

- $b=0 \Rightarrow$ no hay correlación lineal entre X e Y
- $b=cte \Rightarrow$ se elige el valor de b antes de observar la muestra.

VARIANZA del ESTIMADOR :

1) Para bale ($\neq 0$), el estimador de regresión es insesgado y su varianza es:

$$V[\bar{x}_{reg}] = V[\bar{x} + b(\bar{y} - \bar{y})] = V[\bar{x}] + \overset{\substack{\text{bale} \\ \downarrow}}{b^2} \overset{\substack{\text{conocido} \\ \downarrow}}{V[\bar{y} - \bar{y}]} + 2b\text{cov}[\bar{x}, \bar{y} - \bar{y}] = \\ = V[\bar{x}] + b^2 V[\bar{y}] - 2b\text{cov}[\bar{x}, \bar{y}] \quad \leftarrow \text{Expresión general}$$

expresión que se puede calcular para los \neq tipos de muestreo:

$$SR \rightarrow V[\bar{x}_{reg}] = \frac{1-f}{n} (s_x^2 + b^2 s_y^2 - 2b s_{xy})$$

que tiene por estimador insesgado

$$\hat{V}[\bar{x}_{reg}] = \frac{1-f}{n} (\hat{S}_x^2 + b^2 \hat{S}_y^2 - 2b \hat{S}_{xy})$$

Para el total poblacional:

$$V[\hat{X}_{reg}] = N^2 V[\bar{x}_{reg}] = N^2 \frac{(1-f)}{n} (s_x^2 + b^2 s_y^2 - 2b s_{xy})$$

$$\hat{V}[\hat{X}_{reg}] = N^2 \hat{V}[\bar{x}_{reg}] = N^2 \frac{(1-f)}{n} (\hat{S}_x^2 + b^2 \hat{S}_y^2 - 2b \hat{S}_{xy})$$

$$CR \rightarrow V[\bar{x}_{reg}] = \frac{1}{n} (\sigma_x^2 + b^2 \sigma_y^2 - 2b \sigma_{xy})$$

que tiene por estimador insesgado:

$$\hat{V}[\bar{x}_{reg}] = \frac{1}{n} (\hat{\sigma}_x^2 + b^2 \hat{\sigma}_y^2 - 2b \hat{\sigma}_{xy})$$

Para el total poblacional:

$$V[\hat{X}_{reg}] = N^2 V[\bar{x}_{reg}] = \frac{N^2}{n} (\sigma_x^2 + b^2 \sigma_y^2 - 2b \sigma_{xy})$$

$$\hat{V}[\hat{X}_{reg}] = N^2 \hat{V}[\bar{x}_{reg}] = \frac{N^2}{n} (\hat{\sigma}_x^2 + b^2 \hat{\sigma}_y^2 - 2b \hat{\sigma}_{xy})$$

2) Para b MÍNIMA varianza

Como la varianza del estimador de regresión depende de la de b , que usualmente se determina después de observar la muestra, se puede elegir el valor de la de b que haga mínima la varianza del estimador.

$$\underset{b}{\text{MIN}} V(\bar{x}_{reg}) = \phi(b) \rightarrow \begin{aligned} \text{CN: } \phi'(b) &= 0 \Rightarrow b_{\text{MIN}} \\ \text{CS: } \phi''(b) \big|_{b_{\text{MIN}}} &> 0 \Rightarrow \text{MÍNIMO.} \end{aligned}$$

Es un problema de optimización clásica;

$$\phi(b) = V(\bar{x}_{reg}) = V[\bar{x}] + b^2 V[\bar{y}] - 2b \text{COV}[\bar{x}, \bar{y}]$$

$$\phi'(b) = 2b V[\bar{y}] - 2 \text{COV}[\bar{x}, \bar{y}] = 0 \Leftrightarrow \underset{\text{MIN}}{b} = \frac{\text{COV}[\bar{x}, \bar{y}]}{V[\bar{y}]}$$

$$\phi''(b) \big|_{b_{\text{MIN}}} = 2 V[\bar{y}] \big|_{b_{\text{MIN}}} = 2 V[\bar{y}] > 0 \text{ siempre}$$

$$b_{\text{MIN}} = \frac{\text{COV}[\bar{x}, \bar{y}]}{V[\bar{y}]}, \text{ substituyendo esta expresión en}$$

la varianza del estimador:

$$\begin{aligned} V[\bar{x}_{reg}]_{\text{MIN}} &= V[\bar{x}] + \left(\frac{\text{COV}[\bar{x}, \bar{y}]}{V[\bar{y}]} \right)^2 V[\bar{y}] - 2 \frac{\text{COV}[\bar{x}, \bar{y}]}{V[\bar{y}]} \cdot \text{COV}[\bar{x}, \bar{y}] = \\ &= V[\bar{x}] - \frac{\text{COV}^2[\bar{x}, \bar{y}]}{V[\bar{y}]} = V[\bar{x}] (1 - \rho_{xy}^2) \end{aligned}$$

$$\left. \begin{aligned} b_{\text{MIN}} &= \frac{\text{COV}[\bar{x}, \bar{y}]}{V[\bar{y}]} \\ V_{\text{MIN}}(\bar{x}_{\text{req}}) &= V[\bar{x}](1 - \rho_{xy}^2) \end{aligned} \right\} \leftarrow \text{Expresión general}$$

Se puede particularizar para los \neq tipos de muestreo:

$$\text{SR} \rightarrow \begin{cases} b_{\text{MIN}} = \frac{S_{xy}}{S_y^2} = \beta \\ V_{\text{MIN}}(\bar{x}_{\text{req}}) = \frac{1-f}{n} S_x^2 (1 - \rho_{xy}^2) \\ \hat{V}_{\text{MIN}}(\bar{x}_{\text{req}}) = \frac{1-f}{n} \hat{S}_x^2 (1 - \hat{\rho}_{xy}^2) \end{cases}$$

$$\begin{aligned} \text{Por el total: } V_{\text{MIN}}(\hat{x}_{\text{req}}) &= N^2 V_{\text{MIN}}(\bar{x}_{\text{req}}) && \uparrow \text{SR} \\ \hat{V}_{\text{MIN}}(\hat{x}_{\text{req}}) &= N^2 \hat{V}_{\text{MIN}}(\bar{x}_{\text{req}}) && \downarrow \text{CR} \end{aligned}$$

$$\text{CR} \rightarrow \begin{cases} b_{\text{MIN}} = \frac{G_{xy}}{G_y^2} = \frac{S_{xy}}{S_y^2} = \beta \\ V_{\text{MIN}}(\bar{x}_{\text{req}}) = \frac{1}{n} G_x^2 (1 - \rho^2) \\ \hat{V}_{\text{MIN}}(\bar{x}_{\text{req}}) = \frac{1}{n} \hat{S}_x^2 (1 - \hat{\rho}^2) \end{cases}$$

3) En principio, el valor $b_{\text{MIN}} = \frac{S_{xy}}{S_y^2}$ no depende de la muestra obtenida, por lo que puede fijarse de antemano. Pero los parámetros poblacionales S_{xy} y S_x^2 no siempre son conocidos $\rightarrow \hat{b}_{\text{MIN}} = \frac{\hat{S}_{xy}}{\hat{S}_y^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (y_i - \bar{y})^2}$

que es un estimador insesgado de b_{MIN} para muestras grandes

3. COMPARACIONES del ESTIMADOR de ^{REGRESIÓN} ~~RAZÓN~~ con EL ESTIMADOR de EXPANSIÓN y con EL DE RAZÓN

A la hora de comparar dos estimadores, será más preciso aquel con menor varianza.

Necesitamos un tamaño muestral sufic. grande para que las aproximaciones de la varianza de los estm. de razón y de regresión sean válidas.

Utilizaremos el estimador de la media poblacional y su varianza en los diferentes tipos de muestreo.

* En el caso general, las expresiones de la varianzas difieren bastante unas de otras, por lo que resulta más apropiado hacer las comparaciones para los diferentes tipos de muestreo.

Parámetro poblacional: $\theta = \bar{X} \rightarrow$ media poblacional
 Estimadores: $\hat{\bar{X}} = \begin{cases} \bar{X} & \rightarrow \text{media muestral (estm. expansión)} \\ \hat{\bar{X}}_R = \hat{R} \cdot \bar{Y} & \rightarrow \text{estm. razón} \\ \hat{\bar{X}}_{reg} = \bar{X} + b(\bar{Y} - \bar{y}) & \rightarrow \text{estm. regresión} \end{cases}$

cuyas varianzas cambian de expresión al cambiar el tipo de muestreo:

$$\begin{aligned} 1) \text{ SR} &\rightarrow V(\bar{X}) = \frac{(1-f)}{n} S_x^2 \\ V(\hat{\bar{X}}_R) &= \frac{1-f}{n} (S_x^2 + R^2 S_y^2 - 2R S_x S_y \rho_{xy}) \\ V(\hat{\bar{X}}_{reg}) &= V_{\min}(\bar{X}_{reg}) = \frac{1-f}{n} S_x^2 (1 - \rho_{xy}^2) \end{aligned}$$

Es evidente que $V(\bar{x}_{reg}) \leq V(\bar{x})$ porque:

$$1 - r^2_{xy} \leq 1 \rightarrow 0 \leq r^2_{xy} \leq 1 \Rightarrow 1 - r^2_{xy} \leq 1.$$

$$1 - r^2_{xy} = 1 \Leftrightarrow r^2_{xy} = 0 \Leftrightarrow X \text{ e } Y \text{ in correladas}$$

⇒ La estimación por regresión es mejor que la estimación por muestreo aleatorio simple, excepto cuando las var. están in correladas, ~~en~~ caso en que son iguales.

Por otra parte,

$$V(\bar{x}_{reg}) < V(\bar{x}_R) \Leftrightarrow V(\bar{x}_{\cancel{R}}) - V(\bar{x}_{\cancel{reg}}) > 0$$

$$\frac{1}{n} \left[\cancel{S_x^2} + R^2 S_y^2 - 2RS_x S_y r_{xy} - \cancel{S_x^2} + S_x^2 r^2_{xy} \right] \geq 0$$

$$\Leftrightarrow R^2 S_y^2 - 2RS_x S_y r_{xy} + S_x^2 r^2_{xy} \geq 0 \Leftrightarrow$$

cuadrado de una dif

$$(RS_y - S_x r_{xy})^2 \geq 0 \rightarrow \text{SIEMPRE}$$

↳ Se produce la igualdad si $R = \frac{S_x}{S_y} r_{xy} = \beta$.
 ⇒ regresión pasa por origen

⇒ La estimación por regresión es más precisa que la estimación por muestreo en todos los casos, salvo cuando la recta de regresión pasa por el origen, en que las 2 est. son igual de precisas.

$$2) CR \rightarrow V(\bar{X}) = \frac{1}{n} \sigma_x^2$$

$$V(\bar{X}_R) = \frac{1}{n} (\sigma_x^2 + R^2 \sigma_y^2 - 2R \sigma_x \sigma_y \rho_{xy})$$

$$V_{\min}(\bar{X}_{\text{reg}}) = \frac{1}{n} \sigma_x^2 (1 - \rho_{xy}^2)$$

• $V_{\min}(\bar{X}_{\text{reg}}) \leq V(\bar{X})$ porque $1 - \rho_{xy}^2 \leq 1$
 $1 - \rho_{xy}^2 = 1$ si X e Y incorrelados

• $V_{\min}(\bar{X}_{\text{reg}}) \leq V(\bar{X}_R)$ siempre
 $V(\bar{X}_{\text{reg}}) = V(\bar{X}_R)$ si $R = \beta \rightarrow$ recta por el origen

$$V(\bar{X}_R) - V_{\min}(\bar{X}_{\text{reg}}) \geq 0 \Leftrightarrow$$

$$\frac{1}{n} [\cancel{\sigma_x^2} + R^2 \sigma_y^2 - 2R \sigma_x \sigma_y \rho_{xy} - \cancel{\sigma_x^2} + \sigma_x^2 \rho_{xy}^2] \geq 0$$

$$\frac{1}{n} [(R \sigma_y - \sigma_x \rho_{xy})^2] \geq 0 \text{ siempre}$$

$$= 0 \text{ si } R = \frac{\sigma_x}{\sigma_y} \rho_{xy} = \beta.$$

4. REFERENCIA A LA ESTIMACIÓN GENERALIZADA por REGRESIÓN

la estimación por regresión se puede generalizar al caso en el que se dispone de K variables auxiliares que tengan correlación con la var. a estudiar X ,

Sea \vec{Y}_i el vector de variables auxiliares observadas sobre el elemento poblacional U_i .

Para $U_i, i=1 \dots N \longrightarrow \vec{Y}_i = (Y_{1i} \dots Y_{Ki})$

Suponemos que cada valor de la var. a estudiar se puede relacionar de manera lineal con \vec{Y}_i :

$$X_i = \sum_{j=1}^K \beta_j Y_{ji} + \varepsilon_i, \quad i=1 \dots N \longrightarrow X = Y' \beta + \varepsilon$$

donde

$$E[X_i] = \sum_{j=1}^K \beta_j Y_{ji}, \quad i=1 \dots N \quad (\text{sp. } Y \text{ determinista})$$

$$V[X_i] = \sigma_i^2, \quad i=1 \dots N$$

$$\varepsilon_i \sim N(0, \sigma_i^2) \quad i=1 \dots N \rightarrow \text{perturbación o término aleat.}$$

Al estar las var. auxiliares correlacionadas, el modelo incumple la hipótesis de ~~homoscedasticidad~~ ^{homoscedasticidad} ($\sigma_i^2 \neq \sigma^2_{ale}$),

por lo que para estimar el vector de coeficientes de regresión $\vec{\beta} = (\beta_1 \dots \beta_K)^t$ hay que utilizar Mínimos Cuadrados Generalizados (ponderados).

$$\vec{\beta} = \left[\sum_{i=1}^N \frac{\vec{Y}_i \vec{Y}_i^t}{\sigma_i^2} \right]^{-1} \sum_{i=1}^N \frac{\vec{Y}_i X_i}{\sigma_i^2}$$

$$\vec{\beta} = (Y^{*t} Y^*)^{-1} Y^{*t} X^* \quad \text{donde } ()^* = \frac{()}{\sigma_i^2}$$

Al depender de los valores poblacionales $X_1 \dots X_N$, el estimador $\hat{\beta}$ no se puede obtener, por lo que hay que estimarlo a partir de las observaciones muestrales.

En el caso de w.a.s.r.p. desij:

$$\hat{\beta} = \left(\sum_{i=1}^n \frac{\vec{y}_i \vec{y}_i^t}{\sigma_i^2 \pi_i} \right)^{-1}_{K \times K} \cdot \left(\sum_{i=1}^n \frac{\vec{y}_i X_i}{\sigma_i^2 \pi_i} \right)_{K \times 1} = ()_{K \times 1}$$

Por lo que el estimador generalizado de regresión para el total poblacional resulta:

$$\hat{X}_{G,reg} = \hat{X} + \sum_{j=1}^K \hat{\beta}_j (Y_j - \hat{Y}_j)$$

donde:

Y_j \rightarrow total poblacional de la var. auxiliar y_j .

\hat{X} \rightarrow estimador insesgado del total de X , $\hat{X} = \sum_{i=1}^n \frac{x_i}{\pi_i}$

\hat{Y}_j \rightarrow estimador insesgado del total de Y

$\hat{\beta}$ \rightarrow estimador generalizado de regresión para X .

La varianza de $\hat{\beta}$ no se puede calcular, pero se pueden aplicar técnicas de ~~regresión~~ linealización para obtener una aproximación del estimador y de su varianza.