

# Métodos indirectos de estimación: razón, regresión y diferencia

## Contenido

1. Introducción, definición de estimadores indirectos
2. Estimador de razón, propiedades y varianzas. Límites de confianza.
3. Tamaño de la muestra en los estimadores de razón
4. Eficiencia de los estimadores de razón
5. Estimadores de razón en el muestreo estratificado: razón separado y razón combinado.
6. Estimadores de regresión con “b” pre-asignada y con “b” estimada.
7. Comparación entre los estimadores de razón, de regresión y simple.
8. Estimadores de regresión en el muestreo estratificado: separado y combinado.
9. Estimador de diferencia.

# Métodos indirectos de estimación

Es frecuente en el muestreo usar información adicional para mejorar la precisión de las estimaciones, esto se conoce como métodos indirectos y los más comunes son los de razón y los de regresión.

Los métodos indirectos requieren el poder observar una variable auxiliar  $X$  apareada con la variable en estudio  $Y$  ( $y_i, x_i$ ) y además conocer el total  $X$ .

La correlación existente entre  $X$  e  $Y$  es lo que permite el incremento de la precisión al estimar los parámetro de  $Y$ .

Los estimadores indirectos se usan para estimar el total, la media y la razón propiamente.

# Métodos indirectos de estimación

Para el total  $\hat{Y}$  la forma general de los estimadores indirectos es:

$$\hat{Y}_I = \hat{Y} + \alpha(X - \bar{X})$$

Donde  $\alpha$  es el coeficiente de corrección para mejorar el estimador directo  $\hat{Y}$

Si  $\alpha = 0$   $\hat{Y}_I \equiv \hat{Y}$  estimador directo

Si  $\alpha = \hat{R} = \frac{\hat{Y}}{\hat{X}}$   $\hat{Y}_I = \hat{R}X$  estimador de razón

Si  $\alpha = 1$   $\hat{Y}_I = \hat{Y} + X - \hat{X}$  estimador de diferencia

Si  $\alpha = b$   $\hat{Y}_I = \hat{Y} + b(X - \hat{X})$  estimador de regresión  
b coeficiente de regresión lineal Y/X

# Estimador de razón en el Muestreo Aleatorio Simple.

En el mas el estimador de razón del total es:  $\hat{Y}_R = \hat{R} X = \frac{\bar{y}}{\bar{x}} X$

$X_i$  debe ser una variable apareada con  $Y_i$  y correlacionada entre si,  $X_i$  puede ser el valor de  $Y_i$  en una ocasión anterior ( $X_t=Y_{t-s}$ )

Si la razón  $y_i/x_i$  es la misma en todas las unidades de muestreo,

$\bar{y}/\bar{x}$  es estable de muestra en muestra y  $\hat{Y}_R$  será mas precisa. El método se justifica al suponer que  $Y_i$  varía proporcional a  $X_i$   
( $y_i = kx_i \rightarrow \bar{y} = k\bar{x}$ )

El estimador de razón  $\hat{Y}_R$  (o  $\bar{y}_R$ ) es:

1. Consistente en el sentido de Cochran
2. Sesgado, pero el sesgo es despreciable en muestras grandes
3. En muestras grandes se puede utilizar la distribución normal (grandes si  $n > 30$  y  $CV(\bar{x}) < 10\%$  y  $CV(\bar{y}) < 10\%$ ), en muestras medianas la distribución es asimétrica positiva

# Estimadores de razón del total, la media, y la razón.

Estimadores

Sus varianzas (aproximadas)

$$\hat{Y}_R = \hat{R} X$$

$$V(\hat{Y}_R) = \frac{N^2(1-f)}{n} \left[ \frac{\sum (y_i - Rx_i)^2}{N-1} \right]$$

$$\bar{y}_R = \hat{R} \bar{x}$$

$$V(\bar{y}_R) = \frac{1-f}{n} \left[ \frac{\sum (y_i - Rx_i)^2}{N-1} \right]$$

$$\hat{R} = \frac{\bar{y}}{\bar{x}}$$

$$V(\hat{R}) = \frac{1-f}{n\bar{X}^2} \left[ \frac{\sum (y_i - Rx_i)^2}{N-1} \right]$$

Los estimadores son “aproximadamente” insesgados y las pruebas de las varianzas se vieron en el tema 2 para  $\hat{R}$  la de los otros dos sigue inmediato.

# Formas de expresión de la varianza

Como

$$\begin{aligned} \sum^N (y_i - Rx_i)^2 &= \sum^N (y_i - Rx_i - \bar{Y} - R\bar{X})^2 = \sum^N ((y_i - \bar{Y}) - R(x_i - \bar{X}))^2 \\ &= \sum^N (y_i - \bar{Y})^2 + R^2 \sum^N (x_i - \bar{X})^2 - 2R \sum^N (y_i - \bar{Y})(x_i - \bar{X}) \end{aligned}$$

Y

$$S_Y^2 = \frac{\sum^N (y_i - \bar{Y})^2}{N-1}; \quad S_X^2 = \frac{\sum^N (x_i - \bar{X})^2}{N-1} \quad \text{y} \quad \rho = \frac{\sum^N (y_i - \bar{Y})(x_i - \bar{X})}{(N-1)S_X S_Y} = \frac{S_{XY}}{S_X S_Y}$$

Así  $V(\hat{Y}_R) = \frac{N^2(1-f)}{n} (S_Y^2 + R^2 S_X^2 - 2R\rho S_X S_Y)$

$$= (1-f) \frac{\bar{Y}^2}{n} \left( \frac{S_Y^2}{\bar{Y}^2} + \frac{S_X^2}{\bar{X}^2} - \frac{2S_{XY}}{\bar{Y}\bar{X}^2} \right) = (1-f) \frac{\bar{Y}^2}{n} (C_{YY} + C_{XX} - 2C_{XY})$$

donde  $C_{YY}$  y  $C_{XX}$  son los cuadrados de los coeficientes de variación y  $C_{XY}$  es la covarianza relativa.

El cuadrado del coeficiente de variación de  $\hat{Y}_R$  se denomina

varianza relativa  $\left( CV(\hat{Y}_R) \right)^2 = \frac{V(\hat{Y}_R)}{\bar{Y}^2} = \frac{(1-f)}{n} (C_{YY} + C_{XX} - 2C_{XY})$

## Estimaciones de la varianza

Las estimaciones  $\hat{V}(\hat{R})$ ,  $\hat{V}(\hat{Y}_R)$  y  $\hat{V}(\bar{y}_R)$  se obtienen al sustituir  $\frac{\sum^N (y_i - Rx_i)^2}{N-1}$  por  $\frac{\sum^n (y_i - Rx_i)^2}{n-1}$  en las formulas de  $V(\bullet)$ .

Cuando  $\bar{X}$  o  $R$  son desconocidos se sustituyen por sus estimaciones  $\bar{x}$  o  $\hat{R}$ .

## Límites de confianza

Si  $n$  es suficientemente grande y,  $CV(\bar{x})$  y  $CV(\bar{y})$  son ambos  $<0,1$  se puede aplicar razonablemente la aproximación normal.

Obteniendo los límites de confianza para  $Y$ ,  $\bar{Y}$  y  $R$ .

$$R \quad \hat{R} \mp t_{\alpha} \sqrt{\hat{V}(\hat{R})} \quad Y \quad \hat{Y}_R \mp t_{\alpha} \sqrt{\hat{V}(\hat{Y}_R)} \quad \bar{Y} \quad \hat{\bar{Y}} \mp t_{\alpha} \sqrt{\hat{V}(\hat{\bar{Y}})}$$

Cuando las condiciones exigidas no se cumplen la distribución puede ser de notoria asimetría. En estos casos existen métodos alternos para calcular los límites de confianza.

# Tamaño de la muestra en los estimadores de razón

El procedimiento para obtener el tamaño de la muestra para estimar  $R$ , ( $Y$  o  $\bar{Y}$ ) en el muestreo aleatorio simple es igual al ya estudiado; fijado el riesgo ( $\alpha$ ) que estamos dispuestos a afrontar de que el error real supere a  $e$  en  $\Pr\left(\left|\hat{R}-R\right|\geq e\right)=\alpha$  y suponiendo que  $\hat{R} \sim n(R, \sigma_R)$ ,

$$\text{de } e = t_\alpha \sqrt{\hat{V}(\hat{R})} \rightarrow \frac{e^2}{t_\alpha^2} = \frac{N-n}{nN} \frac{s_d^2}{\bar{X}^2} \quad \text{donde } s_d^2 = \frac{\sum (y_i - \hat{R} x_i)^2}{n-1}$$

$$\text{Así } n = \frac{n_o}{1 + \frac{n_o}{N}} \quad \text{donde } n_o = \frac{t_\alpha^2 s_d^2}{e^2 \bar{X}^2}$$

$s_d^2$  es calculada a priori de acuerdo a información disponible.



# Tamaño de la muestra en los estimadores de razón

De igual manera se pueden obtener los tamaños de muestras requeridos para estimar  $Y$  o  $\bar{Y}$  usando estimadores de razón

Obteniendo  $n_o = \frac{t_\alpha^2 S_d^2}{e^2}$  para  $\bar{Y}_R$

y  $n_o = \frac{N t_\alpha^2 S_d^2}{e^2}$  para  $Y_R$ , en este caso  $n = \frac{n_o}{1 + n_o} N$

# Eficiencia de los estimadores.

La eficiencia de los estimadores de razón (en el mas), se mide en relación con la de los estimadores directos. En muestras grandes (mas)

el estimador de razón del total  $\hat{Y}_R$  tiene una varianza menor que el estimador directo  $\hat{Y} = N\bar{y}$  si:  $\rho > \frac{CV(X)}{2CV(Y)}$  así:

$$V(\hat{Y}) > V(\hat{Y}_R) \rightarrow \frac{N^2(1-f)}{n} S_Y^2 > \frac{N^2(1-f)}{n} (S_Y^2 + R^2 S_X^2 - 2R\rho S_Y S_X)$$

$$\rightarrow 2R\rho S_Y S_X > R^2 S_X^2 \quad \text{si} \quad R = \frac{\bar{y}}{\bar{x}} > 0 \quad \rightarrow \rho > \frac{R S_X}{2 S_Y} \rightarrow \rho > \frac{CV(X)}{2CV(Y)}$$

Se demuestra que los estimadores de razón son estimadores óptimos dentro de la clase de estimadores insesgados si:

1. Existe una relación lineal entre X,Y que pasa por el origen y se mantiene para "todo" par  $(X_i, Y_i)$
2. La varianza alrededor de la línea recta es proporcional a  $X_i$ , es decir,  $V(Y/X) = X_i V(\cdot)$

Se pueden comprobar estas 2 condiciones haciendo un grupo de dispersión  $(X_i, Y_i)$  (caso particular  $X_t = Y_{t-1}$ )

# Estimadores de razón en el muestreo estratificado.

En el muestreo estratificado, los estimadores de razón pueden ser más eficientes que los estimadores directos si los tamaños de muestra en cada estrato son lo suficientemente grandes como para que las varianzas sean las apropiadas y se cumple las condiciones de optimización.

Hay dos tipos de estimadores de razón (para  $Y$  o  $\bar{Y}$ ) en el muestreo estratificado: de razón separada y de razón combinada.

# Estimadores de razón separada

En las estimaciones de razón separada, la razón  $R_h$  en cada estrato  $h$  de la población se estima por separado, se pondera por los respectivos totales

$X_h$  y luego se suma sobre todos los estratos  $\hat{Y}_{RS} = \sum \frac{\bar{y}_h}{X_h} X_h$

Si se supone que la verdadera razón permanece constante al pasar de un estrato a otro y además se requiere conocer  $X_h$   $h=1,2,..,L$

Como la selección en cada estrato es aleatoria simple, independiente y tamaño suficientemente grande se cumple que:

$$V(\hat{Y}_{RS}) > \sum \frac{N_h^2(1-f_h)}{n_h} (S_{Y_h}^2 + R_h^2 S_{X_h}^2 - 2R_h \rho_h S_{Y_h} S_{X_h})$$

La prueba es inmediata

Cuando los  $n_h$  son pequeños  $\hat{Y}_{RS}$  puede poseer un sesgo no despreciable

# Estimación de razón combinada

Este estimador consiste en estimar previamente  $\hat{Y}_{st}$  y  $\hat{X}_{st}$  y luego hacer una sola estimación estimada.

Con los datos muestrales se calcula:  $\hat{Y}_{st} = \sum N_h \bar{y}_h$  y  $\hat{X}_{st} = \sum N_h \bar{x}_h$

La estimación de razón combinada del total estimado Y es:

$$\hat{Y}_{RC} = \frac{\hat{Y}_{st}}{\hat{X}_{st}} X = \frac{\bar{y}_{st}}{\bar{x}_{st}} X \quad \text{donde} \quad \bar{y}_{st} = \frac{\hat{Y}_{st}}{N} \quad \text{y} \quad \bar{x}_{st} = \frac{\hat{X}_{st}}{N}$$

La estimación de  $\hat{Y}_{RC}$  no requiere del conocimiento de los totales por estrato  $X_h$  sino del total general X, esta estimación está menos sujeta a riesgo de sesgo que de razón separada.

Se comprueba:  $V(\hat{Y}_{RC}) > \sum \frac{N_h^2(1-f_h)}{n_h} (S_{Y_h}^2 + R^2 S_{X_h}^2 - 2R\rho_h S_{Y_h} S_{X_h})$

para las estimaciones muestrales de estas varianzas se sustituye los  $R_h$  y  $\rho_h$  y los  $S_{Y_h}$   $S_{X_h}$  y  $S_{XY}$  por sus respectivas estimaciones muestrales.

# Estimadores de regresión

Los estimadores de regresión (al igual que los de razón) se justifican en el incremento de la precisión de las estimaciones, que brindan ante los estimadores directos.

Estos estimadores se utilizan propiamente cuando:

- a) la relación entre Y y X es lineal y no pasa por el origen
- b) es menos costoso observar  $x_i$  que  $y_i$
- c) se quieren predecir valores particulares de  $y_i$ .

El estimador de  $\bar{Y}$  de regresión es:  $\bar{y}_{rl} = \bar{y} + b(\bar{X} - \bar{x})$

y para el total  $\hat{Y}_{rl} = N \bar{y}_{rl}$

donde **b** es una estimación del coeficiente de regresión lineal de y en x.

# Estimadores de regresión

El fundamento de esta estimación es que si  $\bar{x}$  está por debajo del promedio, también deberíamos esperar que  $\bar{y}$  lo estuviera por la cantidad  $b(\bar{X} - \bar{x})$  debido a la regresión de  $y_i$  en  $x_i$ . Los estimadores de regresión tienen las mismas propiedades que los de razón, son consistentes en el sentido de Cochran pero suelen ser sesgados, pero la razón del sesgo al error estándar se reduce cuando  $n$  es grande y en este caso se puede usar la aproximación normal.

**b** se obtiene por estimaciones muestrales, o se obtiene de los valores muestrales o con información previa a la muestra, así, se distinguen dos casos:

1. Estimadores de regresión con “**b**” preasignada
2. Estimadores de regresión con “**b**” calculada en la muestra

# Estimador de regresión con b presasignada

Si  $b_0$  es una constante calculada con información pre-muestral, el estimador de regresión lineal es:  $\bar{y}_{rl} = \bar{y} + b_0(\bar{X} - \bar{x})$  el cual es insesgado con varianza:

$$V(\bar{y}_{rl}) = \frac{1-f}{n} \left( \sum [(y_i - \bar{Y}) - b_0(x_i - \bar{X})]^2 \right)$$

$$V(\bar{y}_{rl}) = \frac{1-f}{n} (S_Y^2 - 2b_0 S_{YX} + b_0^2 S_X^2)$$

cuyo estimador insesgado es  $\hat{V}(\bar{y}_{rl}) = \frac{1-f}{n} (s_Y^2 - 2b_0 s_{YX} + b_0^2 s_X^2)$

El valor de  $b_0$  que minimiza  $\hat{V}(\bar{y}_{rl})$  es  $b_0 = B = \frac{s_{YX}}{s_X s_Y}$  calculado sobre la población.



# Estimador de regresión con b estimada en la muestra

Si **b** no es conocida previo a la muestra, debe calcularse con los

datos muestrales por:

$$b = \frac{\sum (y_i - \bar{Y})(x_i - \bar{X})}{\sum (x_i - \bar{X})^2}$$

**b** así calculada, es un estimador, una variable aleatoria y para poderlo utilizar adecuadamente debemos exigirle que la relación entre Y y X sea aproximadamente lineal, la varianza alrededor de la línea de regresión sea constante (homoscedasticidad), que el tamaño de la muestra sea grande y el muestreo sea aleatorio

simple. En este caso:  $\bar{y}_{rl} = \bar{y} + b(\bar{X} - \bar{x})$

cuya varianza es  $V(\bar{y}_{rl}) = \frac{1-f}{n} S_Y^2 (1 - \rho^2)$  donde  $\rho = \frac{S_{YX}}{S_X S_Y}$

cuya estimación muestral válida para muestras grandes es:

$$V(\bar{y}_{rl}) = \frac{1-f}{n(n-2)} \frac{\sum [(y_i - \bar{y}) - b(x_i - \bar{x})]^2}{\sum (x_i - \bar{x})^2}$$

# Comparación entre los estimadores de razón, regresión y simples (n grande)

Si  $n$  es suficientemente grande de modo que  $V(\bar{y}_{rl})$  y  $V(\bar{y}_R)$  aproximadas sean válidas, se cumple que bajo ciertas condiciones  $V(\bar{y}_{rl}) < V(\bar{y})$   $V(\bar{y}_{rl}) < V(\bar{y}_R)$  y como vimos  $V(\bar{y}_R) < V(\bar{y})$

Recordemos 
$$V(\bar{y}_{rl}) = \frac{N-n}{Nn} S_Y^2 (1 - \rho^2) \quad V(\bar{y}) = \frac{N-n}{Nn} S_Y^2$$

$$V(\bar{y}_R) = \frac{N-n}{Nn} (S_Y^2 + R^2 S_X^2 - 2R\rho S_Y S_X)$$

Entonces

$$V(\bar{y}_{rl}) < V(\bar{y}) \quad \text{Se cumple para todo } \rho \text{ y es igual para } \rho = 0$$

$$V(\bar{y}_{rl}) < V(\bar{y}_R) \quad \text{si } -\rho^2 S_Y^2 < R^2 S_X^2 - 2R\rho S_Y S_X \quad \text{o equivalente}$$

$$(\rho S_Y - R S_X)^2 > 0 \quad \text{o} \quad (B - R)^2 > 0$$

así la  $\bar{y}_{rl}$  es mas precisa que  $\bar{y}_R$  a menos que  $B=R$  y esto ocurre si la regresión pasa por el origen.