

MUEST. T13. MÉTODOS SIMPLIFICADOS de ESTIMACIÓN de VARIANZAS en ENCUESTAS COMPLEJAS.

MT. de los GRUPOS ALEATORIOS.

MT. de los CONGLOMERADOS ÚLTIMOS.

MT. de SEMIMUESTRAS REITERADAS.

MT. JACKKNIFE.

MT. BOOTSTRAP.

1 - MT. SIMPLIFICADOS de ESTIMACIÓN de VARIANZAS en ENCUESTAS COMPLEJAS

En la práctica, suelen utilizarse esquemas de muestreo poli-
tápico complejos, en los que se estiman muchas características
poblacionales. En estos casos, las fórmulas ordinarias de
estimación de varianzas son de difícil aplicación, pues
conducen a tediosos cálculos.

Este hecho ha llevado al desarrollo de técnicas más
sencillas de estimación de varianzas, aunque resulten meno
precisas, son de más fácil aplicación.

De entre estas técnicas ya conocemos el mt. de la
muestra interpenetrante, que se utiliza en muestreo siste-
mático y el mt. de los conglomerados últimos, que volve-
remos a ver.

A continuación veremos:

Ø - INTRODUCCIÓN

Muestreo

* Concepto

* Tipos

Muestreo polietápico

* Concepto

* Estimación

* Tur. Rodow.

* Tur. Durbin.

1 - MTS SIMPLIFICADOS

Caso general \rightarrow Tur. Durbin

Caso particular \rightarrow aquí

| \rightarrow Tur. atrás

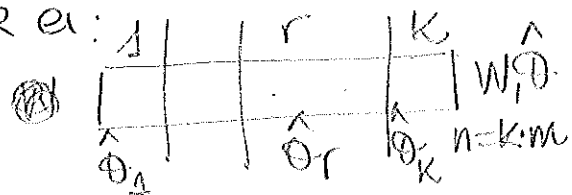
2 - MT. de los GRUPOS ALEATORIOS

A partir de una muestra de tamaño n obtenida de una población finita de N unidades, el mt. consiste en subdividir la muestra en K submuestras de tamaño m , de modo que $n = K \cdot m$. Así, cada grupo aleatorio es una submuestra de la muestra y a su vez una muestra de la población completa, de menor tamaño pero con las mismas propiedades probabilísticas que la muestra concreta.

La formación de los K grupos aleatorios de tamaño m dentro de una muestra W de tamaño n puede realizarse considerando una permutación aleatoria de los n 1, 2, ..., n , y eligiendo el primer m aleatorio formado por los elementos de la muestra que ocupan los lugares definidos por los m primeros números de la permutación. Así sucesivamente.

Si $\hat{\theta}$ es un estimador insesgado de la característica poblacional θ basado en la muestra completa W , y si $\hat{\theta}_r$ es un estimador insesgado de θ basado en el r -ésimo grupo aleatorio, un estimador insesgado de la varianza de $\hat{\theta}$ en muestreo CR es:

$$(R) \hat{V}(\hat{\theta}) = \frac{1}{K(K-1)} \sum_{r=1}^K (\hat{\theta}_r - \hat{\theta})^2$$



(82) (1-f)
El caso más sencillo de aplicación es el de muestreo aleatorio simple con reposición:

El estimador insesgado del total poblacional X es $\hat{X} = N\bar{x}$ con $\bar{x} = \frac{1}{n} \sum_{i=1}^n X_i$, media muestral de la muestra completa W

$$V(\hat{X}) = V(N\bar{x}) = N^2 V(\bar{x}) = N^2 \cdot \frac{\sigma^2}{n}, \text{ con } \sigma^2 = \frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})^2$$

Como la r -ésima muestra, o r -ésimo grupo aleatorio $(X_{r1} \dots X_{rm})$ es también una muestra de la poblac. total, podemos estimar el total X_r :

$$\hat{X}_r = N \bar{X}_r = N \cdot \frac{\sum_{j=1}^m X_{rj}}{m}$$

$$V(\hat{X}_r) = N^2 V(\bar{X}_r) = N^2 \cdot \frac{\sigma^2}{m} = \overset{m \cdot k}{N^2 \cdot \frac{\sigma^2}{n}} = k \cdot N^2 \frac{\sigma^2}{n} = k V(\hat{X})$$

Como el r -ésimo grupo aleatorio $\frac{k}{K}$ es también una submuestra de la muestra completa $W = (X_1 \dots X_K)$:

$$E_W(\hat{X}_r) = E_W\left(N \cdot \frac{\hat{X}_r}{m}\right) = E_W(N \bar{X}_r) = N \bar{X} = \hat{X}$$

luego

$$E[\hat{X}_r] = E E_W(\hat{X}_r) = E[\hat{X}] = X.$$

\hat{X}_r es una copia de \hat{X} para el grupo aleatorio r -ésimo de W , y es un estimador insesgado de X .

$$V(\hat{X}_r) = E(\hat{X}_r - \hat{X})^2 = E(\hat{X}_r - \hat{X} + \hat{X} - X)^2 = E(\hat{X}_r - \hat{X})^2 + E(\hat{X} - X)^2$$

$$\Rightarrow E(\hat{X}_r - \hat{X})^2 = V(\hat{X}_r) - V(\hat{X}) = k V(\hat{X}) - V(\hat{X}) = (k-1) V(\hat{X})$$

Sumando en $r=1 \dots K$ los dos términos de la igualdad:

$$\sum_{r=1}^K E(\hat{X}_r - \hat{X})^2 = \sum_{r=1}^K (k-1) V(\hat{X}) = K(k-1) V(\hat{X}) \Rightarrow$$

por lo que $\Rightarrow V(\hat{X}) = E\left(\frac{\sum_{r=1}^K (\hat{X}_r - \hat{X})^2}{K(k-1)}\right)$

luego $\frac{1}{K(k-1)} \sum_{r=1}^K (\hat{X}_r - \hat{X})^2$ es un estimador insesgado de $V(\hat{X})$.

- Este mt. también se puede utilizar con grupos aleatorios de distinto tamaño m_r / $\sum_{r=1}^K m_r = n$

$$\hat{V}(\hat{\theta}) = \frac{1}{K-1} \sum_{r=1}^K \frac{m_r}{n} (\hat{\theta}_r - \hat{\theta})^2 \text{ estim. insesgado de } V(\hat{\theta})$$

- El mt. también puede aplicarse a:

- Muestreo polietápico, formando los K grupos aleatorios con las n unidades primarias que constituyen la muestra complex.
- Muestreo estratificado, tomando K grupos aleatorios en cada estrato
- Estimadores de razón, siempre que sea insesgados.

En el caso de probabilidades iguales y muestreo sin reposición:

$$V(\hat{X}_r) = N^2(1-f) \cdot \frac{S^2}{m} = N^2(1-f) \cdot \frac{S^2}{n/k} = kN^2(1-f) \frac{S^2}{n} = k(1-f)V(\hat{X})$$

luego

$$\hat{V}(\hat{X}_r) = \frac{1-f}{k(k-1)} \sum_{r=1}^k (\hat{X}_r - \hat{X})^2 \quad \text{es un estim. insesgado de } V(\hat{X})$$

La precisión será mayor cuanto mayor sea el n° de grupos aleatorios

3. MÉTODO de los CONGLOMERADOS ÚLTIMOS

El método de los conglomerados últimos, estudiado en el tema anterior, es un caso particular del método de los grupos aleatorios. Basta hacer $k=n$, es decir tomar cada una de las n observaciones como un grupo aleatorio de tamaño $m=1$.

$$\hat{V}(\hat{\theta}) = \frac{1}{n(n-1)} \sum_{r=1}^n (\hat{\theta}_r - \hat{\theta})^2 \quad \text{estim. insesgado de } V(\hat{\theta})$$

$$\text{siendo } \hat{\theta} = \frac{1}{n} \sum_{r=1}^n \hat{\theta}_r$$

$$V(\hat{\theta}_r) = nV(\hat{\theta}) \rightarrow \text{se cumple en muestr. con repos.}$$

$$E_W(\hat{\theta}_r) = \hat{\theta} \rightarrow \hat{\theta}_r \text{ copia de } \hat{\theta} \text{ para } W_r.$$

$\hookrightarrow E[\hat{\theta}_r] = \hat{\theta}$

Este mt. es más preciso que el de grupos aleatorios, porque $n > k$.

Significa que es preferible utilizar n estimadores $\hat{\theta}_i$ individuales que K estimadores $\hat{\theta}_r$ basados en m unidades cada uno.

El mt. de los conglomerados último debe su nombre a su aplicación en muestreo polietápico, para considerarlo como un caso particular del muestreo monoeetápico de conglomer.

Se denomina conglomer. último al conjunto de unidades muestrales de ~~la~~ última etapa seleccionadas en la misma unidad primaria.

Cada conglomerado ~~último~~ ^{primera etapa} es un grupo aleatorio, y a partir de su conglomerado último se construye un estimador insesgado $\hat{\theta}_i$ para θ , de modo que el estimador insesgado de θ basado en la muestra completa es:

$$\hat{\hat{\theta}} = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_i \quad \text{insesgado de } \theta = \frac{1}{N} \sum_{i=1}^N \theta_i$$

En el caso de muestreo polietápico con reposición en 1^a etapa y probabilidades desiguales, el estimador insesgado del total poblacional X es:

$$\hat{\hat{X}}_{HH} = \sum_{i=1}^n \frac{\hat{x}_i}{n p_i} = \frac{1}{n} \sum_{i=1}^n \frac{\hat{x}_i}{p_i} = \frac{1}{n} \sum_{i=1}^n \hat{x}_{(i)}$$

donde $\hat{x}_i \rightarrow$ estim. insesgado de X_i

$\hat{x}_{(i)} \rightarrow$ estim. insesgado de X basado en el conglomer. último i -ésimo.

$\hat{x}_{(i)}$ son indep. por serlo las muestras de las que proceden

$$\hat{V}(\hat{\hat{X}}_{HH}) = \frac{1}{n(n-1)} \sum_{i=1}^n \left(\frac{\hat{x}_i}{p_i} - \hat{\hat{X}}_{HH} \right)^2 \quad \text{es insesgado de } V(\hat{\hat{X}}_{HH})$$

por ser:

$$V(\hat{\hat{X}}_{HH}) = V\left(\sum_{i=1}^n \frac{\hat{x}_i}{n p_i}\right) = \frac{1}{n^2} \sum_{i=1}^n V\left(\frac{\hat{x}_i}{p_i}\right) = \frac{1}{n} V(X_{(i)}) \Rightarrow V(X_{(i)}) = n V(\hat{\hat{X}}_{HH})$$

$$E_W(\hat{x}_{(i)}) = E_W\left(\frac{\hat{x}_i}{p_i}\right) = \sum_{i=1}^n \frac{\hat{x}_i}{p_i} \cdot p(\hat{x}_i) = \sum_{i=1}^n \frac{\hat{x}_i}{p_i} \cdot \frac{1}{n} = \hat{\hat{X}}_{HH}$$

Aunque este método se utiliza siempre para muestreo con reposición en primera etapa, se puede aplicar de forma aprox. al muestreo SIN reposición, multiplicando por $(1-f)$.

4. MT. de las SEMIMUESTRAS REITERADAS

Supongamos una muestra completa W de tamaño n , de la cual extraemos una submuestra aleatoria de tamaño $n/2$ (op. n par), denominada semimuestra.

Reponiendo la semimuestra, repetimos de manera indep. K veces la selección \rightarrow se obtienen K semimuestras reiteradas de W , y K estimadores insesgados $\hat{\theta}_r$ de θ .

No se trata de "grupos" como en el mt. grupos aleatorios, pues la unión de las semimuestras reiteradas no coincide con la muestra completa.

~~Si se cumplen las condiciones:~~

Si se verifican las condiciones:

$$(1) E_W(\hat{\theta}_r) = \hat{\theta} \quad (\Rightarrow E(\hat{\theta}_r) = \theta)$$

$$(2) V(\hat{\theta}_r) = 2V(\hat{\theta})$$

entonces

$$\hat{V}(\hat{\theta}) = \frac{1}{K} \sum_{r=1}^K (\hat{\theta}_r - \hat{\theta})^2 \text{ es un estim. insesgado de } V(\hat{\theta}).$$

- La 2ª condición es obvia en muestreo con reposición, puesto que cada reiteración es de tamaño $n/2$.

- La 1ª condición tb es inmediata, por ser cada reiteración una muestra aleatoria de la muestra completa W .

Para ver la insesgadez del estimador de la variancia:

$$V(\hat{\theta}_r) = E(\hat{\theta}_r - \theta)^2 = E(\hat{\theta}_r - \hat{\theta} + \hat{\theta} - \theta)^2 = E(\hat{\theta}_r - \hat{\theta})^2 + E(\hat{\theta} - \theta)^2$$

$$\text{porque } E(\hat{\theta}_r - \hat{\theta})(\hat{\theta} - \theta) = E(\hat{\theta}_r \hat{\theta}) - E(\hat{\theta}_r \theta) - E(\hat{\theta}^2) + E(\hat{\theta} \theta) = \\ = E(\hat{\theta}^2) - \theta^2 - E(\hat{\theta}^2) + \theta^2 = 0$$

$$\text{luego } V(\hat{\theta}_r) = E(\hat{\theta}_r - \hat{\theta})^2 + E(\hat{\theta} - \theta)^2 \Rightarrow$$

$$\Rightarrow E(\hat{\theta}_r - \hat{\theta})^2 = V(\hat{\theta}_r) - E(\hat{\theta} - \theta)^2 = V(\hat{\theta}_r) - V(\hat{\theta}) = \\ = 2V(\hat{\theta}) - V(\hat{\theta}) = V(\hat{\theta})$$

Por lo que

$$\sum_{r=1}^K E(\hat{\theta}_r - \hat{\theta})^2 = \sum_{r=1}^K V(\hat{\theta}) \Rightarrow E\left(\sum_{r=1}^K (\hat{\theta}_r - \hat{\theta})^2\right) = K V(\hat{\theta}) \Rightarrow$$

$$\Rightarrow E\left(\frac{1}{K} \sum_{r=1}^K (\hat{\theta}_r - \hat{\theta})^2\right) = V(\hat{\theta})$$

y llegamos a comprobar que el estimador $\frac{1}{K} \sum_{r=1}^K (\hat{\theta}_r - \hat{\theta})^2$ es insesgado de $V(\hat{\theta})$.

- La precisión se puede mejorar tomando las K semimuestras reiteradas SIN reposición (para eliminar repeticiones), y de modo que no aparezcan semimuestras complementarias (unión = W), pues proporcionar la misma información sobre $(\hat{\theta}_r - \hat{\theta})^2$.

En consecuencia, tomando solamente $\frac{1}{2} \binom{n}{n/2}$ semimuestras también se obtiene $V_W(\hat{\theta}_r)$.

5. MT. JACKKNIFE ó de los estimadores herramienta

Jackknife = Jack-knife = navaja multiusos.

Los estimadores herramientales son estimadores de múltiples usos y de fácil manejo.

El mt. fue desarrollado por Quenouille (1949) para eliminar o reducir el sesgo de ciertos estimadores. Su generalización y el nombre se debe a Tukey (1958), que lo desarrolló como método general de estimación.

Se parte de un estimador $\hat{\theta}_n$ sesgado, cuyo valor numérico se obtiene a partir de una muestra completa de tamaño n , $W = (X_1 \dots X_n)$.

A partir de W se consideran n muestras de tamaño $n-1$, suprimiendo para muestra i el dato i -ésimo.

$$\left. \begin{array}{l} W = (X_1 \dots X_n) \longrightarrow \hat{\theta}_n \\ W_{n-1}^{(1)} = (\square X_2 \dots X_n) \longrightarrow \hat{\theta}_{n-1}^{(1)} \\ W_{n-1}^{(2)} = (X_1 \square \dots X_n) \longrightarrow \hat{\theta}_{n-1}^{(2)} \\ \vdots \\ W_{n-1}^{(n)} = (X_1 X_2 \dots \square) \longrightarrow \hat{\theta}_{n-1}^{(n)} \end{array} \right\} \longrightarrow \bar{\theta}^{(i)} \longrightarrow \hat{\theta}_n$$

A partir de los estimadores obtenidos de cada muestra, se definen los n pseudovalores:

$$\bar{\theta}^{(i)} = n \hat{\theta}_n - (n-1) \hat{\theta}_{n-1}^{(i)}, \quad i=1 \dots n$$

y el estimador herramienta:

$$\bar{\theta}_n = \frac{1}{n} \sum_{i=1}^n \bar{\theta}^{(i)} = n \hat{\theta}_n - \frac{n-1}{n} \sum_{i=1}^n \hat{\theta}_{n-1}^{(i)}$$

$\bar{\theta}_n$ = media muestral de los pseudovalores muestrales

De este modo, si el sesgo de $\hat{\theta}_n$ depende del tamaño de la muestra y es del tipo $\frac{b}{n}$ ó $\frac{b_1}{n} + \frac{b_2}{n^2} + \dots$, puede eliminarse o al menos reducirse en orden, como resultado de la aplicación de la herramienta. En efecto:

$$\begin{aligned} E[\hat{\theta}_n] &= \theta + \frac{b}{n} \\ E[\hat{\theta}_{n-1}^{(i)}] &= \theta + \frac{b}{n-1} \end{aligned} \quad \left\{ \begin{aligned} E(\bar{\theta}_n) &= n E[\hat{\theta}_n] - \frac{n-1}{n} \sum_{i=1}^n E[\hat{\theta}_{n-1}^{(i)}] = \\ &= n\left(\theta + \frac{b}{n}\right) - \frac{n-1}{n} \cdot n\left(\theta + \frac{b}{n-1}\right) = \\ &= n\theta + b - (n-1)\theta - b = \theta \end{aligned} \right.$$

Una estimación de la varianza de $\bar{\theta}_n$ es:

$$\hat{V}(\bar{\theta}_n) = \frac{n-1}{n} \sum_{i=1}^n (\bar{\theta}_{n-1}^{(i)} - \bar{\theta}_n)^2$$

- Las expresiones anteriores pueden generalizarse suprimiendo dos o más elementos consecutivos en la muestra inicial.
- El mt. Jackknife tb se puede utilizar para mejorar la varianza muestral como estimador de la var. poblacional ó para mejorar los estimadores de razón, que son sesgados.

6 - Métodos BOOTSTRAP ó de autogeneración

Bootstrap \equiv correa que sirve para ayudar a ponerse las botas.

El mt. de autogeneración (bootstrap) se emplea para la estimación aproximada de riesgos, precisiones, intervalos de confianza, regresiones, etc., generalmente a partir de los datos de una sola muestra.

Se debe a Efron (1982) y se utiliza cuando se desconoce la distribución poblacional.

A partir de una muestra aleatoria con reposición de tamaño n , $W = (X_1 \dots X_n)$, donde X_i se consideran v.a.i.i.d, se construye la función de distrib. empírica de la muestra F_n , se obtiene $\hat{\Theta}(F_n)$, estimador del parámetro Θ con distrib. de probabilidad desconocida.

Como F_n asigna la frecuencia $\frac{1}{n}$ a cada valor muestral, podemos considerar la muestra W como una población donde $X_1 \dots X_n$ son los valores que toma la variable con probab $\frac{1}{n}$, de modo que $\hat{\Theta}(F_n)$ es su parámetro.

A partir de la muestra inicial W se toma una muestra de tamaño n con reposición $W^* = (X_1^* \dots X_n^*)$ y se obtiene la función de distrib. empírica autogenerada F_n^* y el estimador $\hat{\Theta}^*$.

Este proceso se repite, de manera independiente, un gran n.º de veces, $M \rightarrow \hat{\Theta}_1^* \dots \hat{\Theta}_M^*$

El estimador bootstrap se obtiene como promedio de los estimadores de las muestras autogeneradas:

$$\hat{\theta}_{\text{BOOT}} = \frac{1}{M} \sum_{j=1}^M \hat{\theta}_j^*$$

$$V(\hat{\theta}_{\text{BOOT}}) = \frac{1}{M} \sum_{j=1}^M (\hat{\theta}_j^* - \hat{\theta}_{\text{BOOT}})^2$$

que se puede estimar como:

$$\hat{V}(\hat{\theta}_{\text{BOOT}}) = \frac{1}{M-1} \sum_{j=1}^M (\hat{\theta}_j^* - \hat{\theta}_{\text{BOOT}})^2 = \frac{1}{M-1} \left[\sum_{j=1}^M \hat{\theta}_j^{*2} - M \left(\sum_{j=1}^M \frac{\hat{\theta}_j^*}{M} \right)^2 \right].$$

Los experimentos efectuados por simulación y por ut. computación, ponen de manifiesto que la autogeneración presenta propiedades que la hacen deseable en cuanto a precisión y rapidez de cálculo, pero da malos resultados para un pequeño porcentaje de muestras posibles.

Los estimadores bootstrap, al igual que los estimadores jackknife, se basan en la obtención de datos ficticios a partir de datos originales, y estiman la variabilidad de un estimador basándose en su variabilidad sobre los conjuntos de datos ficticios, por lo que dan además una idea de la distrib. en el muestreo del estadístico en estudio.