

MUEST_T14. ALGUNAS TÉCNICAS ESPECIALES de MUESTREO.

MUESTREO DOBLE O BIFÁSICO.

MODELOS de CAPTURA - RECAPTURA.

MUESTREO por CUOTAS.

ESTIMACIONES en ÁREAS PEQUEÑAS.

1. MUESTREO DOBLE O BIFÁSICO

En muchas ocasiones se consiguen mejores estimaciones de la variable en estudio X si disponemos de información anticipada de una variable auxiliar Y correlacionada con X , como es el caso del muestreo estratificado y de los estimadores de razón y de regresión.

Cuando falta esta información sobre la variable auxiliar, puede ser relativamente económico tomar una muestra preliminar barata para medir la Y y hacer estimaciones sobre Y , para luego usarlas en las estimaciones de X , más precisas y más costosas.

El muestreo doble selecciona en una 1ª fase una muestra relativamente grande, en la que a bajo coste pueden observarse una o varias características auxiliares de las unidades. En una segunda fase seleccionamos una submuestra de la primera en la que observamos la característica objeto de estimación \rightarrow muestreo en dos fases, muestreo doble o muestreo bifásico.

Aspectos de notación:

1ª fase \rightarrow muestra grande de tamaño n' para estudiar Y_i a bajo coste.

2ª fase \rightarrow muestra de tamaño n ($n < n'$, generalmente submuestra) para observar X_i con coste mayor.

Esevidente fue la conveniencia de esta técnica depende de los costes.

Si la observación de la característica en estudio X_i no tiene coste, tomaríamos una muestra del tamaño n_0 necesario para la precisión deseada de las estimaciones.

Sp. que disponemos de un presupuesto total C :

~~coste unit. 1ª muestra =~~

muestra coste unit. tamaño muestral

1	c'	n'
2	c	n

Nos planteamos dos opciones:

- 1ª fase muestra $\rightarrow C = C_0 \cdot n_0$
- Dos fases $\rightarrow C = c'n' + cn$

Iguando los costes totales, $C_0 n_0 = c'n' + cn$

$$n_0 = \frac{c'n' + cn}{C_0}$$

por lo que con la técnica en dos fases la observación de X se hace en una muestra de tamaño n , menor que n_0 , con el mismo coste total \Rightarrow al introducir las dos fases el tamaño de la muestra necesario es más pequeño que si hubiera una fase (muestreo aleatorio normal) y hay una pérdida en la precisión de los estimadores.

Para ver si compensa la disminución del tamaño efectivo de la muestra con el incremento de información adquirido en la 1ª fase (\Rightarrow pérdida de precisión en la estim. de X) calculamos la variancia:

- 1 fase: σ^2/n_0 (para la media)
- 2 fases: $n_0 - n = (\frac{c'}{c})n' \Rightarrow$ cuanto menor sea $\frac{c'}{c}$ más favorable es el muestreo doble (más cerca estará n de n_0 y menor será la pérdida de precisión).

El muestreo bifásico sólo es bueno si lo que se gana en precisión con la variable Y compensa con la pérdida en precisión debido a la reducción del tamaño de la muestra.

Muestreo doble para estratificación

Supongamos que la población se estratifica en L estratos.

1ª fase: muestra aleatoria de tamaño n'

$$W_h = \frac{\text{nº elem. poblac. estrato } h}{\text{nº total elem. poblac.}} = \frac{N_h}{N}$$

$$\hat{W}_h = \frac{\text{nº elem. 1ª muestra estrato } h}{\text{nº total elem. 1ª muestra}} = \frac{n_h}{n'} \quad \text{i sesgado de } W_h$$

$$\sum_{h=1}^L n'_h = n' \quad \text{y} \quad \sum_{h=1}^L n_h = n$$

2ª fase: muestra aleatoria estratificada de tamaño n

De cada estrato se elige una submuestra de tamaño $n_h \leq n'_h$ indep.

El estimador usual de la media es $\hat{\bar{X}} = \sum_{h=1}^L W_h \bar{x}_h$,

en muestreo doble es $\hat{\bar{X}} = \sum_{h=1}^L \hat{W}_h \bar{x}_h$ con $\hat{W}_h = \frac{n'_h}{n'}$ y $\bar{x}_h = \frac{x_h}{n_h}$

que es i sesgado para la media.

Para calcular su variancia, se acude al teo de Madon:

$$V(\hat{\bar{X}}) = V(E_{W'}(\hat{\bar{X}})) + E(V_{W'}(\hat{\bar{X}}))$$

la expresión de la variancia del estimador, distinta para el caso SR y CR es tal que n' (tamaño muestra 1ª fase) aparece en el denominador, de modo que cuanto mayor sea n' , la pérdida de precisión por el uso de muestreo doble disminuye. (Como el coste aumenta, valorar los tamaños y la afijación óptimos en función del coste).

Muestreo doble para estimadores de razón

El estimador usual de razón para la media \bar{X} utiliz como información conocida \bar{Y} o Y .

El muestreo doble utiliza

1=fase: muestra de tamaño n' para obtener una buena estimación de \bar{Y} o de Y .

2=fase: muestra de tamaño n para obtener \bar{x} e \bar{y} .

El estimador de razón para el muestreo doble es:

$$\hat{\bar{X}}_R = \frac{\bar{x}}{\bar{y}} \cdot \bar{y}', \quad \text{donde } \bar{y}' \equiv \text{media de la 1ª muestra}$$

El estimador será insesgado si lo es \hat{R} .

Se puede dar expresiones aproximadas de la variancia del estimador (variancia y covariancia para CR y cuasivariancias y cuasicovariancias para SR).

La variancia es igual a la variancia del estim. de razón más una penalización por utilizar muestreo doble (cuanto mayor sea n' , la ~~precisión~~ penalización es más pequeña).

* En estimadores de regresión ocurre lo mismo.

2 - MODELOS DE CAPTURA - RECAPTURA

El mt. de captura y recaptura o de captura-marcado-recaptura y recaptura, consiste en tomar una primera muestra, marcar las unidades capturadas (procurando que no modifique su característica ni condiciones de comportamiento), y ponerlos en libertad en la zona que fueron aprehendidos. Tras un período establecido no muy largo, se toma una segunda muestra y se observan cuántos elementos están marcados. El tamaño de la poblac. se estima por:

$$\hat{N} = \frac{n_1 \cdot n_2}{m_2} \quad (\text{basado en la distib. hipergeométrica})$$

donde $n_i \equiv$ tamaño muestra i ($i=1,2$).

$m_2 \equiv$ número de recapturas

la probab. de obtener en una muestra de tamaño ~~n_2~~ n_2 , m elementos marcados, con n_1 elementos marcados en una poblac. de tamaño N (desconoc) es:

$$P(X=m) = \frac{\binom{n_1}{m} \cdot \binom{N-n_1}{n_2-m}}{\binom{N}{n_2}} = \frac{n_1! n_2!}{m! (n_1+n_2-m)!}$$

$N \equiv n_1 + n_2 - m$

Para estimar N puede utilizarse el mt. de máxima verosimilitud, $\hat{N} = \frac{n_1 n_2}{m}$ es un estimador máx. verosímil de N que se apoya en el supuesto de que la proporción de recapturas en la segunda muestra es igual a la prop. de individuos marcados en la población: $\frac{n_1}{N} = \frac{m_2}{n_2}$

Este supuesto se basa en:

- la poblac. es cerrada en el período de estudio, N cte.
- Los individuos marcados no pueden afectar al muestreo, y se distribuyen de manera aleatoria, con lo que la probab. de ser capturados en la 2ª muestra es constante.

Estos supuestos no son siempre plausibles, por lo que se han propuesto otros tipos de estimación para evitar sesgos basándose en supuestos sobre la población de origen. Deben mencionarse los mt. de eliminación (removal) y de capturas múltiples.

En ocasiones han de estimarse algunos totales, peso, volumen, alimento consumido, etc. lo que requiere estimadores del tipo : $\hat{N} \cdot \bar{X}$, donde ambos factores son variables aleatorias (estimadores del producto).

El mt. de captura y recaptura se aplica a la estimación de poblaciones móviles humanas, en cuyo caso se recomienda ponerse en contacto con los jefes tribales o superiores jerárquicos.

3 - MUESTREO POR COTAS

El muestreo por cuotas, desarrollado en los años 30 por Chenington, Roper, Gallup y Gessley, ha sido adoptado por muchas organizaciones dedicadas a realizar encuestas sobre opinión pública, estudios de mercado, etc.

No existe una definición precisa de muestreo por cuotas, es su punto débil: existen tantas variantes como aplicaciones.

Sp. fue el diseño de la encuesta lo seguido los principios del muestreo probabilístico hasta llegar al momento de seleccionar las personas a entrevistar. En esta etapa se le impone al entrevistador que realice un determinado nº de entrevistas por edad, sexo, nivel económico... y c. de característica sociológica o socioeconómica de interés, y se le deja libertad plena para elegir las personas que cumplan estos requisitos \Rightarrow sesgos que no pueden ser detectados.

Por otro lado, el desconocimiento de las probab. de selección no permite evitar los errores por ponderaciones incorrectas en el proceso de estimación y no se pueden estimar los errores debidos al muestreo.

\Rightarrow No se puede comparar el coste del muestreo por cuotas con el coste por muestreo probabilístico.

Especialmente en encuestas de opinión, la muestra puede representar muy bien a la población para una característica y muy mal para opiniones.

Una muestra por cuotas puede proporcionar estim. muy útiles, pero la dificultad está en conocer a priori hasta qué punto lo son.

En encuestas importantes suele exigirse muestreo probabilístico.

4- ESTIMACIONES en ÁREAS PEQUEÑAS.

En muchas encuestas, la población se subdivide en clases para las que se requiere estimaciones separadas.

Dominió \rightarrow cualquier subpoblación para la que se necesita una estim. separada, antes o después de diseñar el muestreo.

Si la necesidad de estimaciones separadas se conoce antes de hacer el diseño, se puede plantear los dominios como estratos \rightarrow coste adicional, depende del u^2 de dominio.

Si las estimaciones a posteriori se necesitan a posteriori (por desconocimiento, economía...) las unid. muestrales pueden pertenecer al dominio o no, por lo que el u^2 de observaciones muestrales pertenecientes al dominio

- es una v.a.

- puede ser muy pequeño, o nulo.

Surge el problema de estimación en dominios pequeños, ya que los estimadores usuales serán menos precisos al estar basados en menos observaciones, y si no existen observaciones del dominio, no se podrán utilizar.

El término estimación en pequeñas áreas se utiliza cuando los dominios pequeños corresponden a zonas geográficas. Hay varias alternativas:

\rightarrow Estimadores directos:

- Utilizan datos únicamente de las unidades muestrales pertenecientes al dominio de estudio

- El tamaño de la muestra es una v.a.

- El estimador tiene mayor varianza.

Ejemplo: En m.a.s.s.r.j. sea $y_j \equiv u^2$ unid. muestrales dominio (dec).

$$\bar{x}_j = \sum_{k=1}^{N_j} x_{jk} \rightarrow \hat{\bar{x}}_j = \frac{1}{n} \sum_{k=1}^{N_j} x_{jk} \quad \text{Considerando } x'_{jk} = \begin{cases} x_{jk} & \text{si } e \text{ dom } j \\ 0 & \text{si } e \notin \text{dom } j \end{cases}$$

$$X' = \sum_{i=1}^N x'_{i.} = y_j \Rightarrow N\bar{x}' = N \cdot \frac{1}{n} \sum_{i=1}^N x'_{i.} = \bar{x}_j$$

$$V(\bar{x}_j) = (1-f) \cdot \frac{S^2}{n}, \text{ donde } S^2 = \frac{1}{N-1} \sum_{i=1}^N (x'_{i.} - \bar{x}')^2$$

→ Estimadores indirectos (simuléticos)

- utilizan datos de unid. de otro dominio
- Gran variedad en función de los modelos que engende la poblar. finita de estudio
- Estimadores sesgados → varianza no es buena medida de la bondad de los estimadores

Ejemplo: En m.a.s.s.r. → $X \equiv \text{var. estudio}$
 $Y = RY$ $Y \equiv \text{var. correlacionada con } X$ $\hat{R} = \frac{\bar{X}}{\bar{Y}}$
 $\hat{X}_j = \hat{R} \hat{Y}_j$

Consideramos que lo que se verifica en el modelo poblacional tb. se verifica en el dominio. El sesgo será mayor cuanto más falte esta suposición.

→ Estimadores combinados:

- Media ponderada de un estim. directo (basado en el diseño) y un estim. indirecto (basado en el modelo)
- Intentan equilibrar el sesgo potencial del estimador simulético y la varianza del estimador directo.

$$\hat{Y}_d = \alpha \hat{Y}_d^{\text{directo}} + (1-\alpha) \hat{Y}_d^{\text{simulético}} \quad / \alpha \in (0,1)$$