

Muestreo Probabilístico

Giovany Babativa, PhD

Sobre Mi

PhD en Estadística, MSc en Big Data, MSc en Estadística. Como años de experiencia, actual director de analítica en el CNC, miempo del comité de expertos en pobreza en el DANE y consultor de la División de Estadística de la CEPAL. Ex-decano de la Facultad de Estadística USTA, ex-director de operaciones en el ICFES,...

Puedes encontrarme en:

- Soogle scholar
- GitHub. https://github.com/jgbabativam
- in linkedin
- **≥** j.babativamarquez@uniandes.edu.co

2

su calidad y su sentido



MUESTREO PROBABILÍSTICO

Notación



Defina a U un universo¹ de elementos $\{U_1,\ldots,U_N\}$ finito y conocido de antemano con una variable de interés Y que toma valores $\{y_1,\ldots,y_N\}$. Sea el parámetro θ (medida del universo) una función de (y_1,\ldots,y_N) de esta manera a $\theta(y_1,\ldots,y_N)$ se denomina parámetro y se denota θ .

1 En adelante se denominará universo a la población objetivo

Conceptos básicos



Algunos parámetros de interés en un estudio por muestreo:

- ullet Total: $t_y = \sum_U y_k$
- Promedio: $ar{y}_U = rac{1}{N} \sum_U y_k$
- Razón: $R = rac{\sum_{U} y_k}{\sum_{U} z_k} = rac{t_y}{t_z}$

Probabilidades de Inclusión



Se define probabilidad de inclusión de primer orden del elemento k

$$\pi_k = \sum_{k \in s_i} p(s_i)$$

Sea:

$$I_k = \left\{ egin{array}{l} 1 ext{ si } k \in s \ 0 ext{ en otro caso} \end{array}
ight.$$

Entonces $\pi_k = P(I_k = 1)$

Probabilidades de Inclusión



Se define probabilidad de inclusión de segundo orden de los elementos $k \ \mathbf{y} \ l$

$$\pi_{kl} = \sum_{k,l \in s_i} p(s_i)$$

Entonces
$$\pi_{k,l} = P(I_k I_l = 1)$$

Probabilidades de Inclusión



- \bullet $\pi_{kk}=\pi_k$.
- Por definición de muestra probabilística $\pi_k > 0$.
- En muestreo de elementos π_k para todo $k=1,\ldots,N$ son conocidos de antemano.

Ejercicio





- Use el espacio muestral del diseño MAS para calcular la probabilidad de inclusión del elemento 1, verifique que $\pi_k = \frac{n}{N}$.
- Use el espacio muestral del diseño BER para calcular la probabilidad de inclusión del elemento 1, verifique que $\pi_k=\pi$.

10:00

Estadística y Estimador



Sea $\hat{\theta}$ una estadística o estimador entonces bajo el diseño mues $p(\cdot)$ se define:

- ullet Valor esperado: $E_P(\hat{ heta}) = \sum_S p(s_i) \hat{ heta}$
- ullet Sesgo: $B(\hat{ heta}) = E_P(\hat{ heta}) heta$
- ullet Varianza: $V_P(\hat{ heta}) = \sum_S p(s_i) \Big [\hat{ heta} E(\hat{ heta}) \Big]^2$
- ullet Error Cuadrático Medio: $ECM(\hat{ heta}) = V(\hat{ heta}) + B^2(\hat{ heta})$

Estadística y Estimador



Sea $\hat{ heta}$ una estadística o estimador entonces bajo el diseño mues $p(\cdot)$ se define:

- ullet Error estándar: $\sqrt{V_P(\hat{ heta})}$
- Coeficiente de variación (error relativo): $\frac{\sqrt{V_P(\hat{\theta})}}{E_P(\hat{\theta})}$ Coeficiente de variación estimado: $cve(\%)=100*\frac{\sqrt{\hat{V}_P(\hat{\theta})}}{\hat{\theta}}$

Ejercicio

Tu turno

Calcule $E_P(\hat{ heta}) = \sum_S p(s_i) \hat{ heta}$ para:

$$\hat{ heta} = rac{1}{n} \sum_s y_k$$

Use el espacio muestral de un MAS

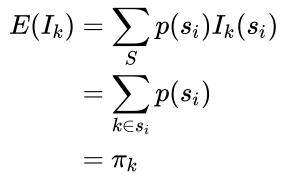


10:00

Estadística I_k

Ejemplo:





Estadística I_k

UNIVERSIDAD EL BOSQUE Por una cultura de la vida, su calidad y su sentido

Ejemplo:

$$egin{aligned} V(I_k) &= \sum_S p(s_i) [I_k(s_i) - E(I_k(s_i))]^2 \ &= \sum_S p(s_i) I_k^2(s_i) - 2 \sum_S p(s_i) I_k(s_i) \pi_k + \sum_S p(s_i) \pi_k^2 \ &= \pi_k - 2 \pi_k^2 + \pi_k^2 \ &= \pi_k (1 - \pi_k) \end{aligned}$$

- ejemplo MAS
- ejemplo BER

Estadística I_k

Ejemplo:



$$egin{aligned} Cov(I_k,I_l) &= \sum_{S} p(s_i) \left(I_k(s_i) - \pi_k
ight) \left(I_l(s_i) - \pi_l
ight) \\ &= \sum_{S} p(s_i) I_k(s_i) I_l(s_i) - \pi_l \sum_{S} p(s_i) I_k(s_i) \\ &- \pi_k \sum_{S} p(s_i) I_l(s_i) + \pi_k \pi_l \sum_{S} p(s_i) \\ &= \sum_{k,l \in s_i} p(s_i) - \pi_l \sum_{k \in s_i} p(s_i) - \pi_k \sum_{l \in s_i} p(s_i) + \pi_k \pi_l \\ &= \pi_{kl} - \pi_k \pi_l - \pi_k \pi_l + \pi_k \pi_l \\ &= \pi_{kl} - \pi_k \pi_l \\ &= \Delta_{kl} \end{aligned}$$



Estadística n_s

Ejemplo:



$$egin{aligned} E(n_s) &= \sum_S p(s_i) n_s \ &= \sum_S p(s_i) \sum_U I_k \ &= \sum_U \sum_S p(s_i) I_k \ &= \sum_U \sum_{k \in s_i} p(s_i) \ &= \sum_U \pi_k \end{aligned}$$

• ejemplo MAS y BER



Parámetros de interés



Para un universo U de tamaño N, sea y la característica de interentonces podríamos estar interesados en:

- Total: $t_y = \sum_{U} y_k$ (personas con cierta enfermedad)
- ullet Media: $ar{y}_U = rac{\sum_U y_k}{N} = rac{t_y}{N}$ (dinero)
- ullet Proporción: $p_U=rac{\sum_U y_k}{N}=rac{t_y}{N}$ para $y_k=\{1,0\}$ (desplazados)
- Razón: $R = \frac{t_y}{t_z}$. Unidades del producto por establecimiento con la intención de venderlo.

Nótese que todos los parámetros pueden ser expresados como función de totales, por tanto hay un particular interés encontrar estimadores para este parámetro.



Estimador de Horvitz-Thompson (1952)



Para un universo U se desea estimar el total de una característic interés y denotado como t_y . Por ejemplo,

Para
$$heta=t_y=\sum_U y_k$$
 se define:

$$\hat{ heta} = \hat{t}_{\,y,\pi} = \sum_{s_i} rac{y_k}{\pi_k}$$

 $\frac{1}{\pi_k}$ se denomina **Factor de expansión**

Cada elemento se representa a sí mismo y a una fracción de la población.

Estimador de Horvitz-Thompson (1952) - Propiedades

UNIVERSIDAD EL BOSQUE Por una cultura de la vida, su calidad y su sentido

Resultado

- $E(\hat{t}_{y,\pi}) = t_y$ demost.
- $ullet \ V(\hat{t}_{y,\pi}) = \sum \sum_U \Delta_{kl} rac{y_k}{\pi_k} rac{y_l}{\pi_l}$

$$ullet \ \widehat{V}(\widehat{t}_{\,y,\pi}) = \sum \sum_{oldsymbol{s}_i} rac{\Delta_{kl}}{\pi_{kl}} rac{y_k}{\pi_k} rac{y_l}{\pi_l}$$

$$\Delta_{kl} = \pi_{kl} - \pi_k \pi_l$$
, además $E(\widehat{V}(\widehat{t}_{\:y,\pi})) = V(\widehat{t}_{\:y,\pi})$

Estimador de Horvitz-Thompson (1952)

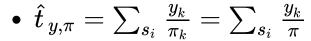
Muestreo Aleatorio Simple

- ullet $\hat{t}_{y,\pi} = \sum_{s_i} rac{y_k}{\pi_k} = rac{N}{n} \sum_{s_i} y_k$
- $ullet V_{MAS}(\hat{t}_{y,\pi}) = rac{N^2}{n} ig(1 rac{n}{N}ig) \, S_{yU}^2$
- $ullet \; \widehat{V}_{MAS}(\hat{t}_{\;y,\pi}) = rac{N^2}{n}ig(1-rac{n}{N}ig)\,S^2_{ys_i}$



Estimador de Horvitz-Thompson (1952)

Muestreo Bernoulli



$$ullet V_{Ber}(\hat{t}_{y,\pi}) = \left(rac{1}{\pi} - 1
ight) \sum_U y_k^2$$

$$ullet ~ \widehat{V}_{Ber}(\hat{t}_{~y,\pi}) = rac{1}{\pi}ig(rac{1}{\pi}-1ig)\sum_{s_i}y_k^2$$

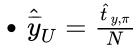
Ejercicio: Para los ejercicios hechos en clase de MAS y Bernoulli calcule $\hat{t}_{y,\pi}$ y $V_{Ber}(\hat{t}_{y,\pi})$, este último vía definición y por el estimador dado por la expresión. Para ambas expresiones compruebe el insesgamiento del estimador.





Estimación de la media poblacional

Muestreo Aleatorio Simple



$$ullet V_{MAS}(\hat{ar{y}}_U) = rac{1}{n} ig(1 - rac{n}{N}ig) \, S_{yU}^2$$

$$ullet \ \widehat{V}_{MAS}(\hat{ar{y}}_U) = rac{1}{n}ig(1-rac{n}{N}ig)\,S^2_{ys_i}$$

Al factor $\left(1-\frac{n}{N}\right)$ se le conoce como factor de corrección para poblaciones finitas.





DISEÑOS CON INFORMACIÓN AUXILIAR

Muestreo Proporcional al Tamaño $PPT(p_k,m)$

Se realizan m eventos aleatorios con probabilidad contante 1/I decir se selecciona uno de los N elementos, se realiza su medicion y_{k_1} , se repone en el universo y se selecciona de nuevo uno de los N elementos con medición y_{k_2} y así sucesivamente hasta completar m elementos. Sea:

- ullet Z: Número de veces que el elemento k aparece en la muestra
- Z = 0, 1, ..., m
- $Z \sim Bin(m, 1/N)$

•
$$P(Z=r)=\binom{m}{r}(1/N)^r(1-1/N)^{m-r}$$

•
$$P(Z=0)=(1-1/N)^m$$

•
$$\pi_k = P(Z \ge 1) = 1 - P(Z = 0) = 1 - (1 - \frac{1}{N})^m$$

Por una cultura de la vida, su calidad y su sentido



Muestreo Proporcional al Tamaño $PPT(p_k,m)$



Una generalización conduce al diseño de Probabilidad Proporcional al Tamaño (PPT) cuando cada elemento se selecciona con probabilidad p_k , tal que $\sum_U p_k = 1$. Para este diseño se debe contar con información auxiliar disponible para todos los elementos de U, este diseño tienes varios beneficios entre los que se encuentra que reduce costos y que su uso es relativamente simple.

Algoritmo de Selección: Acumulativo Total

Ejemplo

Suponga que U=9, para un estudio de mercados se desea seleccionar una muestra de 4 tiendas, pero basado en la información auxiliar se desea darle más peso a las tiendas que tienen una mayor rotación del producto. La información auxiliar para las 9 tiendas es:

Tienda	1	2	3	4	5	6	7	8	9	
Ventas_Unid	32	96	65	140	22	78	65	47	106	•

π -Estimador



La muestra ordenada es de tamaño m mientras que la muestra no ordenada (sin repetición) es de tamaño n.

$$\hat{t}_{\,y,\pi} = \sum_{s_i} rac{y_k}{\pi_k} = \sum_{s_i} rac{y_k}{1-(1-p_k)^m}$$

El estimador cuenta con n sumandos, es decir, que los elementos repetidos solo cuentan una vez dentro del estimador.

Estimador de Hansen-Hurwitz (1943)



$$\hat{t}_{\,y,MCR} = rac{1}{m} \sum_{i=1}^m rac{y_{k_i}}{p_{k_i}}$$

En donde p_{k_i} es la probabilidad de selección del elemento k. Cada elemento se representa así mismo y al resto del universo \to promedio.

Estimador de Hansen-Hurwitz (1943)

$$egin{align} E\left(\hat{t}_{y,MCR}
ight) &= E\left(rac{1}{m}\sum_{i=1}^{m}rac{y_{k_i}}{p_{k_i}}
ight) \ &= rac{1}{m}E\left(\sum_{U}rac{y_k}{p_k}Z_k
ight) \ &= rac{1}{m}\sum_{U}rac{y_k}{p_k}E\left(Z_k
ight) \ &= rac{1}{m}\sum_{U}rac{y_k}{p_k}mp_k \ &= \sum_{U}y_k = t_y \ \end{cases}$$





Estimador de Hansen-Hurwitz (1943)



-
$$V\left(\hat{t}_{y,MCR}
ight) = rac{1}{m} \sum_{U} p_{k} \Big(rac{y_{k}}{p_{k}} - t_{y}\Big)^{2}$$

$$ullet \ \widehat{V}\left(\widehat{t}_{y,MCR}
ight) = rac{1}{m(m-1)} \sum_{i=1}^{m} \left(rac{y_{k_i}}{p_{k_i}} - \widehat{t}_{y,MCR}
ight)^2$$

$$ullet E\left(\widehat{V}\left(\widehat{t}_{y,MCR}
ight)
ight)=V\left(\widehat{t}_{y,MCR}
ight)$$

Ejemplo:



Tu turno

Para el estudio de mercados suponga que $s_o=\{4;7;7;9\};$ $p_k=\{0.215;0.1;0.1;0.163\}$ y que los valores de las ventas semanales recolectados los establecimientos seleccionados fueron $y_k=\{12;10;10;11\}$

- Usando el π —estimador, estime el total de ventas semanales (demanda) del producto en las nueve tiendas.
- Usando el MCR—estimador, estime el total de ventas semanales (demanda) del producto en las nueve tiendas y calcule el cve.

10:00

Muestreo $PPT(p_k,m)$: Notas



- Es fácil comprobar que $V\left(\hat{t}_{y,MCR}\right)$ si $y_k=cp_k,c$ constante. extstyle decir, si y_k es exactamente proporcional a p_k .

En la práctica eso es imposible !!!... primero es la probabilidad y luego la medición.

Muestreo $PPT(p_k,m)$: Notas



- En la práctica, sea x_k una variable auxiliar altamente correlacionada con y_k así $y_k/x_k \doteq c$ entonces se puede determinar

$$p_k = rac{x_k}{\sum_U x_k} = rac{x_k}{t_x}, k = 1, \dots, N$$

Entre más esté correlacionado x_k con y_k entonces $V\left(\hat{t}_{\mathit{u},MCR}\right)
ightarrow 0$

Muestreo $PPT(p_k,m)$: Notas



 $\bullet\,$ Si los p_k no varían casi, me voy con MAS pero si varían mucho me voy con PPT .

Ejercicio



Tu turno

Suponga que se tiene $U=\{1,2,3,4,5\}$ y que se conocen todos los valores de la variable de interés, que están dados por:

Υ	79.0	76.00	54.0	39.00	12.0
pk	0.1	0.15	0.2	0.25	0.3

Para un muestreo con reemplazamiento con m=3, si los valores aleatorios son $\zeta_k=(0.05,058,0.36)$, determine la muestra, calcule la estimación y el cve usando el MCR-estimador. ¿Qué tendría que hacer para determinar si es mejor la estrategia usando π -estimador o MCR-estimador?

12:00

Ejercicio

Tu turno



Resuelva los ejercicios: 3.1.1, 3.3, 3.4, 3.12, 3.14, 3.15 y 3.19 del libro de Gutiérrez, H.A. (2016)



CONFIABILIDAD

Confiabilidad



Sea S el conjunto de todas las muestras posibles s_1,\ldots,s_Q ; para el diseño $P(s_1),\ldots,P(s_Q)$ $\left[\sum_{s_i\in S}p(s_i)=1\right]$; el parámetro θ $\left[\theta:U o\mathbb{R}\right]$ y el estimador $\hat{\theta}(s_i)$.

Confiabilidad

Para cada muestra se construye el intervalo $(LI(s_i), LS(s_i))$ o manera que cada muestra tiene su propio intervalo y la longitud uei intervalo puede variar dependiendo de la muestra.

Se define **confiabilidad** como la suma de las probabilidades de las muestras de las muestras cuyo intervalo de confianza cubre al parámetro:

$$\sum_{ heta \in (LI(s_i), LS(s_i))} P(s_i)$$

Cuando se habla de confiabilidad se habla de la confiabilidad del intervalo NO de la muestra.

su calidad y su sentido

Propuesta de intervalos



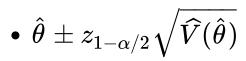
$$ullet$$
 $\hat{ heta} \pm z_{1-lpha/2} \sqrt{V(\hat{ heta})}$, por TLC se tiene que:

$$P\left(heta\in\hat{ heta}\pm z_{1-lpha/2}\sqrt{V(\hat{ heta})}
ight)=1-lpha,$$
 si $E(\hat{ heta})= heta.$

Este intervalo presenta dos inconvenientes:

- Condición de estimador insesgado
- Se necesita la varianza del estimador y esa sólo se conoce en la teoría porque en la práctica solo se conoce cuando es censo y en ese caso no tiene mucho interés.

Propuesta de Intervalos





Este intervalo presenta la cualidad de que en la práctica puede calcularse, lo que interesa es entonces es ¿Qué confiabilidad tiene?.

Puesto que los estimadores vistos son insesgados tanto para el parámetro de interés como para la varianza se espera que para $z_{1-\alpha/2}=1.96$ en un 95% de las muestras el intervalo contenga el parámetro. Pero lo anterior no me garantiza que el parámetro esté en el intervalo de la muestra seleccionada.

Ejercicio

Tu turno



Calcule la confiabilidad para los dos tipos de intervalo usando los datos del gasto recogidos en clase. Hágalo para el caso MAS, de tarea en casa resuelva el diseño BER.

20:00

Sesgo Relativo

La confiabilidad es una medida que se ve directamente afectac el sesgo 1 NO por el margen de error. Es decir, si $E(\hat{\theta}) \neq \sigma$ la cobertura del intervalo decrece y por más de que $z_{1-\alpha/2}=1.96$ la cobertura puede ser muy baja.

$$B_r(\hat{ heta}) = rac{B(\hat{ heta})}{\sqrt{V(\hat{ heta})}} = rac{E(\hat{ heta}) - heta}{\sqrt{V(\hat{ heta})}}$$

1 medibles v no medibles

Sesgo Relativo



Tu turno

Para el ejercicio de clase INCLUYA sesgos de 0;0.05;0.1;0.3;0.5;1;1.5;2 y usando un nivel de confianza del 95% calcule y grafique la confiabilidad real. Use:

$$E(\hat{ heta}) - heta = B_r(\hat{ heta}) \sqrt{V(\hat{ heta})}$$

10:00



GRACIAS!



Referencias

- Gutiérrez, H. A. (2016). Estrategias de muestreo: Diseño de encuestas y estimación de parámetros. Ediciones de la U.
- Lohr, S. L. (2021). Sampling: design and analysis. Chapman and Hall/CRC.
- Särndal, C. E., Swensson, B., & Wretman, J. (2003). Model assisted survey sampling. Springer Science & Business Media.
- Valliant, R., Dever, J. A., & Kreuter, F. (2013). Practical tools for designing and weighting survey samples (Vol. 1). New York: Springer.



Citación y derechos de autor

Este material ha sido creado por Giovany Babativa-Márquez y es de libre distribución bajo la licencia Creative Commons Attribution-ShareAlike 4.0.

Si se copia parcial o totalmente, debe citar la fuente como:

Babativa-Márquez, J.G. *Diapositivas del curso de muestreo probabilístico*. URL: https://jgbabativam.github.io/Muestreo-I/Semana2.html. 2024.