

Estimador asistido por modelos

Muestreo II

Licenciatura en Estadística

2023

set-up para la estimación

Objetivo: construir estimaciones para cantidades poblacionales

- ▶ utilizamos datos recolectados utilizando (generalmente) diseños muestrales complejos
- ▶ también tenemos disponible otras fuentes de datos \mathbf{x} (e.g. censos, registros, etc.)
- ▶ combinando las dos fuentes anteriores de datos, podemos aumentar la **eficiencia** de nuestras estimaciones!

set-up para la estimación

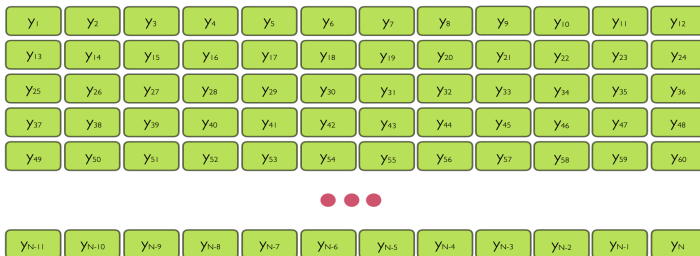
Universo/Elegibles (e.g. hogares de Montevideo)



$$U = \{1, 2, \dots, i, \dots N\}$$

set-up para la estimación

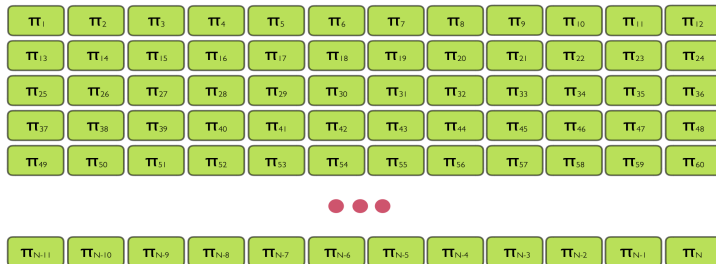
Objetivo: estimar el total de la variable de interés (target) y



$$Y = \sum_{i=1}^N y_i = \sum_{i \in U} y_i$$

set-up para la estimación

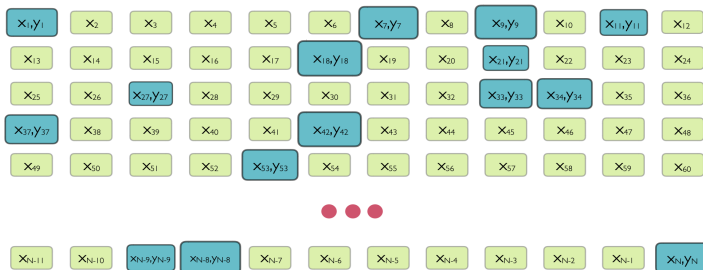
Definimos un diseño muestral en base a los datos disponibles \mathbf{x} (e.g. estratos, MOS, etc.) de forma de asignarle a cada individuo una probabilidad de selección en la muestra (s)



$$\pi_i = \text{Prob}[i \in s] > 0 \quad \forall i \in U$$

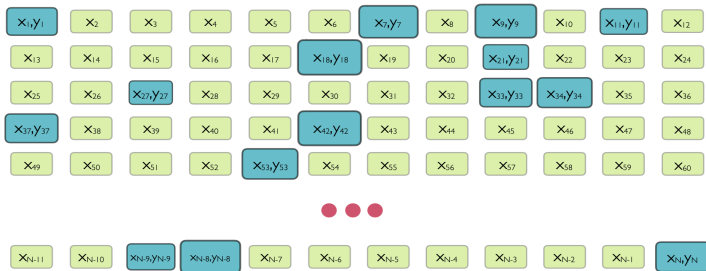
set-up para la estimación

- ▶ se selecciona la muestra
- ▶ se recolecta información de la variable de interés y e información adicional (\mathbf{x})



estimador HT

los estimadores simples utilizan únicamente información de la variable de interés y de las unidades incluidas en la muestra s



$$\hat{Y}^{\text{HT}} = \sum_{i \in s} \frac{y_i}{\pi_i} = \sum_{i \in s} w_i \times y_i$$

propiedades

buenas propiedades teóricas

- ▶ insesgado $\rightarrow E(\hat{Y}^{HT}) = Y$
- ▶ varianza "simple"

$$\text{var}(\hat{Y}^{HT}) = \sum_{i \in U} \sum_{j \in U} (\pi_{ij} - \pi_i \pi_j) (y_i \pi_i^{-1}) (y_j \pi_j^{-1})$$

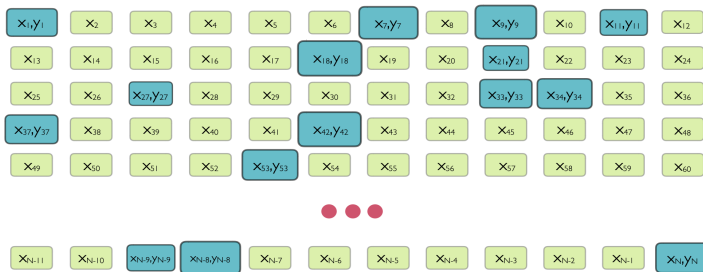
donde $\pi_{ij} = \text{Prob}(i \in s \& j \in s)$

- ▶ un estimador insesgado de la varianza (si el diseño lo permite)

$$\widehat{\text{var}}(\hat{Y}^{HT}) = \sum_{i \in s} \sum_{j \in s} \pi_{ij}^{-1} (\pi_{ij} - \pi_i \pi_j) (y_i \pi_i^{-1}) (y_j \pi_j^{-1})$$

estimador de diferencia

supongamos que tenemos una forma o "método" $m(\cdot)$ para poder predecir la variable y , el cual, no depende de la muestra



con las "salidas" del método $m(\cdot)$ se define el estimador de diferencia:

$$\hat{Y}^{\text{DIFF}} = \sum_{i \in U} m(\mathbf{x}_i) + \sum_{i \in S} w_i (y_i - m(\mathbf{x}_i)) = \sum_{i \in U} m(\mathbf{x}_i) + \text{HT}(y - m)$$

propiedades (1)

el estimador es exactamente insesgado independientemente de la calidad del método $m(\cdot)$

$$E[\hat{Y}^{\text{DIFF}}] = \sum_{i \in U} m(\mathbf{x}_i) + E[HT(y - m)] = \sum_{i \in U} m(\mathbf{x}_i) + Y - \sum_{i \in U} m(\mathbf{x}_i) = Y$$

propiedades (2)

Teniendo en cuenta que $\sum_{i \in U} m(\mathbf{x}_i)$ no depende de la muestra (i.e. NO es aleatorio) y la varianza en el diseño es:

$$\text{var}(\hat{Y}^{\text{DIFF}}) = \sum_{i \in U} \sum_{j \in U} (\pi_{ij} - \pi_i \pi_j) \frac{(y_i - m(\mathbf{x}_i))}{\pi_i} \frac{(y_j - m(\mathbf{x}_j))}{\pi_j}$$

- ▶ la varianza del estimador será mas pequeña que la del estimador HT si los "residuos" $(y_i - m(\mathbf{x}_i))$ tiene una menor variación que los datos puros y_i
- ▶ si el método elegido $m(\cdot)$ tiene un pobre poder predictivo la varianza del estimador de diferencia será similar a la varianza del estimador HT

propiedades (3)

Un estimador insesgado de la varianza es:

$$\widehat{\text{var}}(\hat{Y}^{\text{DIFF}}) = \sum_{i \in s} \sum_{j \in s} \pi_{ij}^{-1} (\pi_{ij} - \pi_i \pi_j) \frac{(y_i - m(\mathbf{x}_i))}{\pi_i} \frac{(y_j - m(\mathbf{x}_j))}{\pi_j}$$

conclusiones

- ▶ si bien el estimador \hat{Y}^{DIFF} resulta atractivo, el problema se encuentra que en la práctica es difícil tener un método $m(\cdot)$ que sea independiente de la muestra y que proporcione buenas "predicciones" de la variable y
- ▶ una alternativa (razonable) es estimar el método $m(\cdot)$ utilizando los datos de la muestra

estimación asistida por modelos

- ▶ el estimador \hat{Y}^{DIFF} necesita un método $m(\cdot)$ que sea independiente de la muestra
- ▶ en la práctica tenemos que utilizar los datos de la muestra para construir el método para predecir
- ▶ las estimaciones asistidas por modelos abordan este problema introduciendo un "working model"

$$y_i = m(\mathbf{x}_i) + \epsilon_i$$

donde ϵ_i tiene media cero.

- ▶ asumimos que $y_i \forall i \in U$ son realizaciones de un modelo **superpoblacional**

estimación asistida por modelos

- ▶ el modelo superpoblacional elegido no tiene porque ser verdadero para la población U
- ▶ es de ayuda que el modelo tenga cierto poder de ajuste respecto a la variable de interés y

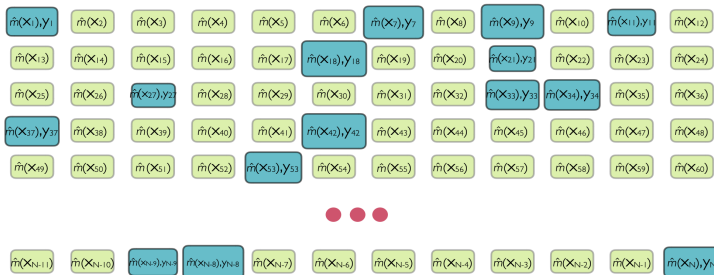
receta

una receta general para la estimación e inferencia utilizando información auxiliar es:

- ▶ si $(y_i, \mathbf{x}_i) \forall i \in U$ fueran observadas (conocidas) para toda la población podríamos computar $m(\cdot)$ utilizando métodos estadísticos (e.g. regresión lineal) que son independientes de la muestra
- ▶ debido a que en la práctica, únicamente tenemos datos de y para la muestra, el modelo $m(\cdot)$ es desconocido y lo tenemos que estimar $\hat{m}(\cdot)$. Notemos que $\hat{m}(\cdot)$ depende de la muestra seleccionada

la idea

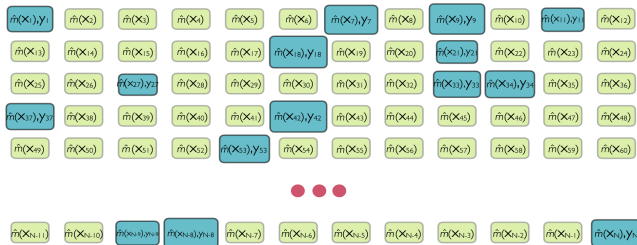
utilizando los datos de la muestra estimamos un modelo $\hat{m}(x_i)$ para predecir la variable de interés y



$$\hat{y}_i = \hat{m}(x_i)$$

la idea

construimos un estimador que sea robusto si el modelo/método elegido no es verdadero



Reemplazamos el método $m(\mathbf{x}_i)$ en \hat{Y}^{DIFF} por su correspondiente estimación $\hat{m}(\mathbf{x}_i)$:

$$\hat{Y}^{\text{PRED}} = \sum_{i \in U} \hat{m}(\mathbf{x}_i) + \sum_{i \in s} w_i (y_i - \hat{m}(\mathbf{x}_i))$$

propiedades

- ▶ es asintóticamente insesgado $\rightarrow E(\hat{Y}^{\text{PRED}}) = Y$
- ▶ no tiene una forma exacta de la varianza
- ▶ el estimador de la varianza se basa en la fórmula del estimador de diferencia

$$\widehat{\text{var}}(\hat{Y}^{\text{PRED}}) = \sum_{i \in s} \sum_{j \in s} \pi_{ij}^{-1} (\pi_{ij} - \pi_i \pi_j) \frac{(y_i - \hat{m}(\mathbf{x}_i))}{\pi_i} \frac{(y_j - \hat{m}(\mathbf{x}_j))}{\pi_j}$$

modelos que asisten al estimador

existe una amplia gama de modelos que se han utilizado para asistir al estimador. Algunos...

- ▶ **Regresión lineal** (Cassel, Sarndal, y Wretman 1976)
- ▶ Regresión logística (Lehtonen y Veijanen 1998)
- ▶ Redes neuronales (Montanari y Ranalli 2005)
- ▶ Árboles de regresión (McConville y Toth 2020)

modelos que asisten al estimador

Cómo elegimos el modelo que asiste al estimador?

Recordemos que el estimador asistido por modelos es asintoticamente insesgado

- ▶ para una amplia gama de métodos/modelos/algoritmos
- ▶ no importa que el modelo elegido sea correcto
- ▶ PERO su **precisión** depende del poder predictivo del modelo!

$$\widehat{\text{var}}(\hat{Y}^{\text{PRED}}) = \sum_{i \in s} \sum_{j \in s} \pi_{ij}^{-1} (\pi_{ij} - \pi_i \pi_j) \frac{(y_i - \hat{m}(\mathbf{x}_i))}{\pi_i} \frac{(y_j - \hat{m}(\mathbf{x}_j))}{\pi_j}$$

estimador de regresión

$$\hat{Y}^{\text{PRED}} = \sum_{i \in U} \hat{m}(\mathbf{x}_i) + \sum_{i \in s} w_i (y_i - \hat{m}(\mathbf{x}_i))$$

El método que asiste al estimador es un modelo de regresión lineal:

$$\begin{aligned} y_i &= m(\mathbf{x}) + \epsilon_i \\ &= \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_J x_{iJ} + \epsilon_i \\ &= \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i \end{aligned}$$

donde $\epsilon_i \sim (0, \sigma_i^2)$

estimador de regresión

el estimador de regresión queda definido como:

$$\begin{aligned}\hat{Y}^{\text{GREG}} &= \sum_{i \in U} \hat{y}_i + \sum_{i \in s} w_i (y_i - \hat{y}_i) \\ &= \sum_{i \in U} \mathbf{x}_i^T \hat{\mathbf{B}} + \sum_{i \in s} w_i (y_i - \mathbf{x}_i^T \hat{\mathbf{B}})\end{aligned}$$

donde

$$\hat{\mathbf{B}} = \left(\sum_{i \in s} w_i \mathbf{x}_i \mathbf{x}_i^T / \sigma_i \right)^{-1} \left(\sum_{i \in s} w_i \mathbf{x}_i y_i / \sigma_i \right),$$

es el estimador de β a nivel muestral utilizando el método de mínimos cuadrados ponderados.

casos especiales

- ▶ **estimador post-estratificado:** una sola variable categórica de entrada x (e.g. tramo de edad, departamento, sexo)
- ▶ **estimador de razón:** una única variable numérica de entrada x en donde el modelo no tiene intercepto (la recta de regresión pasa por el origen)

ventajas del estimador de regresión

Flexibilidad en el requisito de información auxiliar.

únicamente es necesario conocer:

- ▶ $\mathbf{x}_i \forall i \in s$ (i.e. podemos relevar los datos en el formulario de la encuesta).
- ▶ conocer los totales de las covariables $\mathbf{X} = \sum_{i \in U} \mathbf{x}_i$. No es necesario tener la información a nivel micro!

el estimador de regresión se puede expresar como:

$$\hat{Y}^{\text{GREG}} = \mathbf{X}^T \hat{\mathbf{B}} + \sum_{i \in s} w_i (y_i - \mathbf{x}_i^T \hat{\mathbf{B}})$$

propiedades del estimador de regresión

bajo un modelo de regresión lineal, el estimador puede ser expresado como una suma ponderada:

$$\hat{Y}^{\text{GREG}} = \sum_{i \in s} w_i^* y_i,$$

en donde los ponderadores $w_i^* = g_i w_i$ con un factor de ajuste:

$$g_i = 1 + (\mathbf{X} + \hat{\mathbf{X}}^{\text{HT}})^T \left(\sum_{i \in s} w_i \mathbf{x}_i \mathbf{x}_i^T / \sigma_i \right)^{-1} \mathbf{x}_i / \sigma_i,$$

- ▶ los ponderadores w_i^* **no dependen** de la variable y
- ▶ los podemos utilizar para distintas variables de interés!

propiedades del estimador de regresión

el estimador se encuentra **calibrado** a la información auxiliar utilizada:

$$\sum_{i \in s} w_i^* \mathbf{x}_i = \sum_{i \in U} \mathbf{x}_i$$

propiedad atractiva para la producción de estadísticas oficiales debido a que brinda consistencia entre datos de distintas fuentes.

problemas

Existen algunos problemas:

- ▶ puede existir mucha variabilidad en los ponderadores w_i^* lo que repercute en un aumento innecesario de los SEs
- ▶ pueden incluso existir ponderadores w_i^* negativos!
- ▶ si el modelo no es correcto puede ocurrir que no se obtenga una mejora alguna en la precisión en comparación al estimador HT