

Muestreo asistido por modelos

Daniel Czarniewicz

2017

El objetivo es utilizar información auxiliar para construir estimadores más precisos que el estimador \hat{t}_π . Asumimos que tenemos información completa para J variables auxiliares x , las cuales covarían con y . Para el k -ésimo elemento definimos el vector $\mathbf{x}_k = (x_{1;k}; \dots; x_{j;k}; \dots; x_{J;k})'$.

Se extrae una muestra s según el diseño $p(\cdot)$ con $\pi_j > 0$ y $\pi_{kl} > 0 \forall k; l \in U$. A partir de esta muestra se desea estimar $t_y = \sum_U y_k$.

El estimador de diferencia

La idea principal de este estimador es utilizar la información auxiliar para formar un set de N valores proxy para la variable de análisis: $y_1^0; \dots; y_k^0; \dots; y_N^0$, de forma tal que sean buenas aproximaciones para los valores $y_1; \dots; y_k; \dots; y_N$. Una opción es construir estos valores proxy a partir de una combinación lineal de las J variables auxiliares:

$$y_k^0 = \sum_{j=1}^J A_j x_{jk} = \mathbf{A}' \mathbf{x}_k$$

donde A_j son coeficientes conocidos $\forall j$. De esta forma, y_k^0 puede calcularse para toda la población (dado que los valores de las variables auxiliares se asumen disponibles para toda la población). Dado esto, el total poblacional podría expresarse como:

$$t_y = \sum_U y_k = \sum_U y_k^0 + \sum_U (y_k - y_k^0) = \sum_U y_k^0 + \sum_U D_k$$

$\sum_U y_k^0$ es conocida para toda la población, pero $\sum_U D_k$ no lo es, y debe estimarse. Para esto se utiliza la estimación π , de forma de lograr un estimador insesgado, el cual se conoce como estimador de diferencias:

$$\begin{aligned} \star \hat{t}_{y_{diff}} &= \sum_U y_k^0 + \sum_s D_k^\gamma \\ \star E(\hat{t}_{y_{diff}}) &= E\left(\sum_U y_k^0 + \sum_s D_k^\gamma\right) = E\left(\sum_U y_k^0\right) + E\left(\sum_s D_k^\gamma\right) = \\ &= \sum_U y_k^0 + \sum_U D_k = t_y \\ \star V(\hat{t}_{y_{diff}}) &= \sum \sum_U \Delta_{kl} D_k^\gamma D_l^\gamma \\ \star \hat{V}(\hat{t}_{y_{diff}}) &= \sum \sum_s \Delta_{kl}^\gamma D_k^\gamma D_l^\gamma \\ \star E[\hat{V}(\hat{t}_{y_{diff}})] &= \sum \sum_U \Delta_{kl} D_k^\gamma D_l^\gamma \end{aligned}$$

Si el diseño es de tamaño fijo, entonces se cumple que:

$$\begin{aligned} \star V(\hat{t}_{y_{diff}}) &= -\frac{1}{2} \sum \sum_U \Delta_{kl} (D_k^\gamma - D_l^\gamma)^2 \\ \star V(\hat{t}_{y_{diff}}) &= -\frac{1}{2} \sum \sum_s \Delta_{kl}^\gamma (D_k^\gamma - D_l^\gamma)^2 \end{aligned}$$

Alternativamente, el estimador $\hat{t}_{y_{diff}}$ podría considerarse como una mejora sobre el estimador π . Si los valores de la variable proxy se generan como una combinación lineal, entonces:

$$\star \hat{t}_{y_{diff}} = \hat{t}_{y_\pi} + \sum_{j=1}^J A_j (t_{x_j} - \hat{t}_{x_{j\pi}})$$

El estimador de regresión

Si los coeficientes A_j no son conocidos, entonces estos deben ser estimados. Sus estimaciones serán $\hat{B}_1; \dots; \hat{B}_J$. El estimador de regresión será entonces:

$$\star \hat{t}_{y_r} = \hat{t}_{y_\pi} + \sum_{j=1}^J \hat{B}_j (t_{x_j} - \hat{t}_{x_{j\pi}})$$

donde:

$$\hat{\mathbf{B}} = \left(\sum_s \frac{\mathbf{x}_k \mathbf{x}_k'}{\sigma_k^2 \pi_k} \right)^{-1} \left(\sum_s \frac{\mathbf{x}_k y_k}{\sigma_k^2 \pi_k} \right) = \hat{\mathbf{T}}^{-1} \hat{\mathbf{t}}$$

Las estimaciones de los parámetros B_j se obtienen como el resultado de ajustar un modelo con las siguientes características:

- $y_1; \dots; y_N$ son realizaciones de N variables aleatorias iid $Y_1; \dots; Y_N$
- $E_\xi(Y_k) = \sum_{j=1}^J \beta_j x_{jk} \quad \forall k \in U$
- $V_\xi(Y_k) = \sigma_k^2 \quad \forall k \in U$

En un censo, el estimador MCP del vector β sería:

$$\mathbf{B} = \left(\sum_U \frac{\mathbf{x}_k \mathbf{x}_k'}{\sigma_k^2} \right)^{-1} \left(\sum_U \frac{\mathbf{x}_k y_k}{\sigma_k^2} \right) \Rightarrow \mathbf{B} = (\mathbf{X}' \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} (\mathbf{X}' \boldsymbol{\Sigma}^{-1} \mathbf{y}) = \mathbf{T}^{-1} \mathbf{t}$$

donde:

$$\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2; \dots; \sigma_N^2) \Rightarrow \boldsymbol{\Sigma}^{-1} = \text{diag}(\sigma_1^{-2}; \dots; \sigma_N^{-2})$$

Si el muestreo se hace con remplazo, el estimador de regresión puede construirse utilizando el estimador *pwr* en lugar del estimador π .

El modelo ajustado para la muestra s produce la estimación de \mathbf{B} , los valores ajustados $\hat{y}_k = \mathbf{x}_k' \hat{\mathbf{B}}$ para cada elemento de la población, y los residuos $e_{k_s} = y_k - \hat{y}_k$ para cada elemento en la muestra. El estimador de regresión puede expresarse también como:

$$\star \hat{t}_{y_r} = \sum_U \hat{y}_k + \sum_s e_{k_s}^\vee$$

El término de ajuste desaparece en muchas aplicaciones, aún cuando el ajuste del modelo no sea perfecto. Por ejemplo, si $\sigma_k^2 = \lambda' \mathbf{x}_k \Rightarrow \sum_s e_{k_s}^\vee = 0$.

$$\begin{aligned} \sum_s e_{k_s}^\vee &= \sum_s y_k^\vee - \left(\sum_s \frac{\mathbf{x}_k}{\pi_k} \right) \hat{\mathbf{B}} = \sum_s y_k^\vee - \left(\sum_s \frac{\sigma_k^2 \mathbf{x}_k}{\sigma_k^2 \pi_k} \right) \hat{\mathbf{B}} = \sum_s y_k^\vee - \left(\sum_s \frac{\lambda' \mathbf{x}_k \mathbf{x}_k'}{\sigma_k^2 \pi_k} \right) \hat{\mathbf{B}} = \\ &= \sum_s y_k^\vee - \lambda' \left(\sum_s \frac{\mathbf{x}_k \mathbf{x}_k'}{\sigma_k^2 \pi_k} \right) \left(\sum_s \frac{\mathbf{x}_k \mathbf{x}_k'}{\sigma_k^2 \pi_k} \right)^{-1} \left(\sum_s \frac{\mathbf{x}_k y_k}{\sigma_k^2 \pi_k} \right) = \\ &= \sum_s y_k^\vee - \lambda' \left(\sum_s \frac{\mathbf{x}_k y_k}{\sigma_k^2 \pi_k} \right) = \sum_s y_k^\vee - \sum_s \frac{\lambda' \mathbf{x}_k y_k}{\sigma_k^2 \pi_k} = \sum_s y_k^\vee - \sum_s \frac{\sigma_k^2 y_k}{\sigma_k^2 \pi_k} = \\ &= \sum_s y_k^\vee - \sum_s \frac{y_k}{\pi_k} = \sum_s y_k^\vee - \sum_s y_k^\vee = 0 \end{aligned}$$

Este supuesto respecto de la estructura de σ_k^2 se cumplirá si:

- $\sigma_k^2 = \sigma^2 \quad \forall k \in U$ y $x_{1k} = 1 \quad \forall k \in U$
- $\exists x_j / \sigma_k^2 \propto x_{jk} \quad \forall k \in U$
- $\sigma_k^2 \propto \sum_{j=1}^J a_j x_{jk} \quad \forall k \in U$

Sean $\mathbf{t}_x = (t_{x_1}; \dots; t_{x_J})'$ el vector J -dimensional de los totales de las J variables auxiliares, y $\hat{\mathbf{t}}_{x_\pi} = (\hat{t}_{x_1\pi}; \dots; \hat{t}_{x_J\pi})'$ el vector J -dimensional de sus estimadores π , entonces el estimador de regresión puede escribirse como:

$$\begin{aligned} \star \hat{t}_{y_r} &= \hat{t}_{y_\pi} + (\mathbf{t}_x - \hat{\mathbf{t}}_{x_\pi})' \hat{\mathbf{B}} = \sum_s y_k^\vee + (\mathbf{t}_x - \hat{\mathbf{t}}_{x_\pi})' \hat{\mathbf{T}}^{-1} \left(\sum_s \frac{\mathbf{x}_k y_k^\vee}{\sigma_k^2} \right) \Rightarrow \\ &\Rightarrow \hat{t}_{y_r} = \sum_s y_k^\vee \underbrace{\left[1 + (\mathbf{t}_x - \hat{\mathbf{t}}_{x_\pi})' \hat{\mathbf{T}}^{-1} \frac{\mathbf{x}_k}{\sigma_k^2} \right]}_{g_{k_s}} \Rightarrow \boxed{\hat{t}_{y_r} = \sum_s y_k^\vee g_{k_s}} \end{aligned}$$

La última forma de expresar el estimador de regresión utiliza los valores ajustados $y_k^0 = \mathbf{x}_k' \mathbf{B}$ y los residuos del modelo ajustado: $E_k = y_k - y_k^0$. Dado que $y_k = y_k^0 + E_k$, el estimador de regresión toma la forma:

$$\begin{aligned} \star \hat{t}_{y_r} &= \sum_s g_{k_s} (y_k^{0\vee} + E_k^\vee) \\ \sum_s g_{k_s} \mathbf{x}_k^{\vee'} &= \sum_s \left[1 + (\mathbf{t}_x - \hat{\mathbf{t}}_{x_\pi})' \hat{\mathbf{T}}^{-1} \frac{\mathbf{x}_k}{\sigma_k^2} \right] \frac{\mathbf{x}_k'}{\pi_k} = \\ &= \sum_s \mathbf{x}_k^{\vee'} + (\mathbf{t}_x - \hat{\mathbf{t}}_{x_\pi})' \hat{\mathbf{T}}^{-1} \underbrace{\left(\sum_s \frac{\mathbf{x}_k \mathbf{x}_k'}{\sigma_k^2 \pi_k} \right)}_{\hat{\mathbf{T}}} = \sum_s \mathbf{x}_k^{\vee'} + (\mathbf{t}_x - \hat{\mathbf{t}}_{x_\pi})' = \\ &= \hat{\mathbf{t}}_{x_\pi}' + \mathbf{t}_x' - \hat{\mathbf{t}}_{x_\pi}' = \mathbf{t}_x' = \sum_U \mathbf{x}_k' \end{aligned}$$

Post multiplicando por \mathbf{B} obtenemos:

$$\sum_s g_{k_s} \underbrace{\mathbf{x}_k^{\vee'} \mathbf{B}}_{y_k^{0\vee}} = \sum_U \underbrace{\mathbf{x}_k' \mathbf{B}}_{y_k^0} \Rightarrow \sum_s g_{k_s} y_k^{0\vee} = \sum_U y_k^0$$

Por lo tanto,

$$\begin{aligned} \hat{t}_{y_r} &= \sum_s g_{k_s} (y_k^{0\vee} + E_k^\vee) = \sum_s g_{k_s} y_k^{0\vee} + \sum_s g_{k_s} E_k^\vee \Rightarrow \\ &\Rightarrow \boxed{\hat{t}_{y_r} = \sum_U y_k^0 + \sum_s g_{k_s} E_k^\vee} \end{aligned}$$

La varianza del estimador de regresión

El estimador de regresión no es insesgado, pero es aproximadamente insesgado para muestras grandes. Él mismo puede aproximarse por un desarrollo de Taylor de primer orden:

$$\hat{t}_{y_0} = \hat{t}_{y_\pi} + (\mathbf{t}_x - \hat{\mathbf{t}}_{x_\pi})' \mathbf{B} = \sum_U y_k^0 + \sum_s E_k^\vee$$

Demostración:

$$\hat{t}_{y_r} = \hat{t}_{y_\pi} + (\mathbf{t}_x - \hat{\mathbf{t}}_{x_\pi})' \hat{\mathbf{T}}^{-1} \hat{\mathbf{t}} = f(\hat{t}_{y_\pi}; \hat{t}_{x_\pi}; \hat{\mathbf{T}}; \hat{\mathbf{t}})$$

$$\star \frac{\partial f}{\partial \hat{t}_{y_\pi}} = 1$$

$$\star \frac{\partial f}{\partial \hat{t}_{x_\pi}} = -\hat{B}_j \quad \forall j = 1; \dots; J$$

$$\star \frac{\partial f}{\partial \hat{t}_{jj'\pi}} = (\mathbf{t}_x - \hat{\mathbf{t}}_{x_\pi})' (-\hat{\mathbf{T}}^{-1} \mathbf{\Lambda}_{jj'} \hat{\mathbf{T}}^{-1}) \hat{\mathbf{t}} \quad \forall j \leq j' = 1; \dots; J$$

$$\star \frac{\partial f}{\partial \hat{t}_{j_0\pi}} = (\mathbf{t}_x - \hat{\mathbf{t}}_{x_\pi})' \hat{\mathbf{T}}^{-1} \lambda_j \quad \forall j = 1; \dots; J \text{ donde } \lambda_j \text{ es un vector } J\text{-dimensional con valor 1 en la } j\text{-ésima componente, y valor 0 en las demás. } \mathbf{\Lambda}_{jj'} \text{ es una matriz de tamaño } J \times J \text{ con valor 1 en los elementos } (jj') \text{ y } (j'j) \text{ y 0 en los demás componentes. Evaluamos en } \hat{\theta} = \theta \text{ para obtener:}$$

$$\begin{aligned} \star \hat{t}_{y_r} &\doteq \hat{t}_{y_{r_0}} = t_y + \mathbf{1}(\hat{t}_{y_\pi} - t_y) - \sum_{j=1}^J B_j(\hat{t}_{x_{j\pi}} - t_{x_j}) = \\ &= \hat{t}_{y_\pi} + (\mathbf{t}_x - \hat{\mathbf{t}}_{x_\pi})' \mathbf{B} = \sum_U y_k^0 + \sum_s E_k^\vee \\ \star E(\hat{t}_{y_r}) &\doteq E(\hat{t}_{y_{r_0}}) = E\left(\sum_U y_k^0 + \sum_s E_k^\vee\right) = \sum_U y_k^0 + E\left(\sum_s E_k^\vee\right) = \\ &= \sum_U y_k^0 + \sum_U E_k = t_y \\ \star AV(\hat{t}_{y_r}) &\doteq V(\hat{t}_{y_{r_0}}) = V\left(\sum_s E_k^\vee\right) = \sum \sum_U \Delta_{kl} E_k^\vee E_l^\vee \\ \star \hat{V}(\hat{t}_{y_r}) &= \sum \sum_s \Delta_{kl}^\vee E_k^\vee E_l^\vee \\ \star IC_{t_y}^{(1-\alpha)100\%} &= \left[\hat{t}_{y_r} \pm z_{1-\alpha/2} \sqrt{\hat{V}(\hat{t}_{y_r})} \right] \end{aligned}$$