

Regression in complex sampling designs

①

Stat 576

5-27-25

Suppose that the population model is:

$$Y_j = \beta_0 + \beta_1 X_j + \varepsilon_j, \quad j=1, 2, \dots, N$$

Using least squares, the parameter values would be

$$\beta_1 = \frac{S_{XY}}{S_x^2}, \quad \beta_0 = \bar{Y} - \beta_1 \bar{X} = \frac{\bar{Y}}{N} - \beta_1 \frac{\bar{t}_x}{N}$$

where $S_x^2 = \frac{1}{N-1} \sum_{j=1}^N (X_j - \bar{X})^2$

and $S_{XY} = \frac{1}{N-1} \sum_{j=1}^N (X_j - \bar{X})(Y_j - \bar{Y})$

Note: $(N-1)S_x^2 = \sum_{j=1}^N X_j^2 - 2\bar{X} \sum_{j=1}^N X_j + N\bar{X}^2$

$$= \sum_{j=1}^N X_j^2 - N\bar{X}^2$$

$$= \sum_{j=1}^N X_j^2 - N \left(\frac{\bar{t}_x}{N} \right)^2 = t_{x^2} - \frac{\bar{t}_x^2}{N}$$

②

and $(N-1)S_{XY} = \sum_{j=1}^N X_j Y_j - \bar{Y} \sum_{j=1}^N X_j - \bar{X} \sum_{j=1}^N Y_j + N\bar{X}\bar{Y}$

$$= \sum_{j=1}^N X_j Y_j - \bar{Y} \cdot N\bar{X} - \bar{X} \cdot N\bar{Y} + N\bar{X}\bar{Y}$$

$$= \sum_{j=1}^N X_j Y_j - N\bar{X}\bar{Y} = t_{xy} - \frac{\bar{t}_x \bar{t}_y}{N}$$

To estimate ρ_0 and p_i , replace $t_x, t_r, t_{xz}, t_{xy}, \bar{Y}, \bar{X}$
with their estimators.

The Horvitz-Thompson estimator of t_y is

$$\hat{t}_y = \sum_{i=1}^n \frac{y_i}{\pi_i} = \sum_{i=1}^n w_i y_i, \text{ where } w_i = \frac{1}{\pi_i}$$

$$\text{Similarly, } \hat{t}_x = \sum_{i=1}^n w_i x_i, \quad \hat{t}_r = \sum_{i=1}^n w_i x_i^2, \quad \hat{t}_{xz} = \sum_{i=1}^n w_i x_i y_i$$

But it may be that we don't know N , and
we also need to estimate it.

$$E[w_i] = E\left[\frac{1}{\pi_i}\right] = \sum_{j=1}^N \frac{1}{\pi_j} \cdot p_j, \text{ where}$$

p_j is the probability of selecting Y_j on a particular draw.

$$\begin{aligned} \pi_i &= \text{Prob}(Y_i \text{ is in the sample}) \\ &= \sum_{i=1}^n \text{Prob}(Y_i \text{ was drawn on the } i^{\text{th}} \text{ draw}) \\ &= n p_j \quad \Rightarrow p_j = \frac{\pi_i}{n} \end{aligned}$$

$$\text{So } E[w_i] = \sum_{j=1}^N \frac{1}{\pi_j} \frac{\pi_i}{n} = \frac{N}{n}$$

$$\text{Then } E\left[\sum_{i=1}^n w_i\right] = \sum_{i=1}^n \frac{N}{n} = n \frac{N}{n} = N \quad (5)$$

That is, the sum of the sample weights is an unbiased estimator of the population size.

[Note: In SRS, each $\pi_i = \frac{n}{N}$ so $w_i = \frac{N}{n}$
and $\sum_{i=1}^n w_i = N$ exactly.]

$$\begin{aligned} \text{Now, } \hat{\beta}_1 &= \frac{\bar{x}_{yy} - \frac{\bar{x}_y \bar{x}_y}{N}}{\bar{x}_{xx} - \frac{\bar{x}_x^2}{N}} \quad (6) \\ &= \frac{\sum_{i=1}^n w_i x_i y_i - \frac{\left(\sum_{i=1}^n w_i x_i\right)\left(\sum_{i=1}^n w_i y_i\right)}{\sum_{i=1}^n w_i}}{\sum_{i=1}^n w_i x_i^2 - \frac{\left(\sum_{i=1}^n w_i x_i\right)^2}{\sum_{i=1}^n w_i}} \end{aligned}$$

$$\begin{aligned}
 \text{Also, } \hat{\beta}_0 &= \frac{\bar{t}_y}{N} - \hat{\beta}_1 \frac{\bar{t}_x}{N} \\
 &= \frac{\sum_{i=1}^n w_i y_i}{\sum_{i=1}^n w_i} - \hat{\beta}_1 \cdot \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}
 \end{aligned}
 \tag{7}$$

These estimators are exactly the same as those
in Weighted least squares!

In ordinary least squares, the matrix version
of the model is $\mathbf{Y} = \mathbf{X}\vec{\beta} + \vec{\epsilon}$,

where $\vec{\epsilon} \sim N(\vec{0}, \sigma^2 \mathbf{I})$, and the solution is

$$\hat{\beta} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}$$

In weighted least squares, the model is the same, but
we assume $\vec{\epsilon} \sim N(\vec{0}, \sigma^2 \mathbf{V})$, where
 \mathbf{V} is a diagonal matrix of constants.

The solution is $\hat{\beta} = (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{Y}$

The elements of V^{-1} are the weights.

(9)

In our case, $V = \begin{bmatrix} \pi_1 & & 0 \\ & \ddots & \\ 0 & \ddots & \pi_n \end{bmatrix}$ and $V^{-1} = \begin{bmatrix} w_1 & & 0 \\ & \ddots & \\ 0 & \ddots & w_n \end{bmatrix}$

Thus, what we are doing is equivalent to assuming
that $\varepsilon_i \sim N(0, \sigma^2 \pi_i)$

WLS is designed to give smaller weights to y -values that
have larger variances.

This is analogous to the H-T estimator giving smaller
weights to y values that had higher probabilities