

Hi Anthony,

It turns out your idea basically works.

However, there are some subtleties occurring in finite populations with design-based inference. In a standard setting this may not be a big deal, but it could be if inference is on a "small" population such as schools and/or if the number of possible samples from a population under a sampling design is small. A simple solution would be to do the procedure several times and take the average of estimates as the estimate. A nice plus is that this "automatically" happens with replicate weights.

Here follows a (perhaps somewhat overly detailed) description of why.

- Daniel

Definitions

- Let U be the population of units, with unit i any element of U , and $|U|$ the size of the population.
- Let A be a subset of the population U and let \mathcal{A} be the collection of subsets of U that contains all possible samples.
- let the sampling design $p(a)$ assign a probability measure to any sample $a \in \mathcal{A}$, i.e. $p(a) = P(A = a)$ and denote by E_p the expectation over the probability space \mathcal{A} under design $p(a)$.
- Let I_i be an indicator variable for whether unit $i \in A$, and let $\pi_i := E_p(I_i)$
- Assume π_i and $\pi_{ij} := E_p(I_i I_j)$ are known for any $\{i, j\} \in A$

Procedure

- Create weights

$$w_i := \pi_i^{-1} + e_i,$$

where e_i is drawn from an i.i.d. probability distribution ξ , with mean μ_e and standard deviation σ_e . For example $e_i \sim \mathcal{N}(0, \sigma_e^2)$

- Define the linear estimator

$$T_A := \sum_{i \in A} w_i y_i$$

To be shown: Under the definitions above,

T_A is unbiased in the sense that

$$\lim_{|U| \rightarrow \infty} E_p(T_A) = T_U := \sum_{i \in U} y_i$$

Proof of design-unbiasedness in the limit

$$\begin{aligned}
E_p(T_A) &= E_p \left(\sum_{i \in A} w_i y_i \right) \\
&= E_p \left(\sum_{i \in U} I_i w_i y_i \right) \\
&= E_p \left(\sum_{i \in U} I_i (\pi_i^{-1} + e_i) y_i \right) \\
&= E_p \left(\sum_{i \in U} I_i \pi_i^{-1} y_i \right) + E_p \left(\sum_{i \in U} e_i y_i \right) \\
&= \sum_{i \in U} E_p(I_i) \pi_i^{-1} y_i + \sum_{i \in U} E_p(e_i y_i) \\
&= T_U + \sum_{i \in U} E_p(e_i y_i),
\end{aligned} \tag{1}$$

that is, the design bias $E_p(T_A) - T_U$ equals $\sum_{i \in U} E_p(e_i y_i)$.

The bias is therefore close to zero but not exactly equal to zero.

To see this, it is possible to imagine two processes for e_i leading to slightly different results:

- i. e_i is fixed over \mathcal{A} (i.e. drawn once);
- ii. e_i is random over \mathcal{A} (i.e. drawn separately for each possible sample).

For fixed e_i ,

$$\sum_{i \in U} E_p(e_i y_i) = \sum_{i \in U} e_i y_i \tag{2}$$

The design bias can then be recognized as the maximum likelihood estimator (MLE) of the cross-product moment under the probability distribution ξ , that is,

$$\sum_{i \in U} e_i y_i = |U| \hat{\text{Cov}}_\xi(e y_i) - \hat{\mu}_e T_U, \tag{3}$$

where, collecting the parameters $\hat{\theta} := [\hat{\mu}, \hat{\sigma}_e \hat{\text{Cov}}_\xi(e y_i)]$, the population MLE is

$$\hat{\theta} = \arg \max_{\theta} \xi[\theta; \{(y_i, e_i) : i \in U\}]. \tag{4}$$

If the researcher chooses ξ such that $\mu_e = 0$ and $\text{Cov}_\xi(e y_i) = 0$, as is the case for $\xi = \mathcal{N}(0, \sigma_e)$, then

$$\lim_{|U| \rightarrow \infty} [E_p(T_A) - T_U] = 0 \tag{5}$$

by Equation 3 and the standard result in maximum-likelihood theory that the probability limit of a quantity can be obtained by replacing parameter estimates by their probability limits. This demonstrates design-unbiasedness in the limit.

Remarks You can wonder how far away from 0 the design bias $E_p(T_A) - T_U$ is likely to be. Again, applying standard results on ML estimation of covariances, it follows that as $|U| \rightarrow \infty$ the design bias will follow, over ξ , a Normal distribution with variance proportional to the fourth-order central moment of e_i , i.e. $\text{Var}[E_p(T_A) - T_U] = |U|^{-1} E_\xi[(e - \mu_e)^4]$.

Since ξ is likely to be known, it is reasonable to use the result for the finite population. For example, when choosing $\xi = \mathcal{N}(\mu_e, \sigma_e)$,

$$\text{Var}_\xi[E_p(T_A) - T_U] = |U|^{-1} 3\sigma_e^4.$$

What happens when we consider e_i random rather than fixed over \mathcal{A} ? For random e_i , each $E_p(e_i y_i)$ is itself an aggregate of numbers converging to 0 as the number of draws from ξ increases, while $\sum_{i \in U} E_p(e_i y_i)$ is an aggregate of aggregates. For example, if e_i is drawn once from $\mathcal{N}(0, \sigma_e^2)$ for each sample A ,

$$\text{Var}_\xi[E_p(T_A) - T_U] \approx (|U| + |\mathcal{A}|)^{-1} 3\sigma_e^4.$$

Thus, the more draws of e_i , the closer the design bias will be to 0.