



UNITED STATES DEPARTMENT OF COMMERCE
Economics and Statistics Administration
U.S. Census Bureau
Washington, DC 20233-0001

MASTER FILE

December 31, 2002

DSSD A.C.E. REVISION II MEMORANDUM SERIES #PP- 40

MEMORANDUM FOR Donna Kostanich
Chair, A.C.E. Revision II Planning Group

From: William R. Bell (*signed 12/31/02*) *WRB*
Senior Mathematical Statistician for Small Area Estimation

Subject: A.C.E. Revision II: Calculating aggregate data-defined, correct enumeration, and census inclusion rates (for groups that involve aggregation across post-strata)

The attached document discusses how to make aggregate tabulations of A.C.E. Revision II results for the purpose of calculating certain rates related to census coverage for aggregate population groups. While these aggregate rates are not actually used in constructing the A.C.E. Revision II dual system estimates (DSEs), they are useful descriptive statistics for examining approximately how the DSE components differ across aggregate groups.

Calculating aggregate data-defined, correct enumeration, and census inclusion rates (for groups that involve aggregation across post-strata)

William R. Bell

December 31, 2002

This note discusses how to make aggregate tabulations of A.C.E. Revision II results for the purpose of calculating certain rates related to census coverage for aggregate population groups. The purpose of the note is to address questions that were raised regarding how such tabulations should be done given the different post-strata used for the E- and P-samples. The rates considered here are data-defined rates (proportions of census enumerations that are data-defined persons), correct enumeration (CE) rates, and census inclusion rates. Two types of CE rates are considered, one reflecting the proportion of data-defined persons that are correct enumerations, and the other reflecting the proportion of all census enumerations that are correct enumerations. The first is analogous to the usual E-sample CE-rate, the latter we shall call the “census correct enumeration rate.”

Data Defined Rates

Let a denote a geographic area or population subgroup of interest. I can think of two versions of a DD-rate for a that could be of interest. One uses the DDs for a computed from direct tabulations, the other uses the synthetically estimated DDs for a . The rate computed with the former I’ll call the “actual DD-rate.” It makes sense as a descriptive statistic since we know whether all census records are data defined or not. The rate computed with the latter DD tabulation I’ll call the “estimated DD-rate.” It makes sense since it relates to how we compute DDs in computing the synthetic estimates from the DSEs. Below I give formulas for computing both these rates.

Let t index all person records in any given detailed hl post-stratum, including both data-defined persons and non-data-defined persons. To start we need the census tabulations for the intersection of a with the detailed hl post-strata (note that $\sum_{t \in (hl) \cap a}$ indicates summation over all person records t within a that are in the detailed hl post-stratum). We also need the full census tabulation for a .

- $\text{Cen}_{hl,a} = \sum_{t \in (hl) \cap a} 1$
- $\text{Cen}_{\bullet a} = \sum_{hl} \text{Cen}_{hl,a} = \sum_{t \in a} 1.$

We compute the actual DDs in the analogous way:

- $\text{DD}_{hl,a}^{\text{act}} = \sum_{t \in (hl) \cap a} I(t \text{ is data defined})$
- $\text{DD}_{\bullet a}^{\text{act}} = \sum_{hl} \text{DD}_{hl,a}^{\text{act}} = \sum_{t \in a} I(t \text{ is data defined}).$

We compute estimated DDs for a synthetically:

- $DD_{hl,a}^{\text{est}} = r_{DD,hl} \times \text{Cen}_{hl,a}$
- $DD_{\bullet a}^{\text{est}} = \sum_{hl} DD_{hl,a}^{\text{est}} = \sum_{hl} (r_{DD,hl} \times \text{Cen}_{hl,a})$

In the above two expressions $r_{DD,hl} = DD_{hl}^{\text{act}} / \text{Cen}_{hl}$ is the estimated DD-rate for the hl post-stratum that is used in the DSE. Notice that if a were simply an hl post-stratum then the equation for $DD_{hl,a}^{\text{est}}$ would become $DD_{hl}^{\text{est}} = DD_{hl}^{\text{act}}$. This equivalence of the estimated and actual DD tabulations occurs only for the hl post-strata and for their direct aggregates (not for groups that include parts of hl post-strata).

The two DD-rates we can compute are as follows:

- Actual DD-rate $_a = DD_{\bullet a}^{\text{act}} / \text{Cen}_{\bullet a}$
- Estimated DD-rate $_a = DD_{\bullet a}^{\text{est}} / \text{Cen}_{\bullet a}$.

One final point to note is that for the synthetic calculations that follow we will need the $\text{Cen}_{hl,a}$, but we do not really need to retain the detailed hl post-stratum tabulations of $DD_{hl,a}^{\text{act}}$ or $DD_{hl,a}^{\text{est}}$.

Correct Enumeration Rates

Since CEs are identified only for the E-sample cases, and not the whole census, we can only estimate aggregate CE rates synthetically. The synthetic estimate of CEs for area or subgroup a is calculated as follows:

- $\text{CE}_{hl,a} = r_{ce,i} \times DD_{hl,a}^{\text{est}} = r_{ce,i} \times r_{DD,hl} \times \text{Cen}_{hl,a}$
- $\text{CE}_{\bullet a} = \sum_{hl} \text{CE}_{hl,a} = \sum_{hl} (r_{ce,i} \times r_{DD,hl} \times \text{Cen}_{hl,a})$

where i is the E-sample post-stratum determined by the more detailed post-stratum hl . The corresponding CE-rate is then

- $\text{CE-rate}_a = \text{CE}_{\bullet a} / DD_{\bullet a}^{\text{est}}$.

This rate reflects the estimated proportion of data-defined persons in area or subgroup a that were correct enumerations. It is analogous to the estimated CE-rates ($r_{ce,i}$) for the E-sample post-strata. In fact, if a is an E-sample post-stratum i , the above rate reduces to $r_{ce,i} \times DD_{\bullet i}^{\text{est}} / DD_{\bullet i}^{\text{est}} = r_{ce,i}$.

Note that the denominator of CE-rate_a is the estimate of data defined persons, not the census count. This is because in the DSE the E-sample CE-rates $r_{ce,i}$ are applied to the estimates of data defined persons. We can also calculate the following ratio:

- census CE-rate $_a = \text{CE}_{\bullet a} / \text{Cen}_{\bullet a}$

which can be called the “census correct enumeration rate” (for area or subgroup a). It reflects the estimated proportion of all census enumerations for area or subgroup a that were correct. Since the census count is the sum of data-defined and non-data-defined persons, $\text{Cen}_{\bullet a} \geq \text{DD}_{\bullet a}^{\text{est}}$ and census $\text{CE-rate}_a \leq \text{CE-rate}_a$. Note that census CE-rate_a reflects for the area or subgroup a the combined effects of the CE- and DD-rates ($r_{ce,i} \times r_{DD,hl}$).

Census Inclusion Rates and Undercount Rates

Census inclusion rates estimate the proportion of the true population actually included in the census. The true population for a is estimated by the synthetic estimate for a from the DSEs including any adjustment for correlation bias. Let \hat{N} denote such synthetic estimates (reserving “DSE” for conventional DSEs computed without correlation bias adjustments). Under the “two-group” model for adjusting for correlation bias the \hat{N} are given by

$$\begin{aligned} \bullet \hat{N}_{hl,a} &= c_k \frac{\text{CE}_{hl,a}}{r_{m,j}} = \frac{c_k}{r_{m,j}} \times r_{ce,i} \times r_{DD,hl} \times \text{Cen}_{hl,a} = \text{CCF}_{hl} \times \text{Cen}_{hl,a} \\ \bullet \hat{N}_{\bullet a} &= \sum_{hl} \hat{N}_{hl,a} = \sum_{hl} \left(\frac{c_k}{r_{m,j}} \times r_{ce,i} \times r_{DD,hl} \times \text{Cen}_{hl,a} \right) = \sum_{hl} (\text{CCF}_{hl} \times \text{Cen}_{hl,a}) \end{aligned}$$

where j is the P-sample post-stratum determined by hl and c_k is the correlation bias adjustment factor for age-race group k (with k also determined by hl). Note that the c_k are 1 for children and adult females. The coverage correction factors for the hl post-strata are:

$$\text{CCF}_{hl} = \frac{c_k}{r_{m,j}} \times r_{ce,i} \times r_{DD,hl}.$$

The census inclusion rate for a is

$$\bullet \text{Census inclusion rate}_a = \text{CE}_{\bullet a} / \hat{N}_{\bullet a}.$$

This rate uses as its numerator $\text{CE}_{\bullet a}$, because this represents our estimate of the number of correct census inclusions from the true population. Since $c_k \geq 1$ and $r_{m,j} \leq 1$, $\hat{N}_{\bullet a} \geq \text{CE}_{\bullet a}$, and this inclusion rate is ≤ 1 . This rate reflects the effects of dividing in the DSEs by the P-sample match rates and then multiplying the DSEs by the correlation bias factors c_k from the two-group model. If a were a P-sample post-stratum that included the full age-sex detail then the census inclusion rate would be simply $r_{m,j}/c_k$. If we were not adjusting for correlation bias in the DSEs (as is the case for children and adult females) then the census inclusion rate for a P-sample post-stratum is simply the match rate for that post-stratum.

If we replace $\text{CE}_{\bullet a}$ in the numerator of the above rate by the census count we get a net census coverage “rate,” which would exceed 1 for net overcounts. More familiar is the related census undercount rate for a :

$$\bullet \text{Census undercount rate}_a = 1 - (\text{Cen}_{\bullet a} / \hat{N}_{\bullet a})$$

which is negative for overcounts.

Note: Multiplying any of the above rates by 100 expresses them as percentages.