## **Constraining Arrays of Numbers in Demography**

## Charles D. Coleman Administrative Records Methodology and Research Branch U.S. Census Bureau, Mail Stop 8800 Washington, DC 20233-8800 charles.d.coleman@census.gov

Demographic work often requires constraining arrays of numbers to controls in one or two dimensions. For example, subnational population projections may be constrained to a national projection. If the results are allowed to take on any nonnegative values, raking solves the problem in one dimension and two-way iterative raking solves it in two dimensions. The problem is more complicated in one dimension if the data can be of any sign, the so-called "plus-minus" problem, as simple raking may produce unacceptable results. This problem is addressed by generalized raking, which preserves the structure of the data at a cost of a nonunique solution. Since demography is concerned with people, which come in whole units, data often have to be rounded to integers. The Cox-Ernst algorithm accomplishes an optimal controlled rounding in two dimensions. In one dimension, the Greatest Mantissa algorithm applies a simplified version of the Cox-Ernst algorithm.

This paper combines into one place all of the above-mentioned problems and techniques. It is written to provide practical guidance to demographers and other practitioners who work with these problems. The demographic variables that can be input to these procedures include population, proportions and rates. The introductory section briefly describes each type of problem and points the reader to the relevant section. Each section begins with practical advice, followed by a technical discussion.

In all of these methods, final data should, in some way, preserve the structure of the original data. In one dimension, this means that if one initial data element is greater than another, then the transformed elements should preserve this relationship. Optimally, the ratios between elements should be preserved, but this can only be done in the specific case of controlling a vector whose nonzero elements are of the same sign as the control value, with the result left unrounded. Controlled rounding then destroys the ratios, but preserves the order. The effect on the ratios depends on the magnitudes of the initial data elements and unit of rounding. The effect of controlled rounding on these ratios in vectors with large elements relative to the unit of rounding is minimal. On the other hand, initial vectors with elements about the same magnitude as the unit of rounding will find these ratios greatly perturbed. "Generalized raking" of a vector of mixed sign or to zero or control of opposite sign to the nonzero data preserves the order of the original elements, while destroying the ratios. The two-dimensional equivalent of raking minimizes a function that measures the distortion from the original matrix.

One-dimensional "raking" multiplies a vector of data by the ratio of the control to the sum of the initial data. Damage to the structure of the original data is avoided only when the control is of the same sign as the nonzero initial data. When the data are of mixed sign or the control is zero or opposite sign of the nonzero data, "generalized raking" takes a weighted average of the ordinarily raked data and their projection onto the hyperplane defined by the control. The result, except in the case of a zero control, is

Charles D. Coleman, "Constraining Arrays of Numbers in Demography" extended abstract for PAA 2006, Page 1 of 2

nonunique. Generalized raking has several advantages over the earlier Akers-Siegel procedure: a continuous transformation using arithmetic operations is used instead of separate rakes for positive and negative data, it easily handles zeroes, and there is never any need to arbitrarily shift and then rake the data, when it is impossible to apply the A-S method to the original data. Generalized raking does fail when the original data sum to zero: a simple workaround is proposed. When the original data span a wide range of magnitudes, it is possible to simply "stuff" (that is, add) the difference between the control and the sum of the original data to the element of largest absolute value. The choice between stuffing and the generalized rake in these instances is up to the analyst: if the distortion caused by stuffing is minimal, then stuffing may be advised to simplify programming. Stuffing after rounding is necessary due to current lack of an algorithm to do controlled rounding of mixed-sign data.

All of the one-dimensional raking procedures are geometrically illustrated. The ordinary and generalized rakes define straight lines when viewed as element-wise operations, while the Akers-Siegel procedure defines two rays that intersect at the origin. The ordinary and generalized rakes are also viewed as projections of a point onto the hyperplane defined by the control.

Two-dimensional raking, a.k.a. "iterative proportionate fitting" and the "RAS algorithm," is a much-rediscovered method for constraining a nonnegative matrix to positive row and column controls. The sums of the row and column controls (or "marginals") must be equal for it to work. It proceeds by alternately raking row data in parallel to row controls and column data to column controls until convergence. It minimizes a function that measures the distortion of the data. The result is unique. A sufficient condition for feasibility is that the original matrix be positive. When this does not obtain, an algorithm based on linear programming can be used to determine feasibility. Some practical guidance for handling zeroes and "low" (i.e., values too low to be reported) is given to speed convergence. This procedure does not generalize to three or more dimensions: a positive array with positive marginals can still be infeasible.

The Cox-Ernst controlled rounding procedure assures that integers are unchanged and nonintegers are rounded to one of their closest integers. A cost to unconventional rounding (that is, rounding in the opposite direction of conventional rounding) is defined and minimized. Because of its complexity, this algorithm is not described in detail. For vectors, the Greatest Mantissa algorithm performs controlled rounding by rounding up numbers in order of their mantissas. This simplification of the Cox-Ernst algorithm is simple to describe and program.

The Appendix contains plug-in SAS macros to implement each procedure. These macros enable the user to use these algorithms immediately without additional programming.

Extended abstract prepared for submission to the 2006 meetings of the Population Association of America.

Charles D. Coleman, "Constraining Arrays of Numbers in Demography" extended abstract for PAA 2006, Page 2 of 2