

BUREAU OF THE CENSUS
STATISTICAL RESEARCH DIVISION REPORT SERIES
SRD Research Report Number: CENSUS/SRD/RR-84/17

MAKING TABLES ADDITIVE IN THE PRESENCE OF ZEROS

by

James Fagan and Brian Greenberg
Statistical Research Division
U.S. Bureau of the Census

Room 3587, F.O.B. #3
Washington, D.C. 20233

(301)763-7530

This series contains research reports, written by or in cooperation with staff members of the Statistical Research Division, whose content may be of interest to the general statistical research community. The views reflected in these reports are not necessarily those of the Census Bureau nor do they necessarily represent Census Bureau statistical policy or practice. Inquiries may be addressed to the author(s) or the SRD Report Series Coordinator, Statistical Research Division, Bureau of the Census, Washington, DC 20233.

Recommended by: Paul Biemer
Report completed: August 7, 1984
Report issued: August 7, 1984

Making Tables Additive in the Presence of Zeros

James Fagan and Brian Greenberg

Given a two-way contingency table of non-negative reals in which the internal entries do not sum to the corresponding marginals, there is often the need to adjust internal entries to achieve additivity. In general, the objective is to have the revised table, in some sense, close to the original table and to have zero entries remain zero and positive entries remain positive. Not all two-way contingency tables can be adjusted to achieve additivity subject to the constraints above and in this paper we present a procedure that will determine whether a given table can be so adjusted.

I. INTRODUCTION

Given a two-way contingency table of non-negative reals in which the internal entries do not sum to the corresponding marginals, there is often the need to adjust internal entries to achieve additivity. In general, the objective is to have the revised table, in some sense, close to the original table and to have zero entries remain zero and positive entries remain positive. Not all two-way contingency tables can be adjusted to achieve additivity subject to the constraints above and in this paper we present a procedure that will determine whether a given table can be so adjusted, and such adjustable tables will be called feasible.

The most frequently used procedure for adjusting tables that are not additive is iterative proportional fitting, often called raking. The raking algorithm alternately scales rows and columns to achieve respective additivity, and if a table is feasible the algorithm converges. This algorithm is frequently used to reconcile tabular data when the marginals and internal entries arise from different sources, for example see [7].

Raking has been extensively used for over forty years, and its statistical properties have been well-studied, see [1]. However, there has never been a satisfactory answer to the following question: given an arbitrary non-additive table, is it feasible? That is, there was no known procedure to rigorously determine whether raking or any other table adjustment methodology that preserves zeros and leaves positive entries positive will converge for an arbitrary non-additive table. In this paper we present such a procedure.

In Section II we introduce terminology and provide an analytical formulation of the problem. In the next section we attack the problem using the classical transportation problem of operations research. We describe a finite iterative procedure which can be applied to an arbitrary non-additive table, and by examining the outcome of the final

iteration, one can determine if the original table is feasible. The final section briefly discusses alternative methods for table adjustment.

II. FEASIBLE TABLES

By a **contingency table** we mean a triple $A = \{(a_{ij}), r, c\}$ of arrays of non-negative reals where (a_{ij}) is an $R \times C$ matrix, $r = (r_1, \dots, r_R)$, $c = (c_1, \dots, c_C)$, and

$$\sum_{i=1}^R r_i = \sum_{j=1}^C c_j .$$

We say that A is **additive** if

$$\sum_{j=1}^C a_{ij} = r_i \quad i = 1, \dots, R$$

$$\sum_{i=1}^R a_{ij} = c_j \quad j = 1, \dots, C .$$

The table A is said to be **feasible** if there exists an $R \times C$ matrix (b_{ij}) such that $b_{ij} = 0$ if and only if $a_{ij} = 0$ and such that $B = \{(b_{ij}), r, c\}$ is additive. That is, A is feasible if and only if there exists an $R \times C$ matrix (x_{ij}) such that $B = \{(b_{ij}), r, c\}$ is additive, where $(b_{ij}) = (x_{ij} a_{ij})$, and $x_{ij} > 0$ whenever $a_{ij} > 0$. In particular, A is feasible if there exist x_{ij} satisfying the following system:

$$(1) \quad \sum_{j=1}^C a_{ij} x_{ij} = r_i \quad i=1, \dots, R$$

$$(2) \quad \sum_{i=1}^R a_{ij} x_{ij} = c_j \quad j=1, \dots, C$$

$$(3) \quad x_{ij} > 0 \text{ if } a_{ij} > 0 \quad i=1, \dots, R \text{ and } j=1, \dots, C .$$

By way of examples, **Table 1** is clearly feasible and **Tables 2 or 3** are clearly not.

0	1	5
1	1	4
3	6	

Table 1

0	1	5
1	1	4
6	3	

Table 2

1	0	1	4
1	0	1	1
0	1	1	4
0	1	1	4
5	4	4	

Table 3

Note that **Table 2** fails conditions (1) and (2) above, while **Table 3** does satisfy these conditions letting:

$$x_{11} = 4 \quad x_{21} = 1, \quad x_{33} = x_{34} = x_{43} = x_{44} = 2, \quad \text{and} \quad x_{13} = x_{23} = 0;$$

yet fails the joint conditions (1), (2) and (3).

If some r_q for $q=1, \dots, R$ (or c_p for $p=1, \dots, C$) equals zero, then for a contingency table to be additive or feasible it is necessary that $a_{qj} = 0$ for all $j = 1, \dots, C$ ($a_{ip} = 0$ for all $i = 1, \dots, R$). That is, the entire row (or column) must be zero, and hence can be removed from the table. Thus, we can assume without loss of generality that both r and c are positive.

The objective of this paper is as follows. Given an arbitrary (non-additive) table $A = \{(a_{ij}), r, c\}$ find a finite iterative procedure that will determine if A is feasible. That is, determine if there exists an $R \times C$ matrix (x_{ij}) such that (1)-(3) are satisfied. In the next section we apply the classical transportation problem to obtain a finite step-by-step procedure that will solve the problem stated above.

III. A PROCEDURE TO DETERMINE FEASIBILITY

A. The Transportation Problem

A well studied and frequently used construct in the realm of operations research is the **transportation problem**. The objective (in its purest form) is to minimize the cost of shipping a commodity from a number of origins to various destinations. We assume that there are R origins and C destinations, $r_i > 0$ units are to be shipped from the i^{th} origin for $i=1, \dots, R$ and $c_j > 0$ units are to be received at the j^{th} destination for $j=1, \dots, C$, and the cost of shipping a unit from origin i to destination j is c_{ij} . One usually defines $C = (c_{ij})$ to be the **cost matrix**. In the classical transportation problem one further assumes that

$$\sum_{i=1}^R r_i = \sum_{j=1}^C c_j,$$

and seeks to minimize the function

$$(4) \quad \sum_{i=1}^R \sum_{j=1}^C c_{ij} x_{ij}$$

subject to the constraints:

$$(5) \quad \sum_{j=1}^C x_{ij} = r_i \quad i=1, \dots, R$$

$$(6) \quad \sum_{i=1}^R x_{ij} = c_j \quad j=1, \dots, C$$

$$(7) \quad x_{ij} \geq 0 \quad i=1, \dots, R \text{ and } j=1, \dots, C$$

where x_{ij} is the number of units shipped from origin i to destination j .

Given the transportation problem (4)-(7), if r_i for $i=1, \dots, R$ and c_j for $j=1, \dots, C$ are integers, there exists an RC -dimensional vector, (Z_{ij}) , such that (Z_{ij}) minimizes (4) subject to (5)-(7) and (Z_{ij}) has integer components, see [4] for a discussion. Given a table $A = \{(a_{ij}), \mathbf{r}, \mathbf{c}\}$, we can scale \mathbf{r} and \mathbf{c} by the same factor and assume henceforth that \mathbf{r} and \mathbf{c} are integer vectors.

B. The General Case

If we have a table $A = \{(a_{ij}), \mathbf{r}, \mathbf{c}\}$, we can form the table $M = \{(m_{ij}), \mathbf{r}, \mathbf{c}\}$ where

$$m_{ij} = \begin{cases} 0 & \text{if } a_{ij} = 0 \\ 1 & \text{if } a_{ij} \neq 0. \end{cases}$$

It is clear that A is feasible if and only if M is feasible. Looking back to (1), (2), and (3), $M = \{(m_{ij}), \mathbf{r}, \mathbf{c}\}$ is feasible if there exists x_{ij} such that

$$(8) \quad \sum_{j=1}^C m_{ij} x_{ij} = r_i \quad i=1, \dots, R$$

$$(9) \quad \sum_{i=1}^R m_{ij} x_{ij} = c_j \quad j=1, \dots, C$$

$$(10) \quad x_{ij} > 0 \text{ if } m_{ij} > 0 \quad i=1, \dots, R \text{ and } j=1, \dots, C.$$

Given the table M , consider the following sequence of transportation problems indexed by positive integers, q :

Minimize (11) $\sum_{i=1}^R \sum_{j=1}^C c_{ij}^q x_{ij}$

subject to

$$(12) \quad \sum_{j=1}^C x_{ij} = r_i \quad i=1, \dots, R$$

$$(13) \quad \sum_{i=1}^R x_{ij} = c_j \quad j=1, \dots, C$$

$$(14) \quad x_{ij} \geq 0,$$

where

$$c_{ij}^1 = \begin{cases} T & \text{if } m_{ij} = 0 \\ 0 & \text{otherwise,} \end{cases}$$

$$T = \sum_{i=1}^R r_i = \sum_{j=1}^C c_j,$$

and for $q > 1$,

$$c_{ij}^{q+1} = \begin{cases} 1 & \text{if } c_{ij}^q = 1 \text{ or } x_{ij}^q \neq 0 \text{ and } m_{ij} \neq 0 \\ T & \text{if } m_{ij} = 0 \\ 0 & \text{otherwise,} \end{cases}$$

where (x_{ij}^q) minimizes (11) subject to (12)-(14).

Denote the region determined by constraints (8)-(10) by Ω_M^* and note that Ω_M^* is not empty if and only if M is feasible. Define $\Omega_M = \{(y_{ij}) : (y_{ij}) \in \Omega_M^* \text{ and if } m_{ij} = 0 \text{ then } y_{ij} = 0\}$. Clearly, $\Omega_M \subset \Omega_M^*$ and if $\Omega_M \neq \emptyset$, then $\Omega_M^* \neq \emptyset$, so M is feasible if and only if $\Omega_M \neq \emptyset$. Denoting the region determined by the constraints (12)-(14) by Ω_T , we observe that $\Omega_M \subset \Omega_T$.

Notation: Denote by $R \times C$ the set $\{(i, j) : i=1, \dots, R, j=1, \dots, C\}$, and by C^q the minimal value of (11) subject to (12)-(14).

Lemma 1: There exists a positive integer k such that $C^k \geq T$.

Proof: If $C^k < T$, there exists $(t, s) \in R \times C$ such that $c_{ts}^k = 0$ and $x_{ts}^k > 0$. For if not, whenever $x_{ij}^k > 0$, then $c_{ij}^k \geq 1$, so

$$C^k = \sum_{i=1}^R \sum_{j=1}^C c_{ij}^k x_{ij}^k \geq \sum_{i=1}^R \sum_{j=1}^C x_{ij}^k = T.$$

Accordingly, if $C^k < T$ there exists $(t,s) \in R \times C$ such that $c_{st}^k = 0$ and $c_{st}^{k+1} = 1$. Hence, since the set $R \times C$ is finite, for some positive integer, k , $C^k \geq T$.

Notation: Let $N = \min \{ k \in \mathbb{Z}^+ : C^k \geq T \}$.

Lemma 2: If $C^1 \neq 0$, then $C^1 \geq T$ and M is not feasible.

Proof: If $C^1 \neq 0$ there exists an integer array $(w_{ij}) \in \Omega_T$ such that

$$C^1 = \sum_{i=1}^R \sum_{j=1}^C c_{ij}^1 w_{ij}.$$

For some $(t,s) \in R \times C$, $c_{ts}^1 = T$ and $w_{ts} \geq 1$, (otherwise $C^1 = 0$). Thus, if $C^1 \neq 0$, then $C^1 \geq T$.

If M is feasible, let $(y_{ij}) \in \Omega_M$ and observe that

$$C^1 = \sum_{i=1}^R \sum_{j=1}^C c_{ij}^1 y_{ij} = 0.$$

That is, if $C^1 \neq 0$ then M is not feasible.

Lemma 3: If $C^1 = 0$, then $C^N = T$, and C^k is a non-decreasing function of k for $k=1, \dots, N$.

Proof: Note that $C^1 = 0$ if and only if there exists $(x_{ij}^1) \in \Omega_T$ such that for all $(i,j) \in R \times C$ if $m_{ij} = 0$, then $x_{ij}^1 = 0$. Thus, if $x_{ij}^1 \neq 0$ then $m_{ij} = 1$, so $c_{ij}^k \leq 1$, for all $k=1, \dots, N$ and so

$$C^k \leq \sum_{i=1}^R \sum_{j=1}^C c_{ij}^k x_{ij}^1 \leq \sum_{i=1}^R \sum_{j=1}^C x_{ij}^1 = T.$$

Hence, if $C^1 = 0$, then $C^k \leq T$ for all $k=1, \dots, N$. It is clear that C^k is a non-decreasing function of k .

Theorem 1: Suppose $C^1 = 0$ and N is as above. Then M is feasible if and only if $c_{ij}^{N+1} > 0$ for all $(i,j) \in R \times C$.

Proof: (only if) Suppose M is feasible and there exists $(t,s) \in R \times C$ such that $c_{ts}^{N+1} = 0$, and note that $c_{ts}^{N+1} = 0$ implies that $c_{ts}^N = 0$.

Choose $(y_{ij}^N) \in \Omega_M$ and note

$$C^N \leq \sum_{i=1}^R \sum_{j=1}^C c_{ij}^N y_{ij}^N < \sum_{i=1}^R \sum_{j=1}^C v_{ij}^N = T.$$

The strict inequality holds because: (1) if $y_{ij}^N \neq 0$ then $c_{ij}^N \leq 1$ by the definition of Ω_M , and (2) $y_{ts}^N > 0$, yet $c_{ts}^N = 0$. But this contradicts the fact that $C^N = T$.

(if) For each $(t,s) \in R \times C$ such that $m_{ts} = 1$, there exists a q such that $(x_{ij}^q) \in \Omega_T$ and $x_{ts}^q > 0$ because $c_{ij}^{N+1} > 0$ for all $(i,j) \in R \times C$.

Let

$$(z_{ij}) = \sum_{k=1}^N (x_{ij}^k) / N.$$

Since $(x_{ij}^k) \in \Omega_T$, for all $k=1, \dots, N$, then $(z_{ij}) \in \Omega_T$ because Ω_T is a convex set. Also, if $m_{ij} = 1$, then $z_{ij} \geq 0$, so $(z_{ij}) \in \Omega_M^*$. Thus M is feasible.

Iterative Procedure to Determine Feasibility: Given a contingency table $A = \{(a_{ij}), r, c\}$, to determine whether or not it is feasible proceed as follows. Scale r and c so that they are integer and form $M = \{(m_{ij}), r, c\}$ as above. Solve the first transportation problem above, obtaining C^1 . If $C^1 \neq 0$, then A is not feasible. If $C^1 = 0$, form C^2, C^3 , etc., until $C^N = T$, and examine the cost matrix (c_{ij}^{N+1}) . If $c_{st}^{N+1} = 0$ for any $(s,t) \in R \times C$, then A is not feasible, otherwise A is feasible.

C. Non-degenerate Solution

Recall that when given an R by C transportation problem, we say that an optimal solution is non-degenerate if there are exactly $R+C-1$ non-zero variables in the solution. In this case, by reordering the rows and columns of the underlying matrix, we can start at the upper left corner and traverse (more or less) staircase fashion to the bottom right corner stopping only at positive cells, see [4]. The following result enables one to possibly shorten the **Iterative Procedure** outlined above. That is, if any of the transportation problems above has a non-degenerate solution with optimum less than or equal to T then A is feasible. Thus, if $C^k \leq T$ one needs only count the non-zero variables in the solution vector. If that count is equal to $R+C-1$ then A is feasible, otherwise proceed to the next iteration and continue as indicated in the **Iterative**

Procedure with this addendum at each juncture.

Theorem 2: If $\mathbf{A} = \{(a_{ij}), \mathbf{r}, \mathbf{c}\}$ is a contingency table, then \mathbf{A} is feasible if the following transportation problem has a non-degenerate optimal solution of value less than or equal to T .

$$(15) \text{ Minimize } C^q = \sum c_{ij}^q y_{ij}$$

subject to

$$(16) \quad \sum_{j=1}^C y_{ij} = r_i \quad \text{for } i=1, \dots, R$$

$$(17) \quad \sum_{i=1}^R y_{ij} = c_j \quad \text{for } j=1, \dots, C$$

$$(18) \quad y_{ij} \geq 0 \quad \text{for } i=1, \dots, R \text{ and } j=1, \dots, C$$

where c_{ij}^q is defined as earlier.

Proof: Since \mathbf{A} is feasible if and only if \mathbf{M} (as above) is feasible we can focus our attention on \mathbf{M} . Assume (z_{ij}) is a non-degenerate optimal solution to (15)-(18) such that $C^q \leq T$, and suppose $a_{k\ell} \neq 0$ and $z_{k\ell} = 0$ for some $(k, \ell) \in R \times C$. Form a closed path starting and ending at $a_{k\ell}$, transversing only positive elements z_{ij} such that no three consecutive path elements are in the same row or column. That is, form the (+,-) path used in updating feasible non-degenerate solutions of the transportation problem (usually used in conjunction with the so-called Northwest corner solution). Let z be the minimal positive value for the cells in the path, and starting with the (k, ℓ) - position alternately add and subtract $z/2$ from each z_{ij} in the path updating the values of the z_{ij} . Repeat this procedure for all (i,j) positions such that $a_{ij} \neq 0$ and $z_{ij} = 0$. When there are no such cells remaining, then conditions (8) - (10) are satisfied by letting $(x_{ij}) = (z_{ij})$ and hence \mathbf{M} is a feasible table.

IV. STATISTICAL AND PRACTICAL CONSIDERATIONS

If $\mathbf{A} = \{(a_{ij}), \mathbf{r}, \mathbf{c}\}$ is a contingency table, we can let

$$\pi_{ij} = \frac{a_{ij}}{\sum_{i=1}^R \sum_{j=1}^C a_{ij}}$$

$$\pi_{i\cdot} = \frac{r_i}{\sum_{i=1}^R r_i}$$

$$\pi_{\cdot j} = \frac{c_j}{\sum_{j=1}^C c_j} .$$

Observe that \mathbf{A} is a feasible table if and only if $\Pi = \{(\pi_{ij}), (\pi_{i\cdot}), (\pi_{\cdot j})\}$ is a feasible table. Note further that

$$\sum_{j=1}^C \sum_{i=1}^R \pi_{ij} = \sum_{i=1}^R \pi_{i\cdot} = \sum_{j=1}^C \pi_{\cdot j} = 1,$$

so we enter the realm of probability theory.

Notation: Let $V = \{(i, j) : (i, j) \in R \times C \text{ and } \pi_{ij} \neq 0\}$,
 $V(i) = \{j : (i, j) \in V\}$, and
 $V(j) = \{i : (i, j) \in V\}$.

Given a feasible table of probabilities, $\Pi = \{(\pi_{ij}), (\pi_{i\cdot}), (\pi_{\cdot j})\}$ we seek an additive table $\mathbf{P} = \{(p_{ij}), (\pi_{i\cdot}), (\pi_{\cdot j})\}$ such that $p_{ij} = x_{ij} \pi_{ij}$ and $x_{ij} > 0$ for all $(i, j) \in V$. We say that \mathbf{P} is **derived** from Π , write (1) - (3) as

$$(19) \quad \sum_{j \in V(i)} x_{ij} \pi_{ij} = \pi_{i\cdot} \quad i=1, \dots, R$$

$$(20) \quad \sum_{i \in V(j)} x_{ij} \pi_{ij} = \pi_{\cdot j} \quad j=1, \dots, C$$

$$(21) \quad x_{ij} > 0 \quad (i, j) \in V,$$

and note that \mathbf{P} is also a table of probabilities.

For some tables, there is a unique derived table (the deterministic case), for example:

.13	.10	0	.35
0	.38	.07	.60
0	0	.32	.05
.10	.40	.50	

Table 4

If there is more than one derived table, there are infinitely many since any convex combination of derived tables is also a derived table.

In general, given the feasible table Π , we seek a derived table P such that P is close to Π . Of course, the notion of "close" is not unique, and for every criterion of closeness a different objective function must be optimized subject to (19) - (21). Listed below are three objective functions which are candidates for a criterion of closeness. Each objective function is convex up as is easily seen by examining the Hessians, see [5]. Thus, if the original table Π is feasible it is not too hard to see each has a unique minimum subject to (19)-(21). There has been interesting work to discover iterative procedures which will allow a user to start with a feasible table Π and proceed to an additive table P optimizing the objective functions below.

(i) Iterative Proportional Fitting:

$$\text{Minimize } \sum_{(i,j) \in V} p_{ij} \ln \frac{p_{ij}}{\pi_{ij}} \text{ over } (p_{ij}), \text{ which is equivalent to,}$$

$$\text{Minimize } \sum_{(i,j) \in V} \pi_{ij} x_{ij} \ln x_{ij} \text{ over } (x_{ij}).$$

(ii) Maximum Likelihood:

$$\text{Minimize } \sum_{(i,j) \in V} \pi_{ij} \ln \frac{\pi_{ij}}{p_{ij}} \text{ over } (p_{ij}), \text{ which is equivalent to,}$$

$$\text{Minimize } - \sum_{(i,j) \in V} \pi_{ij} \ln x_{ij} \text{ over } (x_{ij}).$$

(iii) Minimum Chi-Square:

Minimize $\sum_{(i,j) \in V} (p_{ij} - \pi_{ij})^2 / p_{ij}$ over (p_{ij}) , which is equivalent to,

Minimize $\sum_{(i,j) \in V} \pi_{ij} / x_{ij}$ over (x_{ij}) .

It is proved in both [3] and [7] that given a feasible table, the "raking algorithm" (alternately scaling rows and columns to achieve respective additivity) converges to a table minimizing the objective function for iterative proportional fitting. This algorithm has been put to many uses and the reader is referred to [6] for further discussion and extensive bibliography. It has been known for quite a while that raking converges when all entries in the contingency table are positive. If there are zeros in a table, and raking appears not to converge, adjustments are made to internal entries so that raking will converge for the revised table.

When raking does converge for some table, it does so rapidly, and less than ten iterations usually suffice so that successive internal values are within a reasonable tolerance. Accordingly, the practice has been to presume that raking will not converge for a table if it fails to converge within a prescribed number of iterations, and at that time, zero cells are promoted to non-zero status or cells are collapsed. We have presented here a procedure that can be used to test for feasibility if raking seems not to converge. That is, if there is no convergence after a fixed number of iterations one can now draw upon the procedures described above to determine if raking does fail to converge, or if it just needs more time.

For maximum likelihood and minimum chi-square, algorithms have been proposed for positive tables, see [2] for more details. It would be interesting to see proofs that these algorithms (iterative procedures) do, in fact, converge to additive tables; although it is easy to see that when they do converge to additive tables, they converge to tables optimizing the respective objective functions. It would be even more interesting to learn something about the existence and convergence of algorithms for tables containing zeros.

This paper will be presented at the 1984 Annual Meeting of the American Statistical Association in Philadelphia, Pennsylvania and will appear in the Proceedings of the Section on Survey Research Methods.

REFERENCES

1. Bishop, Y., Finberg, S., and Holland, P., (1975) Discrete Multivariate Analysis: Theory and Practice. MIT Press. Cambridge, Mass.
2. Causey, B. (1983) Estimation of Proportions for Multinomial Contingency Tables Subject to Marginal Constraints. *Communications in Statistics (A)*. 12, 22.
3. Darroch, J.M. and Ratcliff, D., (1972) Generalized Iterative Scaling for Log-Linear Models. *The Annals of Mathematical Statistics*. 43, 5.
4. Gass, S. (1975) Linear Programming. McGraw-Hill. New York.
5. Luenberger, D. (1973) Introduction to Linear and Nonlinear Programming. Addison-Wesley. Reading, MA.
6. Oh, H.L. and Scheuren, F.J. (1982) Some Unresolved Application Issues in Raking Ratio Estimation. *Proceedings of the Section on Survey Research Methods, American Statistical Association*.
7. Thompson, J. (1981) Convergence Properties of the Iterative 1980 Census Estimator. *Proceedings of the Section on Survey Research Methods, American Statistical Association*.