



---

Components of a Difference Between Two Rates

Author(s): Evelyn M. Kitagawa

Source: *Journal of the American Statistical Association*, Vol. 50, No. 272 (Dec., 1955), pp. 1168-1194

Published by: American Statistical Association

Stable URL: <http://www.jstor.org/stable/2281213>

Accessed: 31/03/2009 14:31

---

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/page/info/about/policies/terms.jsp>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Please contact the publisher regarding any further use of this work. Publisher contact information may be obtained at <http://www.jstor.org/action/showPublisher?publisherCode=astata>.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

JSTOR is a not-for-profit organization founded in 1995 to build trusted digital archives for scholarship. We work with the scholarly community to preserve their work and the materials they rely upon, and to build a common research platform that promotes the discovery and use of these resources. For more information about JSTOR, please contact [support@jstor.org](mailto:support@jstor.org).



American Statistical Association is collaborating with JSTOR to digitize, preserve and extend access to *Journal of the American Statistical Association*.

<http://www.jstor.org>

## COMPONENTS OF A DIFFERENCE BETWEEN TWO RATES\*

EVELYN M. KITAGAWA

*University of Chicago and Scripps Foundation*

WHEN comparing the incidence of some phenomenon in two or more groups, social researchers place much emphasis on the need for holding constant those related factors that would tend to distort the comparison. For example, before comparing the death rates for the residents of two areas, demographers frequently control the factors of differences between the areas in age, sex and race composition. A technique commonly used to accomplish this is "standardization" of the rates for the two areas by relating them both to a standard population with specified age-sex-race composition. By applying the schedule of age-sex-race specific death rates for each of the groups to the age-sex-race composition of the standard population, then noting the total death rate that results, it is possible to compare the death rates for the areas with reasonable confidence that differences in age, sex and race composition do not explain the differences between the rates for the areas that still remain after they have been standardized. Controlling the effect of related factors by this method is termed direct standardization.<sup>1</sup>

It is often noted that such standardized rates are "artificial." For example, an age-sex-race-standardized death rate indicates what the total (or crude)<sup>2</sup> death rate of a population would be *if* it had the age-sex-race composition of the standard population while retaining its own age-sex-race-specific death rates. While such a measure may not be very useful for descriptive purposes, it is an important analytical device.

Since the crude rate of any population is its "real" or "observed" rate in the (descriptive) sense that it is the total rate which results from the particular composition and specific rates which prevail in that population, a systematic statement of relationships between crude and standardized rates for two or more groups may help to bridge the gap

---

\* Expanded version of a paper read at the annual meeting of the American Statistical Association held in Chicago, December 27-30, 1952. The preparation of this manuscript was sponsored jointly by the Population Research and Training Center, University of Chicago, and Scripps Foundation, Miami University—the latter through funds provided by the Rockefeller Foundation for the study of population distribution. The author is indebted to Philip M. Hauser, Donald J. Bogue, O. Dudley Duncan, Beverly Duncan and J. J. Feldman for a careful reading of the paper and many suggestive comments and criticisms.

<sup>1</sup> For a description of the standardization procedure, the assumptions involved, and some of the limitations see A. J. Jaffe, *Handbook of Statistical Methods for Demographers* (Washington: U. S. Govt. Printing Office, 1951), Chapter III. For handling cases where a complete schedule of specific rates is not available for the particular areas being compared, but a cross-tabulation of the population by the variables is available, an alternative procedure termed "indirect standardization" has been developed.

<sup>2</sup> The terms "crude," "total," and "unstandardized" are used interchangeably in this paper.

between the "observed" crude (or total) rates and the "artificial" standardized rates. As yet, very little attention has been directed to the problem of formalizing the analysis of standardized rates, and of systematically explaining which factors account for the differences between standardized rates in comparison with corresponding differences between their unstandardized rates. If standardization alters a difference between two total rates, it should be possible to measure the amount of change, and to break it up into components attributable to the various factors for which the data were standardized. Formalizing the process of making inferences from standardized data, and establishing a technique whereby the change accomplished by standardization can be interpreted in terms of the factors involved, are the objectives of this paper.

The technique presented here is called "components of a difference between two rates." It is a revision and refinement of a mode of analysis utilized at the University of Chicago since 1948.<sup>3</sup> The purpose of the technique is to explain the difference between the total rates of two groups in terms of differences in their specific rates and differences in their composition. Thus, the components framework is broader in scope than that of standardized rates, since the framework of the latter is designed to summarize and compare differences in two (or more) sets of specific rates.

The basic concept of separation into components has been used in research for some time. Implicitly, it has been used whenever the size and direction of the difference between standardized rates for two populations was compared with the size and direction of the difference between their crude rates, and the inference made that the "difference between these two differences" is the result of the different compositions of the two populations. Explicitly, it has been used in research under various headings. For example, in the 1920's Ogburn determined what part of the difference between the per cent of the U. S. population married in 1890 and the per cent married in 1920 was due to changes in the age composition of the population between these two dates.<sup>4</sup> The chapter on standardization in Jaffe's *Handbook of Statistical Methods for Demographers* includes a section titled "Removing the influence of changing occurrence rates" which discusses a similar procedure. In his

---

<sup>3</sup> Earlier statements were prepared by Ralph H. Turner and the writer, with the counsel of Philip M. Hauser: Evelyn Kitagawa, "A Method of Analyzing the Influence of Several Non-Quantitative Factors on a Result," (Dept. of Sociology, University of Chicago, March, 1948, hectographed); Ralph Turner, "Whites and Negroes in the Labor Force" (unpublished Ph.D. dissertation, University of Chicago, September, 1948); Turner, "The Expected Cases Method Applied to the Nonwhite Male Labor Force," *American Journal of Sociology*, LV (September, 1949), 145-56.

<sup>4</sup> E. R. Groves and W. F. Ogburn, *American Marriage and Family Relationships* (New York: Henry Holt & Co., 1928), 160-2.

example, the difference between the proportion of the total population 5 to 20 years of age attending school in 1940 and the proportion attending school in 1910 is, in effect, separated into two parts, one called the "influence of changes in age distribution" and the other the "influence of changes in occurrence rates."<sup>5</sup> A similar procedure is applied to an analysis of the 1890-1930 change in the proportion of persons gainfully occupied, in an article by Wolfbein and Jaffe.<sup>6</sup>

Edwin C. Goldfield's "method of multiple standardization with allocation of interactions" is concerned with the same general problem but carries the analysis further to include consideration of several factors simultaneously, and to evaluate the net influence of each of the factors as well as their interaction. This method is used, and very briefly described for particular examples, in Durand's study of the labor force.<sup>7</sup>

The revised components framework described in the present article was first used in a study of labor mobility.<sup>8</sup> Any comparison of this definition of components with previous definitions requires a statement of the rationale underlying the components analysis, as well as some standard terminology and algebraic notation. These will be developed in the course of presenting the writer's revised approach, and a comparison of the various definitions will be made later.

As has already been noted, components are closely related to standardized rates. Because the method of standardizing rates is familiar to research workers, and because components may be defined by subtracting a standardized rate difference from a crude rate difference, the components analysis will be approached by putting standardized rates in a components framework.

The analysis will be developed first with two factors (*I* and *J*) controlled. Formulas for one factor (*I*) will be given later, and an abbreviated extension to control three or more factors will also be discussed. The factors included in the analysis may be either quantitative or non-quantitative.

---

<sup>5</sup> Jaffe, *op. cit.*, pp. 44-6. While Jaffe's discussion of the problem emphasizes its interpretation as a method of "holding constant changes in the occurrence rates," in contrast to the conventional method of computing standardized rates where composition is held constant, there is a clear separation into the two parts (or components) mentioned above, in Table 4, p. 46.

<sup>6</sup> S. L. Wolfbein and A. J. Jaffe, "Demographic Factors in Labor Force Growth," *American Sociological Review*, XI (August, 1946), 393-6.

<sup>7</sup> John D. Durand, *The Labor Force in the United States, 1890-1960* (New York: Social Science Research Council, 1948), Appendix B. Goldfield's method is very similar to the components analysis described in the 1948 manuscripts of Kitagawa and Turner.

<sup>8</sup> Evelyn M. Kitagawa, "The Relative Importance—and Independence—of Selected Factors in Job Mobility, Six Cities, 1940-49" (Chicago Community Inventory, University of Chicago, hectographed report, 1953).

CONVENTIONAL STANDARDIZATION IN A COMPONENTS FRAMEWORK

Suppose we have observed a difference between the crude rates of two groups,  $p$  and  $P$ .<sup>9</sup> Also, suppose that data for each group are cross-classified by two factors,  $I$  and  $J$ .

Let

$n_{ij}$  = number of persons in both the  $i$ th category of  $I$  and the  $j$ th category of  $J$  in population  $p$

$N_{ij}$  = number of persons in both the  $i$ th category of  $I$  and the  $j$ th category of  $J$  in population  $P$

$t_{ij}$  = rate for persons in the  $i$ th category of  $I$  and the  $j$ th category of  $J$  in population  $p$

$T_{ij}$  = rate for persons in the  $i$ th category of  $I$  and the  $j$ th category of  $J$  in population  $P$

Also

$n_{..}$  and  $N_{..}$  = total number of persons in populations  $p$  and  $P$ , respectively.

$t_{..}$  and  $T_{..}$  = rate for total persons in populations  $p$  and  $P$ , respectively (i.e., crude rates of  $p$  and  $P$ ).

The  $IJ$ -composition of  $p$  and  $P$  is designated by  $n_{ij}/n_{..}$  and  $N_{ij}/N_{..}$  respectively. That is,  $IJ$ -composition is a proportionate distribution in which the number of persons in each  $IJ$  cell has been divided by the total number of persons in the group. In conventional summation notation

$$n_{..} = \sum_i \sum_j n_{ij} \quad \text{and} \quad N_{..} = \sum_i \sum_j N_{ij}$$

$$t_{..} = \sum_i \sum_j \frac{n_{ij}}{n_{..}} \quad \text{and} \quad T_{..} = \sum_i \sum_j T_{ij} \frac{N_{ij}}{N_{..}} \text{ by definition.}^{10}$$

If we let  $p$  represent the group with the higher crude rate, the differ-

<sup>9</sup> The term crude rate is used here to refer to the over-all unstandardized rate of any specified group—that is, we may be concerned with an analysis of the components of the difference between crude rates of two age groups. Although such rates are usually called age-specific rates, in this context they are total or crude rates if the two age groups are the population groups with which we are concerned.

The concept of a rate in the components analysis includes percentages and means. Medians cannot be used because they do not have the algebraic properties of means and percentages which are utilized in the components method.

<sup>10</sup> Since the crude rate of any population may be regarded as the result of its own  $IJ$ -specific rates weighted by its own  $IJ$ -composition.

ence,  $t_{..} - T_{..}$ , will be positive. Expressing this difference in the equation

$$t_{..} - T_{..} = \sum_i \sum_j t_{ij} \frac{n_{ij}}{n_{..}} - \sum_i \sum_j T_{ij} \frac{N_{ij}}{N_{..}}$$

makes explicit the fact that the difference between the crude rates of  $p$  and  $P$  is due to differences in both their  $IJ$ -specific rates and their  $IJ$ -composition.

As has already been mentioned, conventional standardization techniques can be utilized to compute  $IJ$ -standardized rates for  $p$  and  $P$ , which will summarize differences in their  $IJ$ -specific rates *holding constant* differences in their  $IJ$ -composition. In the notation outlined above, the difference between  $IJ$ -standardized rates for  $p$  and  $P$  may be expressed as follows:

$$\sum_i \sum_j \frac{n_{ij}}{n_{..}} (t_{ij} - T_{ij}) \text{ if group } p \text{ is used as standard}^{11}$$

$$\sum_i \sum_j \frac{N_{ij}}{N_{..}} (t_{ij} - T_{ij}) \text{ if group } P \text{ is used as standard}$$

$$\sum_i \sum_j \frac{n'_{ij}}{n'_{..}} (t_{ij} - T_{ij}) \text{ if a third group, } p' \text{ is used as standard}$$

Thus, the difference between the standardized rates of two groups is a weighted average of differences in their  $IJ$ -specific rates, with the  $IJ$ -composition of the standard population used as the weights.

The objective of the components framework is to allocate the difference between two crude rates into components which reflect differences in specific rates of the two groups, on the one hand, and differences in their composition, on the other hand. The equations above suggest that the difference between the  $IJ$ -standardized rates of two groups might be used as the "component due to differences in specific rates," since this difference may be considered a measure of their differences in specific rates. Furthermore, if the difference between two standardized rates is subtracted from the corresponding difference between two crude rates, it is easily demonstrated that the result is a weighted average of differences between the composition of the two groups. For example, when the difference between the  $IJ$ -standardized rates of groups  $p$  and  $P$ , with group  $p$  used as the standard, is subtracted from

<sup>11</sup> Since  $\sum_i \sum_j t_{ij}(n_{ij}/n_{..}) = t_{..}$  = both the crude and  $IJ$ -standardized rate for  $p$ , and since  $\sum_i \sum_j T_{ij}(n_{ij}/n_{..})$  = the  $IJ$ -standardized rate for  $P$ , when  $p$  is used as the standard population.

the difference between their crude rates, the result is

$$\begin{aligned}
 (t_{..} - T_{..}) &= \sum_i \sum_j \frac{n_{ij}}{n_{..}} (t_{ij} - T_{ij}) \\
 &= \sum_i \sum_j t_{ij} \frac{n_{ij}}{n_{..}} - \sum_i \sum_j T_{ij} \frac{N_{ij}}{N_{..}} - \sum_i \sum_j t_{ij} \frac{n_{ij}}{n_{..}} + \sum_i \sum_j T_{ij} \frac{n_{ij}}{n_{..}} \\
 &= \sum_i \sum_j T_{ij} \left( \frac{n_{ij}}{n_{..}} - \frac{N_{ij}}{N_{..}} \right) = \text{weighted average of differences in} \\
 &\qquad\qquad\qquad \text{IJ-composition of } p \text{ and } P, \text{ with the} \\
 &\qquad\qquad\qquad \text{IJ-specific rates of } P \text{ as the weights.}
 \end{aligned}$$

In this case, the difference between the crude rates of *p* and *P* may be expressed as the sum of two major components as follows:

$$t_{..} - T_{..} = \sum_i \sum_j \frac{n_{ij}}{n_{..}} (t_{ij} - T_{ij}) + \sum_i \sum_j T_{ij} \left( \frac{n_{ij}}{n_{..}} - \frac{N_{ij}}{N_{..}} \right)$$

where

$$\begin{aligned}
 \sum_i \sum_j \frac{n_{ij}}{n_{..}} (t_{ij} - T_{ij}) &= \text{component due to differences in } IJ\text{-} \\
 &\qquad\qquad\qquad \text{specific rates} \\
 &= \text{difference between } IJ\text{-standardized rates} \\
 &\qquad\qquad\qquad (p \text{ as standard})
 \end{aligned}$$

$$\sum_i \sum_j T_{ij} \left( \frac{n_{ij}}{n_{..}} - \frac{N_{ij}}{N_{..}} \right) = \text{component due to differences in } IJ\text{-} \\
 \qquad\qquad\qquad \text{composition}$$

However, these equations show that different standard populations are used as weights in the two components. That is, while *p* may have been purposely selected as the standard population to provide weights to measure the difference in *IJ*-standardized rates, the net result—from the components framework—is that the other group, *P*, is the standard population which provides weights for summarizing differences in *IJ*-composition of *p* and *P*.

Similar results are obtained if the second population, *P*, is selected as the standard for computing *IJ*-standardized rates for *p* and *P*. Thus, it is also true that

$$t_{..} - T_{..} = \sum_i \sum_j \frac{N_{ij}}{N_{..}} (t_{ij} - T_{ij}) + \sum_i \sum_j t_{ij} \left( \frac{n_{ij}}{n_{..}} - \frac{N_{ij}}{N_{..}} \right)$$

where

$$\sum_i \sum_j \frac{N_{ij}}{N_{..}} (t_{ij} - T_{ij}) = \begin{array}{l} \text{component due to difference in } IJ\text{-specific} \\ \text{rates} \\ = \text{difference between } IJ\text{-standardized rates} \\ \text{(} P \text{ as standard)} \end{array}$$

$$\sum_i \sum_j t_{ij} \left( \frac{n_{ij}}{n_{..}} - \frac{N_{ij}}{N_{..}} \right) = \begin{array}{l} \text{component due to differences in } IJ\text{-com-} \\ \text{position.} \end{array}$$

In this case, when the  $IJ$ -composition of  $P$  is selected to weight differences in  $IJ$ -specific rates, the  $IJ$ -specific rates of  $p$  are implied as weights for the component due to differences in  $IJ$ -composition.

When a third population,  $p'$ , is used as the standard for computing  $IJ$ -standardized rates for  $p$  and  $P$ , the following components are the result:

$$(t_{..} - T_{..}) = \sum_i \sum_j \frac{n'_{ij}}{n'_{..}} (t_{ij} - T_{ij}) + \left[ \sum_i \sum_j t_{ij} \left( \frac{n_{ij}}{n_{..}} - \frac{n'_{ij}}{n'_{..}} \right) + \sum_i \sum_j T_{ij} \left( \frac{n'_{ij}}{n'_{..}} - \frac{N_{ij}}{N_{..}} \right) \right]$$

where

$$\sum_i \sum_j \frac{n'_{ij}}{n'_{..}} (t_{ij} - T_{ij}) = \begin{array}{l} \text{component due to differences in } IJ\text{-specific} \\ \text{rates of } p \text{ and } P \\ = \text{difference between } IJ\text{-standardized rates of} \\ p \text{ and } P \text{ (} p' \text{ as standard)} \end{array}$$

and

$$\sum_i \sum_j t_{ij} \left( \frac{n_{ij}}{n_{..}} - \frac{n'_{ij}}{n'_{..}} \right) + \sum_i \sum_j T_{ij} \left( \frac{n'_{ij}}{n'_{..}} - \frac{N_{ij}}{N_{..}} \right) = \begin{array}{l} \text{component due to differences in } IJ\text{-composition of } p \text{ and } P. \end{array}$$

The last equation makes explicit the weights implied in the "composition component" when a third population,  $p'$ , is used as the standard for computing standardized rates. Specifically, the difference in  $IJ$ -



composition between  $p$  and  $P$  is broken into two parts—the difference between  $p$  and  $p'$  or

$$\left( \frac{n_{ij}}{n_{..}} - \frac{n'_{ij}}{n'_{..}} \right)$$

and the difference between  $p'$  and  $P$  or

$$\left( \frac{n'_{ij}}{n'_{..}} - \frac{N_{ij}}{N_{..}} \right);$$

and  $IJ$ -specific rates of  $p$  are implied as weights for the first part, while  $IJ$ -specific rates of  $P$  are the weights for the second part.

Thus, the statement that a population,  $p'$ , is used as a standard population for computing  $IJ$ -standardized rates for  $p$  and  $P$  means specifically that the  $IJ$ -composition of  $p'$  is used to weight differences in  $IJ$ -specific rates of  $p$  and  $P$ . If, in interpreting the results of standardization, the difference between the standardized rate difference and the crude rate difference is attributed to the different composition of the two groups, it should be recognized that the  $IJ$ -specific rates of the two groups themselves are the weights, or standard, for summarizing differences in their  $IJ$ -compositions.

#### MAJOR COMPONENTS (TWO FACTORS, $I$ AND $J$ )

The components analysis starts directly and explicitly from the perspective of allocating a crude rate difference into parts attributable to differences in composition and specific rates. Suppose the crude rates refer to the two groups,  $p$  and  $P$ . Also, suppose that their specific rates and composition are classified by two factors,  $I$  and  $J$ . Then, an unambiguous allocation of the crude rate difference into two major components—one reflecting differences in  $IJ$ -composition only, and the other differences in  $IJ$ -specific rates only—will be obtained if two sets of weights,  $w_{ij}$  and  $w'_{ij}$ , are selected which satisfy the equation

$$t_{..} - T_{..} = \sum_i \sum_j w_{ij} \left( \frac{n_{ij}}{n_{..}} - \frac{N_{ij}}{N_{..}} \right) + \sum_i \sum_j w'_{ij} (t_{ij} - T_{ij}).$$

The first component on the right side of the equation represents a weighted sum of differences in  $IJ$ -composition, and the second component is a weighted sum of differences in  $IJ$ -specific rates. These components will be called "Combined  $IJ$ " and "Residual  $IJ$ ," respec-

tively.<sup>12</sup> A meaningful interpretation of the weights is obtained if the  $w_{ij}$  are a set of  $IJ$ -specific rates, and the  $w'_{ij}$  define an  $IJ$ -composition. In this case, the crude rate difference (which is due to differences in both  $IJ$ -composition and  $IJ$ -specific rates) is expressed as the sum of two components: (1) Combined  $IJ$ , or differences between the  $IJ$ -composition of  $p$  and  $P$ , with  $IJ$ -specific rates held constant; and (2) Residual  $IJ$ , or differences between the  $IJ$ -specific rates of  $p$  and  $P$ , with  $IJ$ -composition held constant.

Thus, the chief distinction between the components perspective and that of conventional standardized rates is in the specification of a standard population. In the components framework, the standard population must include a set of  $IJ$ -specific rates to be used as weights for the composition component, as well as an  $IJ$ -composition to be used as weights for the "specific rates" component.

The results of the preceding section indicate that the  $IJ$ -specific rates and the  $IJ$ -composition of the *same population* will not satisfy the requirements of a standard population for the components framework *if* the objective is an unambiguous two-component solution.<sup>13</sup> For example, when the  $IJ$ -composition of  $p$  was used as the standard for computing standardized rates, the  $IJ$ -specific rates of  $P$  were implied as the weights for the component due to differences in composition. Therefore the three sets of components discussed in the preceding section are excluded as possible two-component solutions unless the standard population is conceived as a population with the  $IJ$ -composition of one group and the  $IJ$ -specific rates of another group; for example, the composition of  $p$  and the specific rates of  $P$ . Although for particular problems a researcher may be willing to define such a standard population, the two-component solution proposed below would seem to be more generally useful from the components perspective.

Since the standard population for a two-component solution inevitably involves some characteristic of both  $p$  and  $P$ , the possibility of using the average composition and the average specific rates of these two groups to define the standard population is suggested. In our notation

---

<sup>12</sup> The term "Residual  $IJ$ " is assigned to the component due to differences in  $IJ$ -specific rates, since it measures the difference between the total rates of the two groups *after*  $IJ$ -composition is held constant, while the crude rate difference measures the difference between total rates of the two groups *without* taking into account differences in  $IJ$ -composition.

<sup>13</sup> This statement is, of course, limited to the three populations discussed in the previous section as possible standards for computing standardized rates—namely,  $p$ ,  $P$  or a third "observed" population,  $p'$ . We have not yet discussed the use of certain "hypothetical" populations. Also, we have thus far required a set of components to allocate the crude rate difference into two parts only, one due to differences in specific rates and the other to differences in composition.

$$\frac{t_{ij} + T_{ij}}{2} = \text{average } IJ\text{-specific rates of } p \text{ and } P$$

$$\frac{\frac{n_{ij}}{n_{..}} + \frac{N_{ij}}{N_{..}}}{2} = \text{average } IJ\text{-composition of } p \text{ and } P$$

A little algebraic manipulation will show that

$$t_{..} - T_{..} = \text{Combined } IJ + \text{Residual } IJ$$

where

$$\begin{aligned} \text{Combined } IJ &= \sum_i \sum_j \frac{t_{ij} + T_{ij}}{2} \left( \frac{n_{ij}}{n_{..}} - \frac{N_{ij}}{N_{..}} \right) \\ &= \text{component due to differences in } IJ\text{-composition (with} \\ &\quad \text{average } IJ\text{-specific rates of } p \text{ and } P \text{ as weights)} \end{aligned}$$

$$\begin{aligned} \text{Residual } IJ &= \sum_i \sum_j \frac{\frac{n_{ij}}{n_{..}} + \frac{N_{ij}}{N_{..}}}{2} (t_{ij} - T_{ij}) \\ &= \text{component due to differences in } IJ\text{-specific rates (with} \\ &\quad \text{average } IJ\text{-composition of } p \text{ and } P \text{ as weights)} \\ &= \text{difference between } IJ\text{-standardized rates of } p \text{ and } P \\ &\quad \text{(with average } IJ\text{-composition of } p \text{ and } P \text{ as the} \\ &\quad \text{standard)} \end{aligned}$$

Thus, a standard population having the average *IJ*-composition and the average *IJ*-specific rates of *p* and *P* does allocate the difference between their crude rates into two major components, one reflecting differences in their *IJ*-composition and the other differences in their *IJ*-specific rates.<sup>14</sup>

---

<sup>14</sup> It may be noted that the use of a weighted average will not yield the symmetric results of the simple average. For example, if the specific rates of *p* are assigned a weight of 2 and those of *P* a weight of 1 to obtain weighted average specific rates as weights for the "composition component," then in the weights for the "rates component" the composition of *P* will have a weight of 2 and the composition of *p* a weight of 1. That is

$$\begin{aligned} \text{If Combined } IJ &= \sum_i \sum_j \frac{2t_{ij} + T_{ij}}{3} \left( \frac{n_{ij}}{n_{..}} - \frac{N_{ij}}{N_{..}} \right). \\ \text{Then Residual } IJ &= \sum_i \sum_j \frac{\frac{n_{ij}}{n_{..}} + 2 \frac{N_{ij}}{N_{..}}}{3} (t_{ij} - T_{ij}). \end{aligned}$$

This solution is proposed here as the most meaningful one for general purposes when the components framework is used. However, an alternative three-component solution may be considered, especially for certain types of comparisons. For example, if the two crude rates refer to the same population at two different dates, the following questions might be asked: (1) How much change would there have been in the crude rate between the two dates if the  $IJ$ -composition of the population changed as it did but the  $IJ$ -specific rates had remained constant (as of the earlier date)? (2) How much change would there have been in the crude rate if the  $IJ$ -specific rates changed as they did, but the  $IJ$ -composition had remained constant (as of the earlier date)? (3) If the changes measured in (1) and (2) do not add to the total change in the crude rate, by how much do they fail to do so? If we let  $p$  represent the population at the later date, and  $P$  the population at the earlier date, the total change in the crude rate between the two dates may be expressed as the sum of three components, as follows:

$$t_{..} - T_{..} = \sum_i \sum_j T_{ij} \left( \frac{n_{ij}}{n_{..}} - \frac{N_{ij}}{N_{..}} \right) + \sum_i \sum_j \frac{N_{ij}}{N_{..}} (t_{ij} - T_{ij}) \\ + \sum_i \sum_j (t_{ij} - T_{ij}) \left( \frac{n_{ij}}{n_{..}} - \frac{N_{ij}}{N_{..}} \right).$$

The first (or composition) component measures changes in  $IJ$ -composition assuming no change in  $IJ$ -specific rates; the second (or rates) component measures changes in  $IJ$ -specific rates assuming no change in  $IJ$ -composition; and the third (or interaction) component involves changes in both  $IJ$ -composition and  $IJ$ -specific rates.<sup>15</sup>

An exchange of comments with several readers of this paper revealed some differences in preference for the two- and three-component solutions, particularly when the two crude rates refer to the same population at different dates (or to any comparison where one set of events may be considered to precede another). In such a comparison, it may be argued, a three-component solution, with the initial population providing the weights for both the "rates" and "composition" components, appears to be a logical approach. The interpretation of components in this case is described above.

On the other hand, it is the writer's opinion that a good case can be

---

<sup>15</sup> It may be noted that these three components reduce to a two-component solution in either of two ways. If the first and third components are added, the result is a composition component with the specific rates of  $p$  as weights. If the second and third components are added, the result is a rates component with the composition of  $p$  as weights.

made for the two-component solution in such situations. A brief discussion of the reasons for this preference may clarify the rationale of both solutions. First, the selection of a standard population for a components analysis of the change in a crude rate between two dates can be made from any of several assumptions. The three-component solution derives a rates component assuming no change in composition, a composition component assuming no change in specific rates, and an interaction component reflecting changes in both rates and composition. However, changes in rates and composition are seldom independent—rather, a change in one is likely to affect the other. It may be argued, therefore, that since both were changing during the period, a logical set of weights for summarizing changes in specific rates, for example, would be the average composition of the population during the period; similarly, the weights for the composition component might be the average specific rates experienced by the population throughout the period. The two-component solution uses averages for the two dates as weights for these components. While such averages are not equivalent to averages for regular intervals throughout the period, they are often the only averages which can be obtained; also, even when data are available, the computation of annual averages, or averages for any other regular interval, would be a very laborious task. It may be noted in this connection that the simple average for the beginning and end of the period will equal the annual (or other regular interval) average if it is assumed that changes were distributed uniformly throughout the period.

Second, the selection of a standard population for a components analysis of crude rates for two dates is in many respects comparable to the problem of selecting weights for index numbers. Economists give careful consideration to alternative sets of weights when computing index numbers for two dates, and an average of appropriate values for the two years is frequently used.<sup>16</sup> For example, to compute an index of the physical volume of manufacturing production for 1947 relative to 1939, the Bureau of the Census defined change in physical volume as the “change in value of net output, or value added [by manufacture], at constant prices,” and used the average prices for the two years as weights to be applied to the quantities of each product included in the index.<sup>17</sup> An advantage of the use of average prices, as compared with prices for the first year, is that it avoids overweighting products whose

---

<sup>16</sup> The average used in Fisher's “ideal” index is the geometric mean, while the simple arithmetic mean is used in the Marshall-Edgeworth formula.

<sup>17</sup> Bureau of the Census, *Census of Manufactures: 1947, Indexes of Production*, pp. 2-4.

prices have shown a considerable "relative" decrease between the two years, and also avoids underweighting products whose prices have shown a considerable "relative" increase during the period.

In the general case—when the crude rates refer to two different groups, neither of which may be considered to precede the other—there is no "initial state" from which to measure change and, therefore, less emphasis on the desirability of using the rates and composition of one group as weights for both components. An unambiguous allocation into two components, using the average rates and average composition as weights, is a possible solution, and appears more logical than the various alternatives considered.<sup>18</sup>

SUBCOMPONENTS OF COMBINED *IJ*

Combined *IJ*, the composition component of the two-component solution, may be further divided into three subcomponents as follows:<sup>19</sup>

Combined *IJ* = Net *I<sub>J</sub>* + Net *J<sub>I</sub>* + Joint *IJ* where

$$\begin{aligned}
 \text{Net } I_J &= \sum_i \sum_j \left[ \left( \frac{t_{ij} + T_{ij}}{2} \right) \left( \frac{\frac{n_{\cdot j}}{n_{\cdot\cdot}} + \frac{N_{\cdot j}}{N_{\cdot\cdot}}}{2} \right) \right] \left( \frac{n_{ij}}{n_{\cdot j}} - \frac{N_{ij}}{N_{\cdot j}} \right) \\
 \text{Net } J_I &= \sum_i \sum_j \left[ \left( \frac{t_{ij} + T_{ij}}{2} \right) \left( \frac{\frac{n_{i\cdot}}{n_{\cdot\cdot}} + \frac{N_{i\cdot}}{N_{\cdot\cdot}}}{2} \right) \right] \left( \frac{n_{ij}}{n_{i\cdot}} - \frac{N_{ij}}{N_{i\cdot}} \right) \\
 \text{Joint } IJ &= \sum_i \sum_j \left[ \frac{t_{ij} + T_{ij}}{2} \right] \\
 &\quad \frac{\left( \frac{N_{ij}}{N_{i\cdot}} \frac{n_{i\cdot}}{n_{\cdot\cdot}} - \frac{n_{ij}}{n_{i\cdot}} \frac{N_{i\cdot}}{N_{\cdot\cdot}} \right) + \left( \frac{N_{ij}}{N_{\cdot j}} \frac{n_{\cdot j}}{n_{\cdot\cdot}} - \frac{n_{ij}}{n_{\cdot j}} \frac{N_{\cdot j}}{N_{\cdot\cdot}} \right)}{2}
 \end{aligned}$$

<sup>18</sup> If it is desired, a range of variation for each component, as a result of the different weights which might be used in its computation, might be estimated by computing each component twice—once with one of the two populations providing the weights, and a second time with the other population providing the weights. Thus, the two values for each component defined by the two equations

$$t_{\cdot\cdot} - T_{\cdot\cdot} = \sum_i \sum_j T_{ij} \left( \frac{n_{ij}}{n_{\cdot\cdot}} - \frac{N_{ij}}{N_{\cdot\cdot}} \right) + \sum_i \sum_j \frac{n_{ij}}{n_{\cdot\cdot}} (t_{ij} - T_{ij})$$

and

$$t_{\cdot\cdot} - T_{\cdot\cdot} = \sum_i \sum_j t_{ij} \left( \frac{n_{ij}}{n_{\cdot\cdot}} - \frac{N_{ij}}{N_{\cdot\cdot}} \right) + \sum_i \sum_j \frac{N_{ij}}{N_{\cdot\cdot}} (t_{ij} - T_{ij})$$

might be used to estimate the range of variation for each component. It may be noted that each major component in the two-component solution proposed in this paper is an average of the two values obtained for this component from these equations.

In these formulas,<sup>20</sup>  $n_{i.}$  and  $N_{i.}$  represent the total number of persons in the  $i$ th category of factor  $I$  in groups  $p$  and  $P$ , respectively; and  $n_{.j}$  and  $N_{.j}$  represent the total number of persons in the  $j$ th category of factor  $J$  in  $p$  and  $P$ , respectively.

Net  $I_J$  is a weighted sum of differences in  $I_J$ -composition ( $I$ -composition *within* the  $J$  subgroups)<sup>21</sup> of  $p$  and  $P$  and, therefore, may be interpreted as that part of the difference between their crude rates which is attributable to differences in net  $I_J$ -composition, or  $I$ -composition *independent* of  $J$ . The average  $IJ$ -specific rates and average gross  $J$ -composition of  $p$  and  $P$  are used as weights in this subcomponent.<sup>22</sup> That is, Net  $I_J$  measures differences in  $I_J$ -composition applied to a standard population which has the average  $IJ$ -specific rates and the average gross  $J$ -composition of  $p$  and  $P$ .

Similarly Net  $J_I$  is a weighted sum of differences in  $J_I$ -composition ( $J$ -composition *within* the  $I$  subgroups) and may be interpreted as that part of the crude rate difference which is due to differences in net  $J_I$ -composition, or  $J$ -composition *independent* of  $I$ . The average  $IJ$ -specific rates and average gross  $I$ -composition of  $p$  and  $P$  are used as weights for this subcomponent. That is, Net  $J_I$  measures differences in  $J_I$ -composition applied to a standard population which has the average  $IJ$ -specific rates and the average gross  $I$ -composition of  $p$  and  $P$ .

Joint  $IJ$  is that part of the Combined  $IJ$  component which cannot be allocated to differences in net  $I_J$ -composition or to differences in net  $J_I$ -composition. It represents the part of the crude rate difference which is accounted for by differences in combined  $IJ$ -composition but which cannot be allocated independently to  $I$  or  $J$ . Its equation shows it to be a weighted sum of differences in  $IJ$ -composition which result from combining the net composition of one of the two populations with the gross composition of the other,<sup>23</sup> using the average  $IJ$ -specific rates of the two populations as weights.

<sup>19</sup> Because the composition component of the three-component solution is not a total measure of differences in  $IJ$ -composition—part of the composition differences are included in the interaction component—there is no discussion of its subcomponents.

<sup>20</sup> Note that  $n_{i.} = \sum_j n_{ij}$ ,  $n_{.j} = \sum_i n_{ij}$ , etc.

<sup>21</sup> In the equation for this subcomponent,  $(n_{ij}/n_{.j} - N_{ij}/N_{.j})$  represents the differences in  $I$ -composition *within* the subgroups of factor  $J$ .

<sup>22</sup> The weights are enclosed in brackets, [ ], in the equations defining the subcomponents. Gross  $J$ -composition refers to the per cent distribution of a group when it is classified by factor  $J$  only. For example, in group  $p$  gross  $J$ -composition is defined by  $n_{.j}/n_{..}$

<sup>23</sup> Specifically, the term

$$\frac{N_{ij}}{N_{i.}} \frac{n_{i.}}{n_{..}} - \frac{n_{ij}}{n_{i.}} \frac{N_{i.}}{N_{..}}$$

expresses the difference between (1) the net  $J_I$ -composition of  $P$  combined with the gross  $I$ -composition of  $p$ , and (2) the net  $J_I$ -composition of  $p$  combined with the gross  $I$ -composition of  $P$ . Similarly, the term

Thus, a complete components analysis allocates the difference between the crude rates of  $p$  and  $P$  into four parts:

$$t.. - T.. = \text{Net } I_J + \text{Net } J_I + \text{Joint } IJ + \text{Residual } IJ$$

where Net  $I_J$  is the part due to differences in  $I$ -composition independent of  $J$ ; Net  $J_I$  is due to differences in  $J$ -composition independent of  $I$ ; Joint  $IJ$  is due to differences in Joint  $IJ$  composition, or to differences in Combined  $IJ$  composition which cannot be allocated independently to  $I$  or  $J$ ; and Residual  $IJ$  is due to differences in  $IJ$ -specific rates of  $p$  and  $P$ . The standard population used for this purpose is one having the average  $IJ$ -composition and the average  $IJ$ -specific rates of  $p$  and  $P$ .

#### ONE FACTOR COMPONENTS

If data for the two groups,  $p$  and  $P$ , are cross-classified by only one factor,  $I$ , a two-component allocation of the difference between their crude rates may be defined as follows:

$$t. - T. = \text{Gross } I + \text{Residual } I$$

where

$$\begin{aligned} \text{Gross } I &= \sum_i \frac{t_i + T_i}{2} \left( \frac{n_i}{n.} - \frac{N_i}{N.} \right) \\ \text{Residual } I &= \sum \frac{\frac{n_i}{n.} + \frac{N_i}{N.}}{2} (t_i - T_i). \end{aligned}$$

The Gross  $I$  component represents that part of the difference between the crude rates of  $p$  and  $P$  which is due to differences in their  $I$ -composition, and Residual  $I$  the part due to differences in their  $I$ -specific rates. Residual  $I$  is also equal to the difference between the

---


$$\frac{N_{ij} n_{.j}}{N_{.j} n..} - \frac{n_{ij} N_{.j}}{n_{.j} N..}$$

expresses the difference between (1) the net  $I_J$ -composition of  $P$  combined with the gross  $J$ -composition of  $p$ , and (2) the net  $I_J$ -composition of  $p$  combined with the gross  $J$ -composition of  $P$ . These two sets of differences—one for the combination of net  $I_J$  and gross  $J$  compositions, the other for the combination of net  $J_I$  and gross  $I$  compositions—are obtained because the two factors  $I$  and  $J$ , are interchangeable in this subcomponent. Therefore, corresponding differences (that is, differences for the same  $IJ$  cell of the cross-classification) in these two sets are averaged, and the resulting set of average differences are weighted by average  $IJ$ -specific rates of  $p$  and  $P$ . This interpretation of the Joint  $IJ$  subcomponent is presented to state explicitly the differences which are measured. It is not necessary to compute these differences to obtain this subcomponent since it may be obtained more simply by subtracting the sum of subcomponents Net  $I_J$  and Net  $J_I$  from the major component Combined  $IJ$ .



*I*-standardized rates of *p* and *P*, with the average *I*-composition of *p* and *P* as the standard population.

A three-component allocation similar to that described for two factors, is as follows:

$$t. - T. = \sum_i T_i \left( \frac{n_i}{n.} - \frac{N_i}{N.} \right) + \sum_i \frac{N_i}{N.} (t_i - T_i) + \sum_i (t_i - T_i) \left( \frac{n_i}{n.} - \frac{N_i}{N.} \right).$$

In this case, the *I*-specific rates and *I*-composition of group *P* are used as weights for the “composition” and “rates” components, respectively. The third component is an “interaction” component which is due to differences in both composition and rates.

INTERPRETATION OF RESULTS

We shall first consider the two-component solution. The interpretation of the two major components and the subcomponents is relatively simple when all the components are positive, since the difference in crude rates can then be attributed partly to differences in *IJ*-specific rates and partly to differences in *IJ*-composition, *I<sub>J</sub>*-composition, *J<sub>I</sub>*-composition and Joint *IJ*-composition. In this situation, the components may be converted to per cent components, with the crude rate difference as the base. For example,

$$\frac{\text{Combined } IJ}{t.. - T..} \times 100 = \text{per cent of difference between crude rates which is due to differences in } IJ\text{-composition of } p \text{ and } P.$$

Even in this case, however, two qualifications should be kept in mind: (1) the components include the effects of hidden forces behind the factors, *I* and *J*, and nothing in the components technique justifies inferences as to causal relationships—such inferences must be based on knowledge outside the statistical technique itself; (2) a small Residual component may mask larger influences, in opposite directions, of factors not held constant in the analysis. The latter qualification means, for example, that Residual *IJ* may be less than Residual *IJK*; that is, the difference in *IJK*-standardized rates may be greater than the difference in *IJ*-standardized rates. Thus, the results of a components analysis should be interpreted as “applicable within the frame-

work of the particular factors held constant." For example, a Residual *IJ* per cent component of 35 per cent may correctly be interpreted as indicating that "the difference in *IJ*-standardized rates is only 35 per cent as great as the difference in crude rates," and the corresponding Combined *IJ* component of 65 per cent may be interpreted as indicating that "65 per cent of the crude difference is attributable to *IJ*-composition." But, as additional factors (*K*, *L*, etc.) are held constant, the Residual component may not always decrease (and the Combined component increase) with the addition of each new factor—in fact, it will not do so unless all the other factors (*K*, *L*, etc.) operate in the same direction. This characteristic simply points to the fact that the difference between two crude rates is not the equivalent of a concept like total variance of a dependent variable in regression analysis, for example, which will be increasingly "explained" as more independent variables are added to the regression equation.<sup>24</sup>

When not all of the components are positive, their interpretation is more complicated. For example, if Residual *IJ*—the difference between *IJ*-standardized rates—is greater than the difference between the crude rates, Combined *IJ* is negative. If all of the subcomponents have the same sign (that is, are negative), Residual *IJ* might be used as the base for per cent components, which can be interpreted as follows:

$\frac{t.. - T..}{\text{Residual } IJ} \times 100 =$  per cent of the difference between *IJ*-standardized rates which is observed or evident in the difference between their crude rates.

$\frac{- \text{Combined } IJ}{\text{Residual } IJ} \times 100 =$  per cent of the difference between *IJ*-standardized rates which is obscured (in the crude rates) by differences in *IJ*-composition. (This is a positive component because Combined *IJ* is negative.)

$\frac{- \text{Net } I_J}{\text{Residual } IJ} \times 100 =$  per cent of the difference between *IJ*-standardized rates obscured by differences in *I\_J*-composition, etc.

Other combinations of positive and negative components do not

---

<sup>24</sup> Strictly speaking, the addition of a new independent factor need not increase the amount of "explained variance," but it cannot decrease it.

readily lend themselves to meaningful per cent components. In such situations, the components themselves may be used without conversion to per cents. For example, if  $t.. - T.. = 2$ , Residual  $IJ = -25$ , and Combined  $IJ = 27$ , we might simply say that while the crude rate of  $p$  exceeded that of  $P$  by 2 points, the standardized rate of  $P$  exceeded that of  $p$  by 25 points, and that differences in  $IJ$ -composition were responsible for these widely different results. The Net  $I_J$ , Net  $J_I$ , and Joint  $IJ$  subcomponents may be used to indicate the importance of  $I$ -composition independent of  $J$ ,  $J$ -composition independent of  $I$ , and Joint  $IJ$ -composition in obscuring (in the crude rates) the difference between standardized rates.

If a crude rate difference is allocated into three major components, and all are positive, each component might be expressed as a per cent of the crude rate difference. Or, if this solution is used in a comparison over time, it may be desired to use the crude rate of the earlier date as the per cent base. That is, we may write these components as

$$t.. = T.. + \sum_i \sum_j T_{ij} \left( \frac{n_{ij}}{n..} - \frac{N_{ij}}{N..} \right) + \sum_i \sum_j \frac{N_{ij}}{N..} (t_{ij} - T_{ij}) \\ + \sum_i \sum_j (t_{ij} - T_{ij}) \left( \frac{n_{ij}}{n..} - \frac{N_{ij}}{N..} \right).$$

If each of the terms in this equation is expressed as a per cent of  $T..$ , the crude rate at the earlier date, then the per cent which  $t..$  is of  $T..$  is equal to 100 plus the three percentage components. In this case, for example,

$$\frac{\sum_i \sum_j T_{ij} \left( \frac{n_{ij}}{n..} - \frac{N_{ij}}{N..} \right)}{T..} \times 100 = \text{per cent change in the total rate between} \\ \text{the two dates as a result of changes in} \\ \text{composition, assuming no change in} \\ \text{specific rates during the period.}$$

#### AN EXAMPLE OF RESULTS

In a study of labor mobility in six cities made in 1951,<sup>25</sup> this method was used to determine the extent to which city differences in job mo-

<sup>25</sup> This research, conducted by the Chicago Community Inventory of the University of Chicago, was based on data obtained in the Six-City Mobility Study, one of the industrial manpower research studies sponsored by the United States Air Force under Project SCOOP. Findings are summarized in Evelyn M. Kitagawa, "Relative Importance and Independence of Selected Factors in Job Mobility, Six Cities, 1940-49," *op. cit.*

bility rates were due to differences in the composition of the labor force in the various cities. To cite one particular example, the crude mobility rate (mean number of jobs held, 1940-49) of Los Angeles men was 32 per cent higher than that of Philadelphia men, or an average of 3.14 jobs as compared with 2.37 in Philadelphia.

Table 1 presents data on mobility rates and composition, by migrant status and time spent in the labor force from 1940 to 1949, for men in these two cities. In the lower half of the table are the results of a components analysis of the difference between crude mobility rates of Los Angeles and Philadelphia men with these two factors held constant—migrant status (*J*) and time spent in labor force (*I*).

Examination of the specific mobility rates shows that Los Angeles men were consistently more mobile than Philadelphia men, even when data are cross-classified by migrant status and time in the labor force. However, it is also clear that migrants were more mobile than non-migrants, and persons in the labor force 5-9½ years were more mobile than persons in for less or more time. Furthermore, the percentage distributions which describe the composition of men in the two cities indicate higher proportions of Los Angeles men in the high mobility categories—for example, 47 per cent of the Los Angeles men were migrants as compared with 13 per cent of the Philadelphia men. Thus, we expect that part, but not all, of the difference between crude rates for men in these two cities is due to differences in their composition with respect to migrant status and time spent in the labor force. The components analysis quantifies this relationship.

Differences in composition with respect to both migrant status and time spent in the labor force account for 47 per cent of the difference between crude rates of Los Angeles and Philadelphia men. And, the difference between *IJ*-standardized rates for men in these two cities is 53 per cent of their crude rate difference, using their average *IJ*-composition as standard (Residual *IJ*).

Furthermore, migrant composition alone, independent of time spent in the labor force, accounted for 38 per cent of the crude rate difference, with only 1 per cent due to composition by time spent in the labor force independent of migrant composition, and 7 per cent to these two factors jointly.

Similar components were computed for selected pairs of cities, with these and other factors held constant, to determine which factors were most important in accounting for city differentials in crude mobility rates and to what extent these city differentials reflected differences in specific (or standardized) mobility rates. Space does not permit a more detailed analysis, but the purpose here is only to illustrate the use of the method in a specific set of data.

Measures of the sampling variability of per cent components have not been determined. Since each component is based on the difference between two sums of a large number of products, its sampling variance may prove to be too unwieldy to estimate, though further study of this problem might yield some simplifying assumptions which will furnish approximate estimates of sampling variance without overburdening computations.<sup>26</sup> With census data based on complete counts or with very large samples where cross-classifications do not run thin, per cent components should be relatively stable or reliable. But with small samples or with small cell frequencies in complete counts, inferences should be made with caution.

#### COMPARISON WITH PREVIOUS DEFINITIONS OF COMPONENTS

Reference was made in the introductory paragraphs to previous work involving the allocation of a crude rate difference into components. The set of components proposed in this paper differs from those used previously in two respects. First, the earlier approaches selected one of the two populations being compared to provide weights for one of the two major components, considered this population the standard population for the components analysis, and did not make explicit the weights used in the other major component which was obtained as a residual (by subtracting the computed component from the crude rate difference, or an equivalent procedure). The algebraic presentation in this paper has made explicit the set of weights which was implicit in such a two-component solution; for example, when the  $IJ$ -specific rates of one group are used to weight differences in  $IJ$ -composition, the  $IJ$ -composition of the other group is used to weight differences in  $IJ$ -specific rates. The two-component solution proposed here uses a standard population having the average  $IJ$ -composition and the average  $IJ$ -specific rates of the two groups.

Second, the rationale of a set of subcomponents in this paper is more defensible than that in the previous literature.<sup>27</sup> This can best be seen by noting that earlier definitions of the Joint  $IJ$  subcomponent, if

<sup>26</sup> One place in the analysis where tests of significance can be made is in the determination of whether one population's specific rates are *on the whole* larger than another's. The test consists of considering each sub-group in one population with its comparable sub-group in the other population as a four-fold table. One can then sum the chi-squares from the individual four-fold tables. (See Karl Pearson and J. F. Tocher, "On Criteria for the Existence of Differential Death Rates," *Biometrika*, 11 (1916), 159-64; S. A. Stouffer and Clark Tibbitts, "Tests of Significance in Applying Westergaard's Method of Expected Cases to Sociological Data," *Journal of the American Statistical Association*, 28 (1933), 293-302; H. F. Dorn and S. A. Stouffer, "Criteria of Differential Mortality," *Journal of the American Statistical Association*, 28 (1933), 402-13).

<sup>27</sup> To the writer's knowledge, the only previous work involving subcomponents is contained in the cited reference to Goldfield's method and the Turner and Kitagawa references. The comments here are applicable, for the most part, to the Turner and Kitagawa definitions; although Goldfield's definition of a Joint subcomponent was similar, he did not retain it as a measure of Joint composition but allocated it back to the factors involved.

TABLE 1

JOB MOBILITY RATES (MEAN NUMBER OF JOBS HELD 1940-49) AND COMPOSITION (PERCENTAGE DISTRIBUTION) BY MIGRANT STATUS AND TIME SPENT IN THE LABOR FORCE, FOR LOS ANGELES AND PHILADELPHIA MEN: 1940-49\*

Migrant Status and Time in Labor Force 1940-49	Mobility Rates		Composition (%)		Difference (L.A. Minus Phila.)	
	Los Angeles	Philadelphia	Los Angeles	Philadelphia	Rates	Composition
<i>All Men</i>	3.14	2.37	100	100	.77	0
Less than 5 yrs.	2.90	2.42	11	7	.48	4
5 but less than 9½ yrs.	3.82	3.26	30	26	.56	4
9½-10 yrs.	2.84	2.03	59	67	.81	-8
<i>Migrants</i>	3.77	3.13	47	13	.64	34
Less than 5 yrs.	2.89	2.29	6	1	.60	5
5 but less than 9½ yrs.	4.07	3.43	17	4	.64	13
9½-10 yrs.	3.79	3.15	24	8	.64	16
<i>Non-migrants</i>	2.58	2.26	53	87	.32	-34
Less than 5 yrs.	2.92	2.45	5	6	.47	-1
5 but less than 9½ yrs.	3.49	3.23	13	22	.26	-9
9½-10 yrs.	2.20	1.88	35	59	.32	-24

*Components of the Difference Between the Crude Job Mobility Rates of Los Angeles Men (3.14) and San Francisco Men (2.37)*

(Two Factors: *I* = time spent in labor force 1940-49; *J* = migrant status)

Name of Component	Per Cent Component ( $t.. - T.. = 100$ )	Actual Component ( $t.. - T.. = .77$ )
Combined <i>IJ</i>	47	.359
Net <i>J<sub>I</sub></i>	38	.296
Net <i>I<sub>J</sub></i>	1	.008
Joint <i>IJ</i>	7	.054
Residual <i>IJ</i>	53	.411

\* Data refer to a probability sample of 1,313 men in Los Angeles and 1,571 men in Philadelphia who worked one month or more in 1950. Persons residing in each city in 1951 who had resided there 11 years or less (that is, who moved there after the beginning of the 1940-49 decade) were classified as migrants for purposes of this study.

translated into the algebraic notation, do not reduce to a weighted sum of differences in composition resulting from a combination of the net composition of one group with the gross composition of the other. In previous definitions the Net  $I_J$  and Net  $J_I$  subcomponents were defined independently, and then the Joint  $IJ$  subcomponent was computed by subtracting Net  $I_J$  and Net  $J_I$  from Combined  $IJ$ , without determining, algebraically, what such a residual measured. For example, when the specific rates of  $P$  were used as weights for the Combined  $IJ$  component in previous definitions, the subcomponents were defined as follows, if translated into our notation:

$$\begin{aligned} \text{Net } I_J &= \sum_i \sum_j T_{ij} \frac{n_{.j}}{n_{..}} \left( \frac{n_{ij}}{n_{.j}} - \frac{N_{ij}}{N_{.j}} \right) \\ \text{Net } J_I &= \sum_i \sum_j T_{ij} \frac{n_{i.}}{n_{..}} \left( \frac{n_{ij}}{n_{i.}} - \frac{N_{ij}}{N_{i.}} \right) \\ \text{Joint } IJ &= \sum_i \sum_j T_{ij} \left( \frac{N_{ij}}{N_{i.}} \frac{n_{i.}}{n_{..}} - \frac{N_{ij}}{N_{..}} - \frac{n_{ij}}{n_{..}} + \frac{N_{ij}}{N_{.j}} \frac{n_{.j}}{n_{..}} \right). \end{aligned}$$

However, if the rationale of a set of subcomponents in this paper is used to compute subcomponents for the same Combined  $IJ$  component

$$\left( \text{Combined } IJ = \sum_i \sum_j T_{ij} \left[ \frac{n_{ij}}{n_{..}} - \frac{N_{ij}}{N_{..}} \right] \right)^{28}$$

the results could be  
*either*

$$\begin{aligned} \text{Net } I_J &= \sum_i \sum_j T_{ij} \frac{n_{.j}}{n_{..}} \left( \frac{n_{ij}}{n_{.j}} - \frac{N_{ij}}{N_{.j}} \right) \\ \text{Net } J_I &= \sum_i \sum_j T_{ij} \frac{N_{i.}}{N_{..}} \left( \frac{n_{ij}}{n_{i.}} - \frac{N_{ij}}{N_{i.}} \right) \\ \text{Joint } IJ &= \sum_i \sum_j T_{ij} \left( \frac{N_{ij}}{N_{.j}} \frac{n_{.j}}{n_{..}} - \frac{n_{ij}}{n_{i.}} \frac{N_{i.}}{N_{..}} \right) \end{aligned}$$

---

<sup>28</sup> Two similar sets of subcomponents could be obtained if the rates of  $p$  (instead of  $P$ ) were used as weights in the Combined  $IJ$  component; equations for these two sets can be written by interchanging the roles of  $p$  and  $P$  in the weights of the equations above. Thus, four sets of subcomponents might be defined for the two Combined  $IJ$  components. Averaging the four values for each subcomponent would give a single set of subcomponents which is identical to the set defined in this paper.

or

$$\text{Net } I_J = \sum_i \sum_j T_{ij} \frac{N_{\cdot j}}{N_{\cdot\cdot}} \left( \frac{n_{ij}}{n_{\cdot j}} - \frac{N_{ij}}{N_{\cdot j}} \right)$$

$$\text{Net } J_I = \sum_i \sum_j T_{ij} \frac{n_{i\cdot}}{n_{\cdot\cdot}} \left( \frac{n_{ij}}{n_{i\cdot}} - \frac{N_{ij}}{N_{i\cdot}} \right)$$

$$\text{Joint } IJ = \sum_i \sum_j T_{ij} \left( \frac{N_{ij}}{N_{i\cdot}} \frac{n_{i\cdot}}{n_{\cdot\cdot}} - \frac{n_{ij}}{n_{\cdot j}} \frac{N_{\cdot j}}{N_{\cdot\cdot}} \right).$$

#### COMPARISON WITH WESTERGAARD'S METHOD OF EXPECTED CASES

The logic and mechanics of Westergaard's method of expected cases is described by Woodbury in a study of infant mortality rates.<sup>29</sup> Jaffe describes the method as "essentially an elaboration of the conventional standardization technique which can be used in a situation where the investigator wishes to isolate the influence of a single factor from that of other associated factors."<sup>30</sup>

Basically, the method involves the computing of ratios of the actual number of cases in a particular group to the expected number of cases in the same group should it retain its own composition while being exposed to the specific rates of a standard population (with the total of the various subgroups usually used as standard). For example, Woodbury in the study cited above, uses the method to isolate the influence of order of birth on infant mortality holding constant age of mother and earnings of father. If we let factor *I* represent age of mother and factor *J* earnings of father, then in our notation we might say

$n_{ij}^{(K)}$  = number of births in both the *i*th category of *I* and the *j*th category of *J* in the population of *k*-order births

$N_{ij}$  = number of births in both the *i*th category of *I* and the *j*th category of *J* in the population of total births

$t_{ij}^{(K)}$  = infant mortality rates for births in the *i*th category of *I* and the *j*th category of *J* in the population of *k*-order births

$T_{ij}$  = infant mortality rates for births in the *i*th category of *I* and the *j*th category of *J* in the population of total births

His results are summarized in a table similar to Table 2. A comparison of the Westergaard ratios in the second column is used to indicate the extent of the variation in infant mortality by order of birth with the

<sup>29</sup> R. M. Woodbury, "Westergaard's Method of Expected Deaths as Applied to the Study of Infant Mortality," *Journal of American Statistical Association*, 18 (1923), 366-76, and reproduced in Jaffe, *op. cit.*, Chap. III.

<sup>30</sup> Jaffe, *op. cit.*, p. 48.



influence of factors *I* and *J* eliminated. Also, comparison of the ratios in the first and second columns (the former measure crude or unadjusted birth-order differentials in infant mortality) indicates the influence of factors *I* and *J*—that is, *IJ*-composition—on birth-order differentials in infant mortality rates. However, the Westergaard technique does not quantify this aspect of the analysis—that is, it does not measure what part of the crude birth-order differentials is attributable to factors *I* and *J*.

If we compare the Westergaard framework with conventional standardization and the components analysis, several relationships become evident. First, the components framework applied to the analysis of birth-order differentials in mortality might be used to analyze the difference between crude infant mortality rates of any two birth orders into components due to differences in *IJ*-composition on one hand and to differences in *IJ*-specific rates on the other hand. That is to say, it quantifies a relationship which is apparent in a Westergaard analysis

TABLE 2

Order of Birth	Ratio of Original Rate to Average Rate (Col. 1)	Ratio of Actual to Expected <i>I</i> & <i>J</i> Constant (Col. 2)
Total	$\frac{T_{..}}{T_{..}} = 1.00$	$\frac{\sum_i \sum_j T_{ij} N_{ij}}{\sum_i \sum_j T_{ij} N_{ij}} = 1.00 \left( = \frac{T_{..}}{T_{..}} \right)$
First	—	—
—	—	—
—	—	—
<i>k</i> th	$\frac{t_{..(K)}}{T_{..}} = \frac{\sum_i \sum_j t_{ij(K)} \frac{n_{ij(K)}}{n_{..(K)}}}{\sum_i \sum_j T_{ij} \frac{N_{ij}}{N_{..}}}$	$\frac{t_{..(K)} n_{..(K)}}{\sum_i \sum_j T_{ij} n_{ij(K)}} = \frac{\sum_i \sum_j t_{ij(K)} n_{ij(K)}}{\sum_i \sum_j T_{ij} n_{ij(K)}} = \frac{\sum_i \sum_j t_{ij(K)} \frac{n_{ij(K)}}{n_{..(K)}}}{\sum_i \sum_j T_{ij} \frac{n_{ij(K)}}{n_{..(K)}}}$
—	—	—
—	—	—
—	—	—
Tenth	—	—

but not rigorously measured in terms of the "components" concept defined in this paper.

To the writer's knowledge, there is no published work comparing the Westergaard technique with conventionally standardized rates selected to accomplish the same purpose, although the technique has always been recognized as a "standardization" method. Examination of the formulas for Westergaard's ratios in the table above reveals that these ratios are actually equivalent to ratios of indirectly standardized infant mortality rates for each birth order (with the population of total births as standard) to the crude infant mortality rate for total births (i.e., the standard population). For example, the indirectly standardized rate for  $k$ -order births, with total births as the standard population, is given by<sup>31</sup>

$$\frac{\text{actual infant deaths to } k\text{-order births}}{\text{expected infant deaths to } k\text{-order births}} \times \left( \begin{array}{l} \text{crude death rate of} \\ \text{standard population} \end{array} \right)$$

or

$$\frac{\sum_i \sum_j t_{ij}^{(K)} n_{ij}^{(K)}}{\sum_i \sum_j T_{ij} n_{ij}^{(K)}} T..$$

And if this is divided by the crude rate ( $T..$ ) for total births we get the Westergaard ratio.<sup>32</sup>

Thus, if the objective is to compare the relative incidence of mortality from one birth order group to the next, holding constant the disproportionate weighting of the various groups by age of mother and earnings of father, the Westergaard ratios might be considered as approximations in the same sense that indirect standardization approximates the results of direct standardization. Ratios based on directly standardized rates,<sup>33</sup> with the population of total births as standard, would be defined by

$$\frac{\sum_i \sum_j t_{ij}^{(K)} \frac{N_{ij}}{N..}}{T..} \quad \text{or} \quad \frac{\sum_i \sum_j t_{ij}^{(K)} \frac{N_{ij}}{N..}}{\sum_i \sum_j T_{ij} \frac{N_{ij}}{N..}} \quad \text{or} \quad \frac{\sum_i \sum_j t_{ij}^{(K)} N_{ij}}{\sum_i \sum_j T_{ij} N_{ij}}$$

<sup>31</sup> Jaffe, *op. cit.*, pp. 44-8.

<sup>32</sup> The Westergaard ratio for a category may also be described as the correction factor which is applied to the crude rate of the standard population to obtain the indirectly standardized rate for that category.

<sup>33</sup> Such ratios have been used in at least one study—Donald J. Bogue, *A Methodological Study of Migration and Labor Mobility in Michigan and Ohio in 1947* (Oxford: Scripps Foundation for Research in Population Problems, 1952), p. 64.

The only factor which varies in these ratios for different birth orders are the *IJ*-specific rates of the birth orders, while in the Westergaard ratios the composition which appears in the numerator and denominator of each ratio also varies from one birth order to the next. That is to say, Westergaard ratios do not strictly speaking, hold *IJ*-composition constant from one birth order to the next. However, there may be good reasons for computing ratios of indirectly standardized rates such as the Westergaard ratios.<sup>34</sup>

EXTENSION OF FRAMEWORK TO THREE VARIABLES

Extending the components framework to three variables (*I*, *J* and *K*) is simple for the major components—Combined *IJK* and Residual *IJK*. In this case

$$t... - T... = \text{Combined } IJK + \text{Residual } IJK$$

where

$$\text{Combined } IJK = \sum_i \sum_j \sum_k \left( \frac{t_{ijk} + T_{ijk}}{2} \right) \left( \frac{n_{ijk}}{n...} - \frac{N_{ijk}}{N...} \right)$$

and

$$\text{Residual } IJK = \sum_i \sum_j \sum_k \left( \frac{\frac{n_{ijk}}{n...} + \frac{N_{ijk}}{N...}}{2} \right) (t_{ijk} - T_{ijk}).$$

However, expressions for net and joint subcomponents would be quite complex if the two-factor model were extended to three factors. A much simpler solution results if one is willing to consider two of the three factors as a single factor and apply the two-factor framework. For example, if we are primarily interested in the influence of *I* independent of both *J* and *K*, we might consider *I* as one factor and the cross-classification of *J* by *K* as a second factor—the latter will be denoted by (*JK*). Then, in addition to the combined and residual components defined above, we could define the following subcomponents:

$$\text{Net } I_{JK} = \sum_i \sum_j \sum_k \left( \frac{t_{ijk} + T_{ijk}}{2} \right) \left( \frac{\frac{n_{.jk}}{n...} + \frac{N_{.jk}}{N...}}{2} \right) \cdot \left( \frac{n_{ijk}}{n_{.jk}} - \frac{N_{ijk}}{N_{.jk}} \right)$$

---

<sup>34</sup> Two of the most obvious are: (1) when *IJ*-specific rates are available only for total births and not for each birth order, and (2) when enough of the *IJ*-specific rates for the various birth orders are based on too small numbers for stability. See Peter Cox, *Demography* (London: Cambridge University Press, 1950), Chapter 7, for a statement of weights implicit in indirect standardization.

$$\text{Net } (JK)_I = \sum_i \sum_j \sum_k \left( \frac{t_{ijk} + T_{ijk}}{2} \right) \cdot \left( \frac{\frac{n_{i..}}{n_{...}} + \frac{N_{i..}}{N_{...}}}{2} \right) \left( \frac{n_{ijk}}{n_{i..}} - \frac{N_{ijk}}{N_{i..}} \right)$$

$$\text{Joint } I(JK) = \sum_i \sum_j \sum_k \left( \frac{t_{ijk} + T_{ijk}}{2} \right) \cdot \left( \frac{\frac{N_{ijk}}{N_{i..}} \frac{n_{i..}}{n_{...}} - \frac{n_{ijk}}{n_{.jk}} \frac{N_{.jk}}{N_{...}} + \frac{N_{ijk}}{N_{.jk}} \frac{n_{.jk}}{n_{...}} - \frac{n_{ijk}}{n_{i..}} \frac{N_{i..}}{N_{...}}}{2} \right) \cdot$$

In this situation, Net  $I_{JK}$  represents the part of the difference between crude rates attributable to differences in  $I$ -composition independent of both  $J$  and  $K$ , while Net  $(JK)_I$  measures the part attributable to combined  $JK$ -composition independent of  $I$ .

Such an analysis could be made with any pair of the factors considered as a single factor. Also, if four or more factors are combined in any way to reduce to two, a similar approach could be used.