



Taylor & Francis
Taylor & Francis Group



Simpson's Paradox in Real Life

Author(s): Clifford H. Wagner

Source: *The American Statistician*, Feb., 1982, Vol. 36, No. 1 (Feb., 1982), pp. 46-48

Published by: Taylor & Francis, Ltd. on behalf of the American Statistical Association

Stable URL: <https://www.jstor.org/stable/2684093>

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at <https://about.jstor.org/terms>



JSTOR

Taylor & Francis, Ltd. and American Statistical Association are collaborating with JSTOR to digitize, preserve and extend access to *The American Statistician*

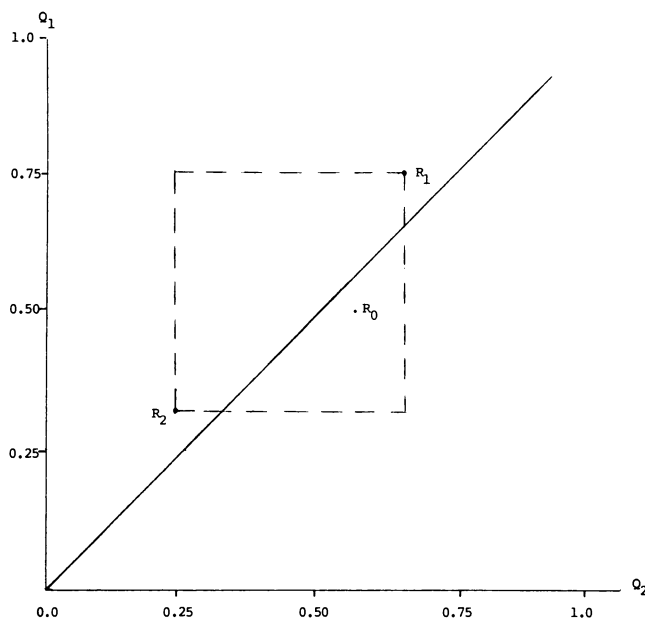


Figure 2. Plot of Table 1—Conditional Probability Coordinates

relevant condition is marginal independence of A and C , that is, independence of A and C when we combine

over B . The distinction reflects that the conditions for confounding differ depending upon the measure of association used.

[Received August 1979. Revised September 1981.]

REFERENCES

- BISHOP, YVONNE M. M.; FIENBERG, STEPHEN E.; and HOLLAND, PAUL W. (1975), *Discrete Multivariate Analysis: Theory and Practice*, Cambridge: MIT Press.
- BLYTH, COLIN R. (1972), "On Simpson's Paradox and the Sure-Thing Principle," *Journal of the American Statistical Association*, 67, 364–366.
- GARDNER, MARTIN (1976), "On the Fabric of Inductive Logic and Some Probability Paradoxes," *Scientific American*, 234, 119–124.
- ROTHMAN, KENNETH J. (1975), "A Pictorial Representation of Confounding in Epidemiologic Studies," *Journal of Chronic Diseases*, 28, 101–108.
- SIMPSON, E. H. (1951), "The Interpretation of Interaction in Contingency Tables," *Journal of the Royal Statistical Society, Ser. B*, 13, 238–241.
- WHITEMORE, ALICE S. (1978), "Collapsibility of Multi-dimensional Contingency Tables," *Journal of the Royal Statistical Society, Ser. B*, 40, 328–340.
- YULE, G. U. (1903), "Notes on the Theory of Association of Attributes in Statistics," *Biometrika*, 2, 121–134.

Simpson's Paradox in Real Life

CLIFFORD H. WAGNER*

1. INTRODUCTION

Simpson's paradox (Blyth 1972) is the designation for a surprising situation that may occur when two populations are compared with respect to the incidence of some attribute: If the populations are separated in parallel into a set of descriptive categories, the population with higher overall incidence may yet exhibit a lower incidence within *each* such category.

An actual occurrence of this paradox was observed (Cohen and Nagel 1934, p. 449) in a comparison of tuberculosis deaths in New York City and Richmond, Virginia, during the year 1910. Although the overall tuberculosis mortality rate was lower in New York, the opposite was observed when the data were separated into two racial categories; in both the white and non-white categories, Richmond had a lower mortality rate. A similar situation involving populations divided into a large number of categories occurred in a well-known study of sex bias in graduate admissions at the University of California, Berkeley (Bickel, Hammel, and

O'Connell 1975). However, this was not a complete instance of Simpson's paradox because, when the data were disaggregated, the overall tendency toward a higher acceptance rate for male applicants was not reversed in each academic department.

Two real-life examples of Simpson's paradox are presented below. They illustrate the paradox in the context of populations composed of several categories and demonstrate how easily the paradox can occur.

2. RENEWAL RATES

Magazine publishers carefully monitor rates of renewal of expiring subscriptions. For example, at *American History Illustrated* in early 1979, the publishers were pleased to note an increase in the overall renewal rate from 51.2 percent in January to 64.1 percent in February. Because renewal rates are highly correlated with established subscription categories, and because one might wish to identify the kinds of subscriptions that account for the increased renewal rate, the data for expiring subscriptions and renewals are tabulated as in Table 1. The paradox of this example is that from Jan-

*Clifford H. Wagner is Assistant Professor, Department of Mathematical Sciences, The Capitol Campus, The Pennsylvania State University, Middletown, PA 17057.

uary to February the renewal rates actually declined in every category. Or, in the terminology of mutually favorable events (Blyth 1973, Chung 1942): overall February is favorable to renewal, while in each category January is favorable to renewal.

The primary cause of the misleading increase in the overall renewal rate is the sharp decrease in the relative importance of the subscription service category (and of its low renewal rate). Letting C_1 , C_2 , C_3 , C_4 , and C_5 designate the five subscriber categories, the overall renewal probability $P(R)$ is a weighted average of the renewal probabilities for the separate categories, in fact

$$P(R) = \sum P(R \cap C_i) = \sum P(C_i)P(R | C_i).$$

For January, the renewal probability is given by

$$.08(.81) + .40(.79) + .06(.60) + .45(.21) + .00(.09),$$

and the corresponding weighted average for February is

$$.10(.80) + .56(.76) + .24(.51) + .09(.14) + .00(.04).$$

Notice the decrease from .45 to .09 in the weight assigned to C_4 , the subscription service category. Changes in the weights as well as changes in the renewal rates of the separate categories determine the change in the overall rate.

3. INCOME TAX RATES

The second example of Simpson's paradox involves a comparison of federal personal income tax rates for different years (see Table 2). Between 1974 and 1978, the tax rate decreased in each income category, yet the overall tax rate increased from 14.1 percent to 15.2 percent. Again, the overall rates are weighted averages, with the tax rate for each category weighted by that category's proportion of total income. Because of inflation, in 1978 there were relatively more persons and consequently relatively more taxable dollars assigned to the higher income (i.e., higher tax rate) brackets. The reader may wish to speculate about the number of legislators who fully understand the effect of Simpson's paradox even though unaware of its official name.

Table 2. Total Income and Total Tax (in thousands of dollars), and Tax Rate for Taxable Income Tax Returns, by Income Category and Year

Adjusted Gross Income	1974			1978		
	Income	Tax	Tax Rate	Income	Tax	Tax Rate
under \$ 5,000	41,651,643	2,244,467	.054	19,879,622	689,318	.035
\$ 5,000 to \$ 9,999	146,400,740	13,646,348	.093	122,853,315	8,819,461	.072
\$ 10,000 to \$14,999	192,688,922	21,449,597	.111	171,858,024	17,155,758	.100
\$ 15,000 to \$99,999	470,010,790	75,038,230	.160	865,037,814	137,860,951	.159
\$ 100,000 or more	29,427,152	11,311,672	.384	62,806,159	24,051,698	.383
Total	880,179,247	123,690,314		1,242,434,934	188,577,186	
Overall Tax Rate			.141			.152

Table 1. Expiring Subscriptions, Renewals, and Renewal Rates, by Month and Subscription Category

Month	Source of Current Subscription					Overall
	Gift	Previous Renewal	Direct Mail	Subscription Service	Catalog Agent	
January						
Total	3,594	18,364	2,986	20,862	149	45,955
Renewals	2,918	14,488	1,783	4,343	13	23,545
Rate	.812	.789	.597	.208	.087	.512
February						
Total	884	5,140	2,224	864	45	9,157
Renewals	704	3,907	1,134	122	2	5,869
Rate	.796	.760	.510	.141	.044	.641

4. CONCLUSION

Simpson's paradox is not a contrived pedagogical example. Because this situation occurs at the level of a purely descriptive data analysis, it can easily bewilder the statistically naive observer. Classroom discussions of descriptive statistics should include examples of anomalies such as Simpson's paradox.

For additional discussion of Simpson's paradox and the more general topic of interactions and collapsibility in three-dimensional contingency tables, see Yule (1903), Simpson (1951), Bishop, Fienberg, and Holland (1975), Fienberg (1977), and Lindley and Novick (1981).

5. ACKNOWLEDGMENT

The author is grateful to James Rietmulder, Director of Planning at Historical Times Incorporated, publisher of *American History Illustrated*, for his cooperation and assistance in obtaining the subscription renewal data presented in this article.

[Received April 1981. Revised July 1981.]

REFERENCES

- BICKEL, P. J.; HAMMEL, E. A.; and O'CONNELL, J. W. (1975), "Sex Bias in Graduate Admissions: Data from Berkeley," *Science*, 187, 398-404.
 BISHOP, Y. M. M., FIENBERG, S. E., and HOLLAND, P. W.

(1975), *Discrete Multivariate Analysis: Theory and Practice*, Cambridge, Mass.: The MIT Press.

BLYTH, COLIN R. (1972), "On Simpson's Paradox and the Sure-Thing Principle," *Journal of the American Statistical Association*, 67, 364-366.

— (1973), "Simpson's Paradox and Mutually Favorable Events," *Journal of the American Statistical Association*, 68, 746.

CHUNG, KAI-LAI (1942), "On Mutually Favorable Events," *Annals of Mathematical Statistics*, 13, 338-349.

COHEN, MORRIS R., and NAGEL, ERNEST (1934), *An Introduction to Logic and Scientific Method*, New York: Harcourt, Brace and World, Inc.

FIENBERG, STEPHEN E. (1977), *The Analysis of Cross-Classified Categorical Data*, Cambridge, Mass.: The MIT Press.

LINDLEY, D. V., and NOVICK, MELVIN R. (1981), "The Role of Exchangeability in Inference," *The Annals of Statistics*, 9, 45-58.

SIMPSON, E. H. (1951), "The Interpretation of Interaction in Contingency Tables," *Journal of the Royal Statistical Society, Ser. B*, 13, 238-241.

The World Almanac and Book of Facts, (1977 and 1981 ed.), New York: Newspaper Enterprise Association, Inc.

YULE, GEORGE U. (1903), "Notes on the Theory of Association of Attributes in Statistics," *Biometrika*, 2, 121-134.

Another Approach to Incomplete Integrals

ANNE CHAO*

One usually writes the incomplete integrals as the sums of discrete probabilities by repeating the procedure of integration by parts. This work provides another approach by employing the binomial expansion.

KEY WORDS: Incomplete gamma; Incomplete beta; Binomial expansion.

In an introductory statistics course, one usually repeats the procedure of integration by parts to establish the relations between the incomplete integrals and discrete probabilities, such as the incomplete gamma and the Poisson

$$\int_0^c \frac{\alpha^r}{\Gamma(r)} x^{r-1} e^{-\alpha x} dx = 1 - \sum_{j=0}^{r-1} e^{-\alpha c} (\alpha c)^j / j! \quad (1)$$

where r is a positive integer, $\alpha > 0$, $c > 0$; or the incomplete beta and the binomial

$$\int_0^c \frac{\Gamma(n+1)}{\Gamma(k)\Gamma(n-k+1)} (1-x)^{n-k} x^{k-1} dx = \sum_{j=k}^n \binom{n}{j} c^j (1-c)^{n-j} \quad (2)$$

where k and n are positive integers, $0 < c < 1$. I have found that most students are more interested in the following direct approach, which employs only a binomial expansion. Note that

$$\int_c^\infty \frac{\alpha^r}{\Gamma(r)} x^{r-1} e^{-\alpha x} dx$$

*Anne Chao is Associate Professor, Institute of Applied Mathematics, National Tsing Hua University, Taiwan, Republic of China.

$$\begin{aligned} &= e^{-\alpha c} \frac{\alpha^r}{\Gamma(r)} \int_0^\infty (y+c)^{r-1} e^{-\alpha y} dy \quad (y = x - c) \\ &= e^{-\alpha c} \frac{\alpha^r}{\Gamma(r)} \sum_{j=0}^{r-1} \binom{r-1}{j} c^j \times \int_0^\infty y^{r-1-j} e^{-\alpha y} dy \\ &= \sum_{j=0}^{r-1} e^{-\alpha c} \frac{\alpha^r}{j!(r-1-j)!} c^j \times \frac{(r-1-j)!}{\alpha^{r-j}} \\ &= \sum_{j=0}^{r-1} e^{-\alpha c} (\alpha c)^j / j! \end{aligned}$$

which proves (1). Similarly, we can write

$$\begin{aligned} &\int_0^c \frac{\Gamma(n+1)}{\Gamma(k)\Gamma(n-k+1)} (1-x)^{n-k} x^{k-1} dx \\ &= \frac{\Gamma(n+1)}{\Gamma(k)\Gamma(n-k+1)} \int_0^1 [(1-c) + cz]^{n-k} c^k (1-z)^{k-1} dz \\ &\quad (x = c - cz) \\ &= \frac{\Gamma(n+1)}{\Gamma(k)\Gamma(n-k+1)} c^k \sum_{i=0}^{n-k} \binom{n-k}{i} c^i (1-c)^{n-k-i} \\ &\quad \times \int_0^1 z^i (1-z)^{k-1} dz \\ &= \sum_{i=0}^{n-k} \frac{n!}{(k-1)!i!(n-k-i)!} c^{k+i} (1-c)^{n-k-i} \\ &\quad \times \frac{(k-1)!i!}{(k+i)!} \\ &= \sum_{i=0}^{n-k} \binom{n}{k+i} c^{k+i} (1-c)^{n-k-i} \end{aligned}$$

which is exactly the right side of (2).