Henry Jay Becker Johns Hopkins University

In the social sciences, attention to the precise measurement of individual variables on individual cases has been supported in principle far more than in practice. However, in several of the disciplines efforts have been made both to educate other researchers of the biasing effects of measurement error and to improve the way in which necessarily "soft" data may be employed in analysis. There have been several approaches taken to the problem reflecting the different academic disciplines and corresponding traditions of research methodology involved.

Sociologists involved in the complexities of multivariate causal models have begun employing both measured but fallible variables and their corresponding unmeasured but "true" variables in their analytic efforts. The causal modeling approach to measurement error is typified by the following causal diagram which postulates measurement error for variable y only, measurement error which is random with respect to x and z, and no direct causal relationships between variables x and z.



These assumptions allow one to estimate the causal paths p_{XX} and p_{ZY} (equivalent in this case to the "true" correlation coefficients) from the observed correlations $r_{XY'}$, $r_{ZY'}$ and r_{XZ} . (Heise, 1969)

In fact, $p_{yx} = \begin{bmatrix} \frac{r_{xz} \cdot r_{xy}}{r_{y'z}} & \text{and} & p_{zy} \end{bmatrix} = \begin{bmatrix} \frac{r_{y'z} \cdot r_{xz}}{r_{xy'}} & \text{and} & p_{zy} \end{bmatrix}$

For example, suppose $r_{XY'}$ =.2, and $r_{Y'Z}$ =.4 and r_{XZ} = .1. If the above causal model is assumed, then p_{YX} = r_{XY} =.22 and p_{ZY} = r_{YZ} = .45. Also $r_{YY'}$ = .89.

Most work by Sociologists in this area has followed this approach or a companion "multiple indicators" approach (Hauser and Goldberger, 1971). The research for the most part has been limited to situations where measurement error for a given variable can be assumed to be randomly distributed with respect to other variables in the causal system and other error terms as well. (For an exception, see Sullivan, 1974.) This is primarily because these methodologists, dealing with many variables simultaneously, have had their hands full just incorporating random measurement error effect into their measurements. In the end, however, problems produced by nonrandom measurement error will necessarily have to be attended to.

The consequences of non-random measurement error have been addressed by the economists, Lansing and Morgan (1971). They have shown, for simple bivariate correlations between continuous variables, the consequences of a number of different <u>types</u> of non-random measurement error. Measurement error is discussed in the context of evaluating its effects on the regression coefficient $({}^{t}y'x')$ where the primes denote observed variables. If the observed values are composed of a "true" value plus an "error" term (e.g., y' = y + v), the estimating equation for β in terms of "true" scores and "error" terms becomes:

$$y'x' = \frac{\cot xy + \cot xv + \cot yu + \cot uv}{\operatorname{var} x + 2\cot xu + \operatorname{var} u} \quad (1)$$

and that for the correlation coefficient:

$$\mathbf{r}_{\mathbf{x}'\mathbf{y}'} = \frac{\operatorname{cov} \mathbf{x}\mathbf{y} + \operatorname{cov} \mathbf{x}\mathbf{v} + \operatorname{cov} \mathbf{y}\mathbf{u} + \operatorname{cov} \mathbf{u}\mathbf{v}}{\sqrt{\operatorname{var} \mathbf{x} + 2\operatorname{cov} \mathbf{x}\mathbf{u} + \operatorname{var} \mathbf{u}} \sqrt{\operatorname{var} \mathbf{y} + 2\operatorname{cov} \mathbf{y}\mathbf{v} + \operatorname{var} \mathbf{v}}}.$$

(a)

An examination of equations (1) and (2) shows the role played by three types of non-random measurement error as well as the role of random measurement error:

(a) Association between the error deviate terms \underline{u} and \underline{v} . (cov uv) Positive association between error terms spuriously raises the reported association between x and y. Negative assocation between error terms attenuates the reported x:y relationship from the true value. A large enough negative assocation between error terms could reverse the reported association between x and y.

(b) Association between the true value of one variable and the error term of the other. (cov xv + cov yu) This has the same kinds of effects as (a); e.g., positive associations spuriously raise correlation and regression coefficients while negative associations reduce them or reverse them.

(c) Association between the true value of one variable and its error deviate. (cov xu; cov yv) With the placement of these factors in the denominator of Eq. (1) and (2), a positive association here would lower the Q or r term from its true value while a negative association would increase the values of these statistics.

(d) <u>Random error</u> (var u; var v) Since these terms, also found in the denominator, can only be positive, their presence always causes regression and correlation coefficients to be underestimated

Categorical variables can also be examined in the above framework. Consider the association between two dichotomies each coded "O" and "L." In this case, the true value is always negatively correlated with its own error term, since a true "O" score will have an error term which is either "O" (accurately measured true score) or "+1" (measured in error) while a true "l" score will have a "-l" error term when falsely reported as "O." Alone, this factor would cause reported associations to be higher than the true association. At the same time, however, the magnitude of the errors--or more precisely, the total variance of the errors of observation in each variable considered separately--is acting to spuriously lower the reported association from the true value. This latter factor is <u>univer-</u> <u>sally</u> more powerful, as can be shown by extensive algebraic manipulation.

Thus, whenever there is no association between the error term of each variable and the true score of the other and where error terms themselves are uncorrelated, the net effect of the other two factors is to spuriously reduce measured association from its true value.

Conversely, wherever error terms are sufficiently positively correlated or where there is a large enough positive correlation between true values of one variable and error terms of the other, these factors can overcome the general tendency for attenuation of relationships. In particular, for the regression coefficient to be spuriously high, the offending factors (cov xv + cov yu + cov uv) must be \forall_{yx} times as great as the factors tending to diminish the measured relationship (cov xu + var u). The inequality with respect to the correlation coefficient is somewhat more complex: (cov xv + cov yu + cov uv)² ~ r_{xy}^2 [var y (cov xu + var u) + var x (cov yu + var v) + (cov xu + var u) (cov yv + var v)].

One of the principal advantages of working with dichotomies and categorical variables in general is the clarity with which statistical phenomena can be exemplified. To examine the effects of non-random measurement error on reported associations between variables, consider the data in Tables 1a and 1b.

Each table shows both the true and reported values of a variable X scored "+" and "-" crosstabulated by a variable Y, whose "types" represent population subgroups. Variable Y is assumed to be know without error. For each "type" (i.e. population subgroup), 90% of the x = "+" cases are known without error to be "+." (The other 10% are called "false negatives.") For two types (B and C), 95% of the x = "-" cases are known without error; only 5% are "false positives." In the other two types (A and D), false positives number 20% of the true x = "-" cases. Overall, more than 88% of the cases are accurately scored on both variables X and Y.

However, in <u>neither</u> case is the result similar to what it would have been if only random measurement error (i.e. equal error rates) were present. If, for example, the error rate were 10% in each cell, the reported percent difference would have been 16% instead of 20%, making the reported result a "conservative" estimate of the true relationship. Instead, in Table 1a, the reported association slightly <u>exaggerates</u> the true association between the variables. (In percentage difference terms, 20% points true difference is increased to 26%.) In Table 1b, the association is markedly <u>understated</u> by the reported data (the 20% true difference is reduced to 5%).

Several earlier papers by biostatisticians involved in public health research have discussed the

TABLE 1: EFFECTS OF DIFFERENTIAL FALSE POSITIVE RATES ON REPORTED ASSOCIATIONS (HYPOTHETICAL DATA)

Table la

	TRUE VALUES	5
	Type A	Type B
+	40%	20%
-	60%	80%
Total	(100)	(100)

%allierence = 20 percentage points
gamma = .45
odds ratio = 2.7

	ERROR	BATES	
True			1

Values	Туре А	Туре В	
+(FN)	10%	10%	
-(FP)	20%	5%	

	REPORTED VALUE	S
	Type A	Туре В
+	48%	22%
-	52%	78%
Total	(100)	(100)
% differenc gamma odds ratio	e = 26 percent =.53 = 3.3	age points

Table 1b

	TRUE VALUES						
	Type C	Type D					
+	40%	20%					
-	60%	80%					
Total	(100)	(100)					
gamma odds ratio	gamma = .45 odds ratio = 2.7						
Traio							
Values	Туре С	Type D					
+(FN)	10%	10%					
-(FP)	5%	20%					
	REPORTED VALU	ES					
	Type C	Type D					
+	39%	34%					
-	61%	66%					
Total	(100)	(100)					
% difference gamma odds ratio	e = 5 percenta = .11 = 1.2	ge points					

issue of non-random error's effects on measured association using cross-classified dichotomies, but their work needs more widespread attention and follow-up.

Keys and Kihlberg (1963) graphed the expected bias on individual variables (true vs. reported "prevalences"--i.e. percent "+") for a variety of false positive and false negative rates. They also graphed the bias of an associational measure (proportionate deviation from a true "relative risk" of 1.0) against that common true risk given particular combinations of the four error rates involved. However, the authors did not calculate bias in the measured assocation under the condition that different (true) rates were present (i.e., true relative risk \neq 1.0). They did explore the complexities of calculating estimated bias with fallibility on both measured variables, but produced no graphical examples as in the case of a single variable measured in error. (Instead of only four parameters -- two false positive and two false negative rates -- that situation involves 16 parameters.)

Goldberg (1975), examining the same type of hypothetical data, found that subgroup differences in false positive rates produces far larger average bias than differences between false negative rates for "prevalences" under 50%; coversely for higher prevalences.

One other point shown in the above example is that the <u>direction</u> of bias depends on whether the group with the greater or the lesser proportion of false positives. If false positives are more frequent in the group with more "+" cases, <u>generally</u> the reported relationships will be an exaggeration of the true relationship. Conversely major underreports of association tend to occur where false positives are concentrated in the group with fewer "+" cases. These results are reversed and apply to false <u>negative</u> rates where "+" proportions are above 50%.

In general, the reported percentage difference is related to the true percentage difference and the various error factors as $(p_1'-p_2') = p_1 (1-FN_1 - FP_1) - p_2 (1-FN_2 - FP_2) + (FP_1 - FP_2)$ where for subgroups i = 1, 2, p_1' are reported proportions "+," p_1 are true proportions "+", FN_1 are false negative rates (proportions of true "+" reported as "-") and FP_1 are false positive rates (proportion of true "-" reported as "+"). Each p_1 is also related to the reported p_1' and the error terms, thusly: $p_1 = \frac{p_1' - FP_1}{1-FN_1-FP_1}$.

Under the assumption that all error rates, FN_i and FP_i \leq e, where $p_1 - p_2 \geq 0$, the maximum and minimum value for $(p_1 - p_2)$ are given by

$$\left(\frac{1}{1-e}\right)\left(p_{1}' - p_{2}'\right) \pm \left(\frac{e}{1-e}\right)$$
.

As shown in Table 2, extremely large variations in the magnitude and direction of bias can be observed due to non-random measurement error even when the absolute proportions of cases in error is rather modest, for example e = .2.

П

In most cases, although not for the maxima and minima shown in Table 2, the magnitudes of p_1' and p_2' as well as their difference $(p_1' - p_2')$ affect the value expected for $(p_1 - p_2)$ under specific error rates. (See Table 3.) The smaller the magnitude of p_1' and p_2' (given $p_1' - p_2' = k$), the greater the difference between true and reported associations.

REPORTED MAX & MIN TRUE VALUES (p1 - p2) VALUES2 FOR MAXIMUM ERROR RATE el								
(p1'-p2')	е	=•3	е	=.2	е	=.1		
	MAX	MIN	MAX	MIN	MAX	MIN		
+•3	.86	.00	.62	.12	.44	.22		
+.2	.71	14	.50	.00	•33	.11		
+.1	•57	29	•38	12	.22	.00		
+.05	.50	 36	.31	19	.17	06		
.00	•43	 43	+.25	25	•11	11		
 excluding effects of sampling error with p₂' > p₁', results are the same except minime and maxime are revensed and signs 								

ABLE 2:	RANGES	OF P	OSSIBLE	TRUE	PERCENTA	GE
	DIFFERE	NCES	GIVEN	CERTAI	N LIMITS	ON
	MEASURE	MENT	ERROR			

Even rather minor differences in the two false positive rates can produce rather major differences in measured association. The example in Table 4 shows that under rather ordinary conditions, a reversal of false positive rates changes the result from nearly no measurement error effect (situation (1)) to a very strong attenuation of

are reversed also.

the true relationship (situation (2)).

The important questions that are raised by this discussion of non-random error concern not the <u>possible</u> effects of differential error rates, but the actual effects caused by real differences in error rates. Here our knowledge is hampered by the paucity of validation studies that report validity for subgroups rather than solely for the complete sample studied. Two sources have been found in the public opinion literature which do report on validity of responses for subgroups.

One is the Denver validity study done in 1949 but most recently analyzed in (Cahalan, 1969). Reworking the results on the basis of the published statistics, interesting and quite major differences between subgroups can be seen for both false negative and false positive rates, with the consequence that many relationships would have been reported erroneously without the validating information.

The Cahalan data is only approximately reproduced in Table 5. (The absence in the original article of sufficient information about the magnitude of "don't knows" for men and women separately made it impossible to exactly reproduce the data.) The data as estimated here, though, exhibit some remarkable differences between cells in error rates for some of the measured associations. Notice that in both cases of large reported associations by sex (the Community Chest contributions and the drivers' licenses), the <u>true</u> relationships are of smaller magnitude--one of them being less than 40% as large as the reported association.

The second article from <u>Public Opinion Quarterly</u> reporting validity coefficients for subgroups is Weiss's (1968) report of interviews with mothers receiving welfare. The results given in Table 6 illustrate the tendency for upward bias in reported percentage differences to occur only where both $p_1 - p_2$ is small and differences in false positive rates exist. For cases where $p_1 = p_2$, bias is also greater when the p_1 approach zero. (See Keys and Kihlberg, 1963.)

Finally, the question remains: given the clearly major consequences of non-random measurement error on reported associations between variables, what now needs to be done to better take account of this problem?

Journal articles most often report findings with due respect for sampling error. Confidence limits, for example, are reported for percentage differences and the result is reported as statistically significant if the null hypothesis of "no difference"lies outside the range of 2 standard errors of the sample statistic. It would seem to make some sense to include in this statement the range of uncertainty that is attributable to possible measurement errors.

That is, the analyst might propose two sets of error rates for a given two-by-two table that would have opposite effects on the bias of the reported result and which are just extreme enough to be plausible given the possible error-producing causes present in the particular data. Those two sets of error rates, in turn, would be applied to the reported data in order to produce "extreme" but plausible "true" sets of sample data. Confidence intervals would then be calculated for these extreme but hypothetically true sample data tables. The summary confidence intervals reported would thus reflect both sampling error and plausible ranges of error caused by measurement inadequacies.

Consider an example, more or less randomly chosen from Rosenberg's The Logic of Survey Analysis (1968). (See Table 7.) In our use of his re-examination of the relationship between vote intention and respondent's education (1948 data), we percentage on education (% with some high school) taking <u>education</u> as the variable possibly reported in error. Vote <u>intention</u> is assumed to be reliably and validly known. Let us make two contrasting suppositions: (t₁) that persons claiming an intention to vote are more

ABLE 3:	EFFECTS OF MAG	NITUDE OF pi'	AND $(p_1' -$	p2') ON BIAS OF
	(p1'-p2') GIVEN	FOLLOWING ER	ROR RATES:	$FN_1 = .10$
	$FN_2 = .10 FP_1$	= .20 FP2 :	10	

REPORTED VALUES		BIAS		
p1'-p2'	p1'	₽2'	p1-p2	(p1-p2)- (p1'-p2')
•30	{•5 •4	.2 .1	•30 •29	.00 01
.20	$\{ .5 \\ .3 \}$.3 .1	.18 .14	02 06
 .10	(.5 .3	.4 .2	.05 .02	05 08
.05	(•5 •3	.45 .25	01 05	06 10
.00	.5 .3	.5 .3	07 11	07 11
05	{ .45 .25	•5 •3	14 18	09 13
10	.4	•5 •3	21 25	11 15
20	•3	•5	 36	16
30	.2	•5	50	20

PARLE 4.	EFFECTS	OF	SMAT.T.	DIFFFERENCES	TN	FALSE.	POSTTTVE	RATES

		T	RUE SCOR	ES		
			A	В		
		%+	25%	15%		
	% dii gamma odds	fference = a = ratio =	= 10 per = .31 = 1.9	centage	points	
ERR	OR RATE	ES (1)		ERRO	OR RATE	5 (2)
	A	В			A	B
FN	12%	12%		FN	12%	12%
FP	15%	10%		FP	10%	15%
REPOR	TED RES	SULTS (1)		REPOR	TED RESI	JLTS (2) B
		B				
%+	33%	22%		%+	30%	26%
% dif: pero gamma odds :	ference centage ratio	e = 11 e points = .27 = 1.7		% dif: pero gamma odds n	ference centage catio	= 4 points = .10 = 1.2

		(77)			ERROR	RATES ²	~
		(1)	KEPORTED-	TRUE	FN	FP	Gammas
Community	Chest Contributions	(percent '	'yes")				
	MEN WOMEN	327 404	85% <u>55%</u> 30%	35% 24% 11%	0%3 0%	76% 41%	r=.65 t=.26
Driver's	License (percent "ye	s")					
	MEN WOMEN	423 497	78% <u>33%</u> 45%	66% <u>32%</u> 34%	3%4 12%4	41% 8%	r=.76 t=.61
Voting in	1946 Cong. Election	s (percent	"yes")	· · · · · · · · · · · · · · · · · · ·			
	MEN WOMEN	382 447	55% 48% 7%	32% <u>33%</u> -1%	5% 9%	37% 27%	r=.14 t=02
Voting in	1948 Pres. Election	s (percent	"yes")				
	MEN WOMEN	423 497	74% <u>72%</u> 2%	61% <u>61%</u> 0%	1.6% 1.3%	36% 31%	$r^{=.05}$ t=.00

4. Unable to reproduce marginals from published data: doubtful accuracy.

PERCENT VOTING IN 1964 ELECTION								
	(N)		 ידו זכויוי	ERROR RATES3		Gammas		
	(11)	THE CRIED THE		FN	FP	FP		
Worked more than 10 years ¹ Worked less than 10 years	189 329	53% 44% 9%	29% <u>28%</u> 1%	0%2 0%	34% 22%	r=.17 t=.02		
Expects children to continue educ. past high school Does not expect	141 191	64% 47% 17%	42% 29% 13%	0% ² 0%	39% 26%	r=.33 t=.29		
Did not get as much education as wanted Got enough education \bigtriangleup	372 147	48% <u>47%</u> 1%	28% <u>28%</u> 0%	0%2 0%	27% 26%	r=.02 t=.00		
 Assumed known without error Assumed to be zero FN = false negative FP = false positive 								

TABLE 6: MEASUREMENT ERROR EFFECTS IN THE WEISS (1968) STUDY OF WELFARE MOTHERS

likely to overreport their education than those not claiming an intention to vote (because of the good citizen image of voting and completing one's education); (t_2) alternatively, that persons not planning to vote, being more alienated from the culture, make more errors in their selfreports (both over-reporting and underreporting) than do those planning to vote. Table 7 shows the the four parameters assumed for each of these two basic models. The data in Table 7 show how the reported relationship (a) is affected by the measurement error assumptions (t_1) and (t_2) producing hypothesized "true" sample results (b) and (c). If the confidence intervals for (p_1-p_2) are calculated from the "extreme" results (b) and (c) instead of from (a), they change the estimate of (p_1-p_2) from $8\% \le (p_1-p_2) \le 20\%$, p = .95, to $5\% \le (p_1-p_2) \le 24\%$, p = .95.

 TABLE 7:
 ILLUSTRATION OF THE USE OF MEASUREMENT

 ERROR HYPOTHESES IN THE ROUTINE

 PRESENTATION OF CONFIDENCE INTERVALS¹

	REPO EDUCA % WIT	RTED TION: H SOME SCHOOL	EXPECTED UNDER ERROR HYPOTHESES						
	(% "	+")	tı		^t 2 ′				
Vote Intention:	(a)	(N)	(b)	(c)				
Positive	59%	(2515)	52%	5	8%				
Negative	43%	(297)	<u>41%</u>	<u>4</u>	0%				
\bigtriangleup	16%	16%		11% 18					
Confidence Intervals:									
p ≕•95	10%≤(p ₁ -p ₂)≤22%	5%≲(p 1- p2)∴24%						
p=.68	13%≤(p ₁ -p ₂)≤19%	8%≤(p ₁ -p ₂)≦21%						
Educa tl FN	tion FP	t ₂		Educa FN	t io n FP				
Vote Int: Pos .01	.15	Vote Int:	Pos	.02	.05				
Neg .03	.05		Neg	.07	.10				
	·	L		• • • • • • • • • • • • • • • • • • •	·				

1. Data from Rosenberg, 1968

One standard error confidence limits for the data based on the measurement error assumptions do not differ much from the two-standard error limits based solely on sampling error: $8\% \le (p_1-p_2) \le 21\%$.

Another procedure that could be more widely utilized is a version of Schuman's suggested "random probe" method of clarifying the meaning of information obtained through interview procedures. Schuman (1966) proposed that respondents to a given survey be asked to give reasons, interpretations, or explanations for their answer to a subset of all closed-ended questions in the interview. These extended responses are used to calculate the proportion of closed-ended responses that were "accurate." A different subset of respondents is used to 'validate" in this way each of several subsets of questions.

This random probe technique has the potential of providing information concerning response error rates for various subgroups being compared in the statistical analysis. Thus real data can be employed as "expected values" in the error rate table, and variances for these error rates can be computed. These variances can be used to produce the "extreme" values of the "true" results which can then be used, along with the sampling distribution of the statistic, to compute reasonable confidence limits for the data.

Obviously, much additional work needs to be done in order to make it possible for measurement error effects to be incorporated into routine statistical procedures to the extent that sampling error has been. But the expanding application of complex statistical methodologies to the "soft" data of the social sciences has not been accompanied by sufficient appreciation of the magnitude and direction of effects caused by unreliability and invalidity of measurement. This is a problem that requires major attention if we are to avoid perpetuating the publication of technically flawless conclusions invalidated by the use of extremely soft and fallible data.

References

Cahalan, Don, "Correlates of Respondent Accuracy in the Denver Validity Survey," Public Opinion Quarterly, Vol. 32, (Winter, 1968), 607-621.

Goldberg, Judith D., "The Effects of Misclassification on the Bias in the Difference Between Two Proportions and the Relative Odds in the Fourfold Table," <u>Journal of American Statistical</u> Association, Vol. 70 (September, 1975), 561-567.

Hauser, Robert M. and Goldberger, Arthur S., "The Treatment of Unobserved Variables in Path Analysis," in <u>Sociological Methodology 1971</u>, Herbert L. Costner, ed. Jossey-Bass, 1971.

Heise, David R., "Problems in Path Analysis and Causal Inferences," in <u>Sociological Methodology</u> <u>1969</u>, Edgar F. Borgatta, ed. Jossey-Bass, 1969.

Keys, Ancel and Kihlberg, Jaakko K., "Effect of Misclassification on Estimated Relative Prevalence of a Characteristic," <u>American Journal of</u> <u>Public Health</u>, Vol. 53 (October 1963), 1656-1665.

Lansing, John B. and Morgan, James N., <u>Economic</u> <u>Survey Methods</u>. Institute for Social Research, University of Michigan, 1971.

Rosenberg, Morris, <u>The Logic of Survey Analysis</u>. Basic Books, 1968.

Schuman, Howard, "The Random Probe: A Technique for Evaluating the Validity of Closed Questions," <u>American Sociological Review</u>, Vol. 31 (1966), 218-222.

Sullivan, John L., "Multiple Indicators: Some Criteria of Selection," in <u>Measurement in the</u> <u>Social Sciences</u>, H.M. Blalock, Jr., ed. Aldine, 1974.

Weiss, Carol H., "Validity of Welfare Mothers' Interview Responses," <u>Public Opinion Quarterly</u>, Vol. 32 (Winter, 1968), 622-633.