

*Measurement error in independent variables produces biased estimates of the coefficients in linear models. These biases can be reduced by obtaining repeated measurements of the variables and then estimating structural equation models with multiple indicators of latent variables. Remeasurement is usually costly, however, raising the question of whether the same benefits can be obtained by remeasuring only a fraction of the sample. Although this strategy has been tried previously, there were no appropriate statistical methods for combining the data in the remeasured subsample and the single measurement subsample. We demonstrate here how recently developed methods for incomplete data provide an attractive solution to this estimation problem. The methodology is illustrated by a reanalysis of Bielby, Hauser, and Featherman's (1977a) study of the OCG-II data.*

## Reducing Bias in Estimates of Linear Models by Remeasurement of a Random Subsample

PAUL D. ALLISON  
*University of Pennsylvania*

ROBERT M. HAUSER  
*University of Wisconsin*

**I**t is well known that random error in the measurement of independent variables in a linear regression model leads to bias in least-squares estimates of the coefficients (e.g., Intriligator 1978, p. 190). The bias may be substantial and may be either positive or negative, depending in complex ways on the true coefficients, the degree of measurement error in each variable, and the pattern of intercorrelations among the independent variables.

In the last fifteen years, there have been major advances in the development of methods for correcting such biases (Jöreskog and Sörbom 1979; Bentler and Weeks 1980). These methods typically require two or more indicators for each variable that is measured with

---

AUTHORS' NOTE: *We thank Peter Bentler for helpful suggestions.*

SOCIOLOGICAL METHODS & RESEARCH, Vol. 19, No. 4, May 1991 466-492  
© 1991 Sage Publications Inc.

error. The indicators are postulated to be linear functions of a smaller set of latent variables that are error-free, plus random error components. Other linear equations may describe structural relationships among the latent variables and other variables measured without error. The entire system of equations is usually estimated by maximum likelihood under the assumption of multivariate normality, although estimation under less restrictive assumptions is now possible (Shapiro and Browne 1987; Bentler and Berkane 1985; Bentler 1989, ch. 10).

While this approach can produce major improvements in the estimation of linear models, the need for multiple indicators can substantially increase the cost of data collection. Ideally, the multiple indicators for a single latent variable should be close replications of one another. To avoid mutual contamination, the replicate measurements are best made at distinct points in time. But this usually means costly and time-consuming reinterviews. Under a fixed budget, the increased cost per case may require a large reduction in sample size. Hence the reduction in bias may be purchased at the cost of greatly increased sampling error.

In this article we describe a method that may improve the terms of this trade-off. The basic idea is to take single measurements for a large sample and repeated measurements for a smaller, random subsample. The small subsample provides the necessary information to eliminate the bias, while the large sample provides for adequate precision.

The idea is not novel. Bielby, Hauser, and Featherman (1977a) used this approach to estimate a linear model for the attainment of occupational status by American males. For a sample of 25,223 nonblack males, Bielby, Hauser, and Featherman (hereafter referred to as BHF) had data on age, education, current occupation, first occupation, father's occupation, father's education, and parents' income. Repeated measurements on these variables were obtained in reinterviews for a random subsample of 578. Using the remeasurement data, they were able to correct for bias due to measurement error in the coefficients of their linear model.

As BHF recognized, however, their estimation method left much to be desired. In a two-step procedure, they first obtained maximum likelihood estimates of the measurement error variances, using the subsample of 578. This step was unproblematic. In the second step,

these estimates were used to correct the variances for the full sample of 25,223. Ordinary least-squares (OLS) estimates of the coefficients were then obtained from the corrected covariance matrix. These estimates were presumably consistent, but their efficiency is unknown and there is no known way to obtain consistent estimates of the standard errors. Thus, although the idea of combining a small remeasurement subsample with a larger sample is attractive, it must be judged unsatisfactory until efficient estimation procedures are available.

We show here how these statistical problems can be solved by simultaneously estimating a model for both the large and small subsamples. This model incorporates both the equations for the dependence of the observed variables on error-free latent variables and the structural equations for the relationships among the error-free variables. Estimation is by maximum likelihood under the assumption of multivariate normality, although less restrictive estimation procedures may also be used. The estimates may be obtained by applying the methods proposed by Allison (1987) and Muthén, Kaplan, and Hollis (1987) for estimation of linear models with incomplete data. We shall use this procedure to obtain efficient estimates for the BHF data using the EQS 3.0 program (Bentler 1989). Allison (1987) used LISREL 6, which required a rather convoluted procedure to incorporate data and models with structured means. This is radically simplified in EQS 3.0, which can directly model mean structures. Such simplified analyses can also be done with LISREL 7 (Jöreskog and Sörbom 1988), using the most recent version (7.17).<sup>1</sup>

#### *THE BIELBY, HAUSER, AND FEATHERMAN STUDY*

A detailed description of the data used by BHF can be found in their (1977a, 1977b) articles. Here we give only the highlights. As part of the U.S. Current Population Survey (CPS) in March 1973, household interviews produced data on age, years of schooling, and current occupation for more than 27,000 male members of the experienced civilian labor force. These data were supplemented with a mail-out, mail-back questionnaire distributed in the fall of that year. Additional

TABLE 1: Variables in the BHF Study

Variable	Occasion		
	March 1973 Household Interview	Fall 1973 Questionnaire	Fall 1973 Remeasurement Interview
1. Father's occupational status (FO)	—	$x_{11}$	$x_{12}$
2. Father's educational attainment (FE)	—	$x_{21}$	$x_{22}$
3. Parent's income (PI)	—	$x_{31}$	$x_{32}$
4. Educational attainment (ED)	$x_{43}$	$x_{41}$	$x_{42}$
5. Occupational status of first job (O1)	—	$x_{51}$	$x_{52}$
6. Current occupational status (OC)	$x_{63}$	—	$x_{62}$
7. Age (AGE)	$x_{73}$	—	—
8. Age squared (AGE2)	$x_{83}$	—	—

variables in this second questionnaire included first occupation, father's education, father's occupation, and parents' income when the respondent was 16 years old. Finally, for a random subsample of about 1,000 respondents (600 nonblack and 400 black), a third interview was conducted by telephone (or in person) about three weeks after the return of the mail questionnaire. This final interview produced repeated measurements of the variables already described.

Table 1 contains a list of the variables, the occasions on which they were obtained, and the symbols and acronyms used to refer to them in this article. Note that each of the occupation reports was scaled using the Duncan (1961) SEI scores to measure socioeconomic status. Educational attainment is coded in exact years of schooling completed, and parents' income is coded as the logarithm of price-adjusted dollars. Age is expressed in years divided by ten, and a quadratic transformation, AGE2, is defined as  $(\text{years} - 40)^2/10$ .

The analysis in BHF (1977a) was restricted to nonblacks, and we shall maintain that restriction here. (For an examination of blacks, see BHF [1977b].) Table 2 gives correlations, standard deviations, and means for 24,645 males who were not remeasured and for 578 males who were remeasured. The matrix for the remeasurement subsample is exactly as it appears in BHF (1977a). BHF also reported correlations, standard deviations, and means for the full sample of 25,223. We have corrected these statistics to approximately what they would





be if the 578 remeasured cases were excluded. As might be expected, these corrections were very small, affecting only the third decimal place.<sup>2</sup>

Assuming no measurement error, the model of interest was a fully recursive system of linear equations in which education, first occupation, and current occupation were endogenous (dependent) variables, and father's education, father's occupation, parental income, and age were exogenous (predetermined) variables:

$$ED = \beta_1 AGE + \beta_2 AGE2 + \beta_3 FO + \beta_4 FE + \beta_5 PI + \epsilon_1 \quad (1)$$

$$O1 = \beta_6 AGE + \beta_7 AGE2 + \beta_8 FO + \beta_9 FE + \beta_{10} PI + \beta_{11} ED + \epsilon_2 \quad (2)$$

$$OC = \beta_{12} AGE + \beta_{13} AGE2 + \beta_{14} FO + \beta_{15} FE + \beta_{16} PI + \beta_{17} ED + \beta_{18} O1 + \epsilon_3 \quad (3)$$

The disturbance terms,  $\epsilon_1$ ,  $\epsilon_2$ , and  $\epsilon_3$  were assumed to be independent of each other and of all the independent variables. All variables are assumed to be measured as deviations from means, which eliminates intercept terms from the equations.

Using the covariance matrix (not shown) for the entire sample of 25,223, OLS estimates were computed for the coefficients in equations (1)-(3). These are presented in Column 1 of Table 3, together with their estimated standard errors. These estimates, which are presumably biased as a result of measurement error, will serve as a basis for comparison with estimates that correct for measurement error. The estimates reported in Table 3 differ slightly from those reported by BHF because they constrained  $\beta_9$ ,  $\beta_{10}$ ,  $\beta_{15}$ , and  $\beta_{16}$  to be zero. This made sense for  $\beta_9$ ,  $\beta_{15}$ , and  $\beta_{16}$  because those coefficients were not significantly different from zero at the .05 level. The justification for setting  $\beta_{10}$  equal to zero was that (a) a negative effect of parental income on occupational status is implausible, and (b) the coefficient was not significantly different from zero in a model for the subsample of 578 (to be discussed in the next section). As we see here, however, the estimate for this coefficient is more than eight times its standard error. Moreover, in later sections we will show that the negative effect becomes even stronger when measurement error bias is corrected.

TABLE 3: Estimates (and Standard Errors) for Structural Coefficients

Variable	Coefficient	(1) Uncorrected OLS	(2) Corrected OLS <sup>a</sup>	(3) ML on Small Subsample	(4) ML using All of the Data
<i>Eq. (1):</i>					
AGE	$\beta_1$	-.058 (.013)	-.034	b	-.029 (.013)
AGE2	$\beta_2$	-.019 (.0011)	-.018	b	-.018 (.0011)
FO	$\beta_3$	.021 (.0008)	.025	.022 (.0062)	.025 (.0015)
FE	$\beta_4$	.183 (.0050)	.175	.162 (.034)	.174 (.0078)
PI	$\beta_5$	2.18 (.044)	2.40	2.43 (.310)	2.47 (.076)
<i>Eq. (2):</i>					
AGE	$\beta_6$	1.75 (.099)	1.73	b	1.60 (.094)
AGE2	$\beta_7$	-.125 (.0081)	-.110	b	-.101 (.0077)
FO	$\beta_8$	.194 (.0061)	.243	.251 (.046)	.240 (.012)
FE	$\beta_9$	.026 (.038) <sup>c</sup>	-.301	-.960 (.257)	-.250 (.051)
PI	$\beta_{10}$	-2.84 (.336)	-.590	.020 (2.40) <sup>c</sup>	-5.84 (.48)
ED	$\beta_{11}$	4.76 (.046)	5.58	5.50 (.326)	5.37 (.099)
<i>Eq. (3):</i>					
AGE	$\beta_{12}$	2.86 (.099)	2.65	b	2.46 (.10)
AGE2	$\beta_{13}$	-.153 (.0081)	-.132	b	-.122 (.0079)
FO	$\beta_{14}$	.074 (.0063)	.063	.103 (.049)	.063 (.011)
FE	$\beta_{15}$	.065 (.038) <sup>c</sup>	-.067	-.186 (.264) <sup>c</sup>	-.033 (.046) <sup>c</sup>
PI	$\beta_{16}$	-.524 (.336) <sup>c</sup>	-1.22	-2.63 (2.41) <sup>c</sup>	-1.31 (.45)
ED	$\beta_{17}$	2.58 (.055)	2.42	2.29 (.446)	2.35 (.15)
O1	$\beta_{18}$	.392 (.0063)	.502	.508 (.053)	.50 (.024)

<sup>a</sup>Standard errors not available.<sup>b</sup>Coefficient set to zero because of unavailable data.<sup>c</sup>Coefficient less than twice its standard error.



For BHF, the first step in getting corrected estimates was to estimate a measurement model for the subsample of 578 with repeated measurements. They tried several different models, all special cases of the confirmatory factor model (Jöreskog 1969) that postulates a set of "true scores" or latent variables underlying the observed variables. Each set of replicate measurements is assumed to depend linearly on a single latent variable plus random error components. The class of models BHF considered can be expressed as

$$\begin{aligned}
 x_{11} &= \lambda_{11} \text{FO} + e_{11} \\
 x_{12} &= \lambda_{12} \text{FO} + e_{12} \\
 x_{21} &= \lambda_{21} \text{FE} + e_{21} \\
 x_{22} &= \lambda_{22} \text{FE} + e_{22} \\
 x_{31} &= \lambda_{31} \text{PI} + e_{31} \\
 x_{32} &= \lambda_{32} \text{PI} + e_{32} \\
 x_{41} &= \lambda_{41} \text{ED} + e_{41} \\
 x_{42} &= \lambda_{42} \text{ED} + e_{42} \\
 x_{43} &= \lambda_{43} \text{ED} + e_{43} \\
 x_{51} &= \lambda_{51} \text{O1} + e_{51} \\
 x_{52} &= \lambda_{52} \text{O1} + e_{52} \\
 x_{62} &= \lambda_{62} \text{OC} + e_{62} \\
 x_{63} &= \lambda_{63} \text{OC} + e_{63}
 \end{aligned} \tag{4}$$

where FO, FE, PI, ED, O1, and OC now refer to error-free latent variables. AGE and its quadratic transform AGE2 are assumed to be measured without error. Correlations among the latent variables are not restricted in any way. The random error terms are assumed to be independent of one another and of the true-score latent variables.

Metrics for the latent variables are established by constraining  $\lambda_{11} = \lambda_{21} = \lambda_{31} = \lambda_{43} = \lambda_{51} = \lambda_{63} = 1.0$ . Normalization of this kind is necessary because the metric of a latent variable is arbitrary; consequently the slope coefficients are identified only relative to each other. Although not necessary for identification, we also constrain all the other  $\lambda$ s equal to 1.0. This is a plausible restriction because the observed indicators of each latent variable were designed to have the same metric. In fact, unconstrained estimates of these coefficients show only small departures from unity.

The measurement model was estimated using the computer program LISREL 6, which does maximum likelihood estimation of

TABLE 4: Error Variance and Reliability Estimates

Variable		Remeasurement Subsample		Full Sample	
		Error Variance	Reliability	Error Variance	Reliability
FO	x <sub>11</sub>	87.76	.85	70.99	.84
	x <sub>12</sub>	63.43	.89	75.87	.83
FE	x <sub>21</sub>	1.25	.92	1.14	.93
	x <sub>22</sub>	.87	.95	.97	.94
PI	x <sub>31</sub>	.021	.88	.020	.88
	x <sub>32</sub>	.007	.95	.008	.94
ED	x <sub>41</sub>	3.15	.73	3.15	.71
	x <sub>42</sub>	.45	.95	.49	.94
	x <sub>43</sub>	.87	.89	.82	.90
OI	x <sub>51</sub>	97.35	.84	80.11	.84
	x <sub>52</sub>	85.57	.85	96.36	.82
OC	x <sub>62</sub>	136.17	.78	148.39	.74
	x <sub>63</sub>	117.92	.81	98.94	.81

confirmatory factor models under the assumption of multivariate normality. Estimates of the error variances and derived reliabilities of the observed variables are given in Table 4. The estimation procedure also produces a likelihood ratio chi-square test of the model against the alternative model that imposes no restrictions on the variance/covariance matrix. This test yielded a chi-square value of 60.1 with 57 degrees of freedom, indicating that the model fits the data reasonably well.

BHF used the error variance estimates to "correct" the variance/covariance matrix for the full sample of 25,223. The correction consisted simply of subtracting the estimated error variance from the observed variance of the corresponding variable. The covariances—which should not be affected by uncorrelated measurement errors—were left unchanged. This corrected matrix was then used as though it were an ordinary covariance matrix to generate pseudo-OLS estimates. The estimates we report in column 2 of Table 3 differ slightly from those given by BHF because BHF forced  $\beta_9$ ,  $\beta_{10}$ ,  $\beta_{15}$ , and  $\beta_{16}$  to be zero for the reasons given above.

Comparing the corrected with the uncorrected estimates, BHF concluded that biases resulting from measurement error are not as

large as some others have previously suggested (Bowles 1972; Bowles and Nelson 1974). The largest discrepancy they found was for  $\beta_{18}$ , whose uncorrected estimate was 22% lower than the corrected estimate. Because they forced  $\beta_{10}$  to equal zero, however, they failed to observe that the uncorrected estimate for  $\beta_{10}$  is less than half as large as the corrected estimate, a result that holds up with the improved estimation methods discussed in the next two sections.

Because the error variance estimates produced by BHF were consistent, it seems quite likely that the coefficient estimates produced by their two-stage procedure are also consistent. However, the efficiency of these estimates is not known. There is also no theory to indicate how one might estimate standard errors, hence none are reported. In the remaining sections we show how to remedy these statistical deficiencies.

### *COMBINING THE STRUCTURAL AND MEASUREMENT MODELS*

In this section we discuss how equations (1)-(4) may be combined into a single model to be estimated by EQS 3.0. The structural equations for the latent variables are specified as

$$\eta = \beta\eta + \Gamma\xi + \zeta \quad (5)$$

where  $\eta$  is a vector of endogenous variables,  $\xi$  is a vector of exogenous variables, and  $\zeta$  is a vector of unobserved disturbances.  $\xi$  and  $\zeta$  are assumed to be uncorrelated.  $\beta$  and  $\Gamma$  are matrices of coefficients. For the BHF data,  $\eta = (\text{ED}, \text{OC1}, \text{OC})'$  and  $\xi = (\text{FO}, \text{FE}, \text{PI}, \text{AGE}, \text{AGE2})'$ . Because this is a recursive model,  $\beta$  is assumed to be lower triangular and  $\text{var}(\zeta)$  is assumed to be diagonal.

We also have equations relating the latent variables to the observed indicators:

$$\begin{aligned} y &= \Lambda_y \eta + \varepsilon \\ x &= \Lambda_x \xi + \delta \end{aligned} \quad (6)$$

Here  $x$  and  $y$  are vectors of observed variables and  $\varepsilon$  and  $\delta$  are vectors representing random measurement error. We let  $y = (x_{41}, x_{42}, x_{43}, x_{51}, x_{52}, x_{62}, x_{63})'$  and  $x = (x_{11}, x_{12}, x_{21}, x_{22}, x_{31}, x_{32}, x_{73}, x_{83})'$ .  $\Lambda_y$  and  $\Lambda_x$  are

matrices of coefficients. Both  $\epsilon$  and  $\delta$  are assumed to be uncorrelated with each other and with  $\eta$ ,  $\xi$ , and  $\zeta$ . We also assume that both  $\text{var}(\epsilon)$  and  $\text{var}(\delta)$  are diagonal matrices. The lambda matrices are given by

$$\Lambda_y = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \quad \Lambda_x = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

To reflect the fact that age and its quadratic transform are assumed to be measured without error, we fix  $\text{var}(\delta_7) = \text{var}(\delta_8) = 0$ .

When the data consist of a complete covariance matrix, estimation of such models is a routine task for LISREL or EQS. The standard method is maximum likelihood, under the assumption that the data come from a multivariate normal distribution. Unfortunately, neither of the correlation matrices given in Table 2 is complete, so standard methods are not appropriate. The matrix for the large subsample is highly incomplete, lacking any of the correlations pertaining to the remeasurement interview. This, of course, will be typical of the research designs we are advocating. More unusual is the fact that the matrix for the small subsample lacks correlations and standard deviations for AGE and AGE2. While these moments could have been computed from the original data, they were not reported by BHF. Nevertheless, the method we propose can accommodate such missing moments.

To illustrate methods for complete data and to provide an additional standard of comparison, we eliminated AGE and AGE2 from the model and estimated it using only the moments for the small subsample. Although corrected for bias due to measurement error, these estimates may suffer some additional bias due to the exclusion of AGE and AGE2 from the three structural equations. Maximum likelihood estimates of the coefficients in equations (1)-(3) are given in Column 3 of Table 3. With a few exceptions, the estimates are very similar to those obtained by BHF with their two-stage procedure. The exceptions

generally occur for those coefficients that had large standard errors in both Columns 1 and 2, and are probably a result of sampling error. A noteworthy exception is  $\beta_{10}$ , which is very much smaller in Column 3 than in either Column 1 or Column 2.

The chi-square goodness-of-fit statistic for this model was 60.1 with 57 degrees of freedom, exactly what was reported earlier for the measurement model alone. This is to be expected because equations (1)-(3) impose no further restrictions on the covariance matrix.

### *ESTIMATION WITH INCOMPLETE DATA*

As seen in the previous section, the type of research design we are considering produces incomplete data as a matter of course. There will usually be complete data for the remeasurement subsample (although not for the example considered here), but the single-measurement subsample will always lack data for all the remeasurements. In this section we show how to combine the complete and incomplete data to produce efficient estimates.

The method we will use was first proposed for confirmatory factor models by Werts, Rock, and Grandy (1979) and later generalized by Allison (1987) and Muthén, Kaplan, and Hollis (1987) to all models of the sort described in the previous section. It is a general method for maximum likelihood estimation of linear models with incomplete data, under the assumptions that the data are missing at random (Rubin 1976) and come from a multivariate normal distribution. Unlike more conventional methods for incomplete data, this method produces direct estimates of the parameters of overidentified models as well as standard errors of those estimates. The missing-at-random assumption is necessarily satisfied when the remeasurement sample is randomly chosen from the larger sample. Although the assumption of multivariate normality may be problematic, the method we will describe can be easily adapted for estimation methods with less restrictive assumptions.

The method makes use of the fact that later versions of LISREL (4 through 7) and EQS 3.0 can simultaneously estimate the same model for two or more independent samples. Individual parameters can be

either constrained equal or allowed to vary across samples. For incomplete data problems, the basic idea is to divide the sample into subsamples, each having a distinct subset of variables present and missing. The model is then fit simultaneously to all subsamples, constraining corresponding parameters to be equal across subsamples.

That simple idea has two complications. First, to correctly maximize the likelihood function it is necessary to incorporate sample means into the analysis and to constrain the corresponding population means to be equal across samples. Some previous attempts to model incomplete data structures failed to recognize the need to structure the means (Hauser and Sewell 1986). In LISREL6, the analysis of structured means required an awkward reformulation of the model. Current versions of both programs allow for direct specification of mean structures. The second complication is specific to LISREL in all versions through 7. Because LISREL is designed to estimate the same model with the same number of variables in each sample, some special tricks are needed both to read in the data and specify the model (see Allison 1987). These techniques are unnecessary in EQS 3.0.

We now apply the method to the published data of BHF to get efficient estimates of the status attainment model of equations (1)-(3). For each of the two subsamples, the data were read into EQS 3.0 as they appear in Table 2. Zeros were substituted for missing correlations and ones were substituted for missing standard deviations, but any other numbers would have done as well. The only reason for including missing rows and columns is to keep the variable numbering consistent across the two samples.

For the remeasurement subsample, the model was specified to correspond to equations (1)-(4), with all the  $\lambda$ s fixed at 1.0. To incorporate the observed means, an intercept term was included in each of the measurement equations. In addition, AGE and AGE2 were specified as latent variables without any indicators. (A complete program listing is given in the Appendix.)

For the single-measurement subsample, the equations for the latent variables were exactly the same as in the remeasurement sample. Measurement equations were specified only for those variables that were actually observed, including AGE and AGE2. For these latter two variables, no error term was included in the equation to reflect

that these variables are assumed to be measured without error. Again, all  $\lambda$  parameters were fixed at 1.0. All parameters that were common to the two samples were constrained to be equal across samples.

Maximum likelihood estimates of the coefficients in equations (1)-(3) are reported in Column 4 of Table 3, together with their standard error estimates. All the estimates are reasonably close to those in Column 2, which were produced by the ad hoc procedure of BHF. In addition, the standard errors are considerably smaller than those in Column 3, which are based solely on the remeasurement subsample. They tend to be somewhat larger, however, than the standard errors in Column 1, which came from the uncorrected least-squares estimates.

For some of the coefficients in Table 3, the new estimates might lead to different conclusions. For example,  $\beta_{16}$  is not significantly different from zero in Columns 1 and 3, but is significant (at the .01 level) in Column 4. Similarly,  $\beta_9$  is not significant in Column 1 but is in Column 4. We have previously mentioned the inconsistent findings in Columns 1-3 for  $\beta_{10}$ , the effect of parental income on the status of the first occupation; in Column 4, the efficient estimate for this coefficient is strongly negative and highly significant ( $t = 11.38$ ). Clearly this unexpected result cannot be easily dismissed, but the interpretation is problematic. Although the estimates in Table 4 indicate that this variable is reliably measured, other research casts doubt on its validity. Featherman (1980) reports that an identical measure of parental income at age 16 had a correlation of only .28 with total personal income of the respondent's father in the census nearest the son's 16th birthday.

The reported value of the chi-square goodness-of-fit statistic for this model was 137.9 with 84 degrees of freedom, for a  $p$  value less than .001. (Unlike LISREL, EQS reports the correct degrees of freedom for incomplete data applications because it only operates on sample moments that are actually observed. With LISREL it is necessary to read in dummy values for missing variances and covariances, which increases the reported degrees of freedom [Allison 1987].) While this suggests rejection of the model, one must recall that the incomplete subsample has over 24,000 cases. Hence any restrictions on the covariance matrix are likely to be rejected at conventional levels. Moreover, the Bentler-Bonnett normed fit index is .998, indi-

cating that the fitted moments correspond quite closely to the observed moments.

There is a more compelling reason for not rejecting the model on the basis of this test. The chi-square statistic can be decomposed into two parts: (a) a part measuring the fit of the model *within* the two subsamples and (b) a part measuring departures from the constraints imposed *across* the two subsamples. For this application, part (a) applies only to the smaller subsample because the model is under-identified for the larger subsample. We have already reported a chi-square of 60.1 with 57 degrees of freedom for fitting the measurement model to the smaller subsample, and this corresponds to part (a). Part (b) may be found by subtracting part (a) from the values for the combined model.<sup>3</sup> Thus the null hypothesis that all the cross-subsample constraints are true yields a chi-square of 77.8 with 27 degrees of freedom. Again the  $p$  value is less than .001. This test may be interpreted as a test of the null hypothesis that the data are missing *completely* at random, against the alternative that they are merely missing at random.<sup>4</sup> For an explanation of the subtleties of this distinction, see Little and Rubin (1987, pp. 14-17). Rejection of the null hypothesis does not impugn the structural model or the parameter estimates because the portion of the likelihood function that varies between the null and the alternative hypothesis is *not* a function of the unknown parameters. Neither does rejection imply that the cross-subsample equality constraints should be relaxed (doing so would lead to biased parameter estimates). What it does suggest is that the remeasurement subsample may not be a true random subsample from the original sample.

Finally, we note that estimates of the error variances (and the derived reliability coefficients) change slightly when the model is fit to the entire sample. The new estimates are reported in Table 4.

#### OPTIMAL ALLOCATION OF CASES

We have demonstrated that a subsample with repeated measurements can be combined with a subsample lacking repeated measurements to produce efficient estimates of linear models. We now consider whether such designs offer any advantages over designs in which



all cases are remeasured. More precisely, what is the optimum allocation of cases between the remeasurement subsample and the subsample without remeasurement? As yet, we have no general, rigorous answer to this question. What we will do in this section is explore the question as it applies to the BHF data. The answer appears to be that, yes, it is more cost-effective to remeasure only a fraction of the sample; but, for this example, the fraction should be substantially higher than 578/25,223.

An obvious factor in determining the optimum allocation of cases is the cost of remeasurement. Let  $p$  denote the ratio of the cost of a remeasured case (including the initial measurement) to an unremeasured case. Clearly when  $p = 1$  (i.e., when there is no additional cost to remeasurement) the best design is to remeasure the entire sample. As  $p$  gets larger, however, we might expect that a smaller fraction of the cases should be remeasured. For the BHF study, the initial measurements were made by mail questionnaire while the remeasurements were done by telephone or personal interview, suggesting that  $p$  should be greater than 2.0. On the other hand, the cost of the initial measurements should include the cost of developing the sample itself, a cost that may well have been less for the remeasurement interview. As a rough approximation, then, we shall assume that  $p = 2.0$ .

Assuming that there is a fixed amount of money for data collection, we shall examine the effect on the standard errors of varying the number of cases allocated to the remeasurement and single-measurement subsamples. Letting  $n$  and  $m$  be, respectively, the number of cases in the remeasurement subsample and the single-measurement subsample, the constant-cost constraint requires that  $pn + m$  be held constant. Thus in the BHF study, we require that  $n$  and  $m$  be chosen so that

$$2n + m = 2(578) + 24,645 = 25,801. \quad (7)$$

If all cases were remeasured, then, an equally costly study would have 12,900.5 cases.

Another factor that should affect the allocation of cases is the degree of measurement of error in the independent variables. In the extreme, when there is no measurement error, there is no advantage to remeasurement. Thus in the case of AGE and AGE2, which are assumed to be error-free, the standard errors of the coefficients should be mini-

mized when no cases are remeasured. On the other hand, variables with lower reliability (i.e., a high ratio of error variance to total variance) might be expected to require a greater proportion of the cases allocated to the remeasurement subsample.

Our method for exploring the effect of differential allocation schemes on the standard errors is somewhat unusual. When the data are read into EQS, the number of cases must be specified for each of the two covariance matrices. Without changing the covariance matrices, we simply varied the number of cases specified for each subsample, subject to the restriction in equation (7), and then examined the estimated standard errors. Of course, this approach has obvious deficiencies. It only gives *estimated* standard errors, and it does not allow for the fact that observed covariance matrices would change if these hypothetical allocation schemes were actually implemented. Nevertheless, the covariance matrix for the large subsample of 24,645 is practically equivalent to the population matrix. Even the small subsample of 578 is large enough that the estimated covariance matrix should not be too far off from the population matrix. Thus we expect that this method should give a roughly accurate picture.

In Table 5, we present the standard errors calculated for each of the eighteen coefficients under ten different allocation schemes. For each coefficient, the minimum standard error is printed in boldface. For the coefficients of AGE and AGE2 in each of the three equations, the pattern is essentially what was expected: With the exception of a few coefficients in the first column, the standard errors are lowest when the fewest cases are assigned to be remeasured. Actually, this pattern is probably exaggerated in the present study, because AGE and AGE2 were missing in the correlation matrices reported by BHF for the remeasurement subsample. Consequently, the remeasured cases contributed no information at all about these two variables.

Although the minimum point varies for the other coefficients, it is generally in the range of 3,000 to 6,000 cases in the remeasurement subsample. This is well above the 578 that were actually remeasured. For many of the coefficients, however, the curve is relatively flat in the middle of the range, rising only at the two extremes. For example, the standard error of  $\beta_{11}$  is a constant .063 in the range of 4,000 to 7,000. With only a few exceptions, the standard errors for the 578

remeasured cases are only slightly larger than the minima. We do not discern any relationship between the minimum point and the degree of measurement error, possibly because there is not that much variation in the degree of measurement error. It must be stressed that the picture might change if  $p$  were to change substantially; a different value of  $p$  would mean that different allocation schemes would be permissible.

Although this method of investigating allocation is certainly not rigorous, in the absence of better methods it might be useful as a practical guide for research design. Using available data and making reasonable guesses, one could construct covariance matrices for the two subsamples. Then, using EQS as we have done here, one could get a rough estimate of the optimum allocation of cases.

### CONCLUSION

In this article we have provided a nearly optimal solution to the estimation problem posed by BHF: how to combine data from a single-measurement subsample and a remeasurement subsample to estimate the coefficients of a linear model. Our method produces efficient coefficient estimates and consistent estimates of their standard errors. It does so without introducing any assumptions beyond those made by BHF, and it can be implemented with widely available computer programs, EQS or LISREL.

We believe that the method described here should make remeasurement more attractive. No longer is it necessary, or even desirable, to remeasure the entire sample; a small subsample will suffice. True, the exploratory results in the previous section suggest that optimal designs may require the remeasurement subsample to be a substantial fraction of the total. Nonetheless, those who merely wish to explore the impact of measurement error on their estimates may be content with suboptimal designs in which the remeasurement subsample is small. They can now proceed with the confidence that good statistical procedures are available.

The principal limitation of this technique as a method for bias reduction lies in the plausibility (or lack thereof) of the measurement

TABLE 5: Standard Error Estimates Under Alternative Sample Allocations<sup>a</sup>

Variable	Coefficient	Number of Remeasurement Cases/Number of Single Measurement Cases									
		100/ 25601	578/ 24645	1000/ 23801	3000/ 19801	4000/ 17801	5000/ 15801	6000/ 13801	7000/ 11801	10000/ 5801	12000/ 1801
AGE	$\beta_1$	.014	.013	.013	.015	.015	.016	.017	.019	.027	.048
AGE2	$\beta_2$	.0010	.0011	.0011	.0012	.0012	.0013	.0014	.0015	.0022	.0039
FO	$\beta_3$	.0027	.0015	.0014	.0013	.0013	.0013	.0013	.0013	.0013	.0014
FE	$\beta_4$	.014	.0078	.0071	.0066	.0066	.0067	.0068	.0069	.0074	.0088
PI	$\beta_5$	.14	.076	.067	.059	.058	.059	.059	.060	.064	.071
AGE	$\beta_6$	.092	.094	.095	.11	.11	.12	.13	.14	.20	.36
AGE2	$\beta_7$	.0083	.0077	.0078	.0086	.0091	.0097	.010	.011	.016	.030
FO	$\beta_8$	.023	.012	.011	.0094	.0093	.0093	.0093	.0094	.0097	.010
FE	$\beta_9$	.077	.051	.049	.048	.048	.049	.050	.051	.055	.066
PI	$\beta_{10}$	.75	.48	.45	.44	.44	.44	.45	.46	.49	.55
ED	$\beta_{11}$	.21	.099	.082	.065	.063	.063	.063	.063	.067	.077
AGE	$\beta_{12}$	.13	.10	.10	.11	.11	.12	.13	.14	.20	.36
AGE2	$\beta_{13}$	.0091	.0079	.0079	.0086	.0091	.0097	.010	.011	.016	.030
FO	$\beta_{14}$	.016	.011	.0099	.0095	.0095	.0096	.0097	.0098	.010	.012
FE	$\beta_{15}$	.051	.046	.046	.047	.048	.049	.050	.051	.057	.070
PI	$\beta_{16}$	.61	.45	.44	.43	.44	.44	.45	.46	.50	.59
ED	$\beta_{17}$	.33	.15	.13	.096	.093	.091	.090	.091	.095	.11
O1	$\beta_{18}$	.054	.024	.019	.013	.012	.012	.012	.011	.012	.014

<sup>a</sup>In each row, the smallest standard error(s) are printed in boldface.

error model. Measurement error is assumed to be random and uncorrelated across time for the same variable. In cases where these assumptions are not met, the strategy of remeasurement will only partially correct for measurement error. Although it is possible to estimate models with correlated errors (BHF), such models usually require additional restrictions of dubious validity.

### APPENDIX

Following is a listing of the EQS 3.0 control statements which produced the estimates in Table 3, Column 4. Variables labeled V are observed variables, and their numbering follows the same order as in Table 4. The Fs are latent variables, and the Es and Ds are random disturbances. V999 is a vector of ones, used to specify intercepts in the measurement equations. Numbers followed by asterisks are starting values for free parameters. Starting values were obtained by approximating the estimates in Table 3, Column 3. A file containing these control cards can be obtained from the first author either by electronic mail (ALLISON@PENNDRLS.BITNET) or on diskette by regular mail (Department of Sociology, University of Pennsylvania, 3718 Locust Walk, Philadelphia, PA 19104-6299).

/TITLE

OCG REMEASURED SUBSAMPLE

/SPECIFICATIONS

VARIABLES=15; CASES=578; ANALYSIS=MOM; GROUPS=2;  
MATRIX=COR;

/EQUATIONS

V1 = 33\*V999 + F1 + E1;  
V2 = 33\*V999 + F1 + E2;  
V3 = 9\*V999 + F2 + E3;  
V4 = 9\*V999 + F2 + E4;  
V5 = 4\*V999 + F3 + E5;  
V6 = 4\*V999 + F3 + E6;  
V7 = 12\*V999 + F4 + E7;  
V8 = 12\*V999 + F4 + E8;  
V9 = 12\*V999 + F4 + E9;  
V10 = 33\*V999 + F5 + E10;  
V11 = 33\*V999 + F5 + E11;



## /STANDARD DEVIATIONS

24.27 23.73 4.19 4.14 .41 .39 3.42 2.93 2.87 24.71 24.15 24.81 25.21 1.0 1.0

## /MEANS

32.96 33.6 8.97 8.96 3.78 3.8 11.98 12.12 12.18 34.6 32.1 39.57 41.34 0 0

## /END

## /TITLE

OCG SINGLE MEASURE SUBSAMPLE

## /SPECIFICATIONS

VARIABLES = 15; CASES= 24645; ANALYSIS=MOM; MATRIX=COR;

## /EQUATIONS

$$V1 = 33 * V999 + F1 + E1;$$

$$V3 = 9 * V999 + F2 + E3;$$

$$V5 = 4 * V999 + F3 + E5;$$

$$V9 = 12 * V999 + F4 + E9;$$

$$V10 = 33 * V999 + F5 + E10;$$

$$V13 = 40 * V999 + F6 + E13;$$

$$V14 = 3.97 * V999 + F7;$$

$$V15 = 16.04 * V999 + F8;$$

$$F4 = 2.25 * F3 + .152 * F2 + .024 * F1 - .03 * F7 - .02 * F8 + D1;$$

$$F5 = 0 * F3 + 0 * F2 + .200 * F1 + 2.0 * F7 - .12 * F8 + 5.8 * F4 + D2;$$

$$F6 = 0 * F3 + 0 * F2 + .050 * F1 + 2.7 * F7 - .13 * F8 + 2.4 * F4 + .52 * F5 + D3;$$

## /VARIANCES

$$F1 = 500 *; F2 = 16 *; F3 = .15 *;$$

$$F7 = 1.5625 *; F8 = 214.037 *;$$

$$E1 = 87 *; E3 = 1.22 *; E5 = .02 *;$$

$$E9 = .93 *; E10 = 92 *; E13 = 101 *;$$

$$D1 = 4.5 *; D2 = 231 *; D3 = 221 *;$$

## /COVARIANCES

$$F1, F2 = 60 *;$$

$$F1, F3 = 4 *;$$

$$F1, F7 = -4.5 *;$$

$$F1, F8 = 4.3 *;$$

$$F2, F3 = .8 *;$$

$$F2, F7 = -1.4 *;$$

$$F2, F8 = 1.5 *;$$

$$F3, F7 = -.124 *;$$

$$F3, F8 = -.15 *;$$

$$F7, F8 = 2.6334 *;$$

## /MATRIX

1.0

0 1.0

.536 0 1.0

0 0 0 1.0

.399 0 .466 0 1.0

0 0 0 0 0 1.0

0 0 0 0 0 0 1.0

0 0 0 0 0 0 0 1.0

.410 0 .470 0 .483 0 0 0 1.0

.391 0 .331 0 .291 0 0 0 .636 1.0

0 0 0 0 0 0 0 0 0 1.0

0 0 0 0 0 0 0 0 0 0 1.0

.325 0 .275 0 .256 0 0 0 .571 .617 0 0 1.0

-.174 0 -.297 0 -.248 0 0 0 -.210 -.067 0 0 .025 1.0

.014 0 .026 0 -.027 0 0 0 -.095 -.114 0 0 -.142 .144 1.0

## /STANDARD DEVIATIONS

20.90 1.0 3.88 1.0 .40 1.0 1.0 1.0 2.91 22.48 1.0 1.0 22.78 1.25 14.63

## /MEANS

31.09 0 8.78 0 3.77 0 0 0 12.07 33.81 0 0 41.11 3.97 16.04

## /CONSTRAINTS

(1,V1,V999) = (2,V1,V999);

(1,V3,V999) = (2,V3,V999);

(1,V5,V999) = (2,V5,V999);

(1,V9,V999) = (2,V9,V999);

(1,V10,V999) = (2,V10,V999);

(1,V13,V999) = (2,V13,V999);

(1,F1,F1) = (2,F1,F1);

(1,F1,F2) = (2,F1,F2);

(1,F1,F3) = (2,F1,F3);

(1,F1,F7) = (2,F1,F7);

(1,F1,F8) = (2,F1,F8);

(1,F2,F2) = (2,F2,F2);

(1,F2,F3) = (2,F2,F3);

(1,F2,F7) = (2,F2,F7);

(1,F2,F8) = (2,F2,F8);

(1,F3,F3) = (2,F3,F3);

(1,F3,F7) = (2,F3,F7);

(1,F3,F8) = (2,F3,F8);



```

(1,F6,F1) = (2,F6,F1);
(1,F6,F2) = (2,F6,F2);
(1,F6,F3) = (2,F6,F3);
(1,F6,F4) = (2,F6,F4);
(1,F6,F5) = (2,F6,F5);
(1,F6,F7) = (2,F6,F7);
(1,F6,F8) = (2,F6,F8);
(1,F5,F1) = (2,F5,F1);
(1,F5,F2) = (2,F5,F2);
(1,F5,F3) = (2,F5,F3);
(1,F5,F4) = (2,F5,F4);
(1,F5,F7) = (2,F5,F7);
(1,F5,F8) = (2,F5,F8);
(1,F4,F1) = (2,F4,F1);
(1,F4,F2) = (2,F4,F2);
(1,F4,F3) = (2,F4,F3);
(1,F4,F7) = (2,F4,F7);
(1,F4,F8) = (2,F4,F8);
(1,D1,D1) = (2,D1,D1);
(1,D2,D2) = (2,D2,D2);
(1,D3,D3) = (2,D3,D3);
(1,E1,E1) = (2,E1,E1);
(1,E3,E3) = (2,E3,E3);
(1,E5,E5) = (2,E5,E5);
(1,E9,E9) = (2,E9,E9);
(1,E10,E10) = (2,E10,E10);
(1,E13,E13) = (2,E13,E13);
(1,F7,F7) = (2,F7,F7);
(1,F8,F8) = (2,F8,F8);
(1,F7,F8) = (2,F7,F8);

/PRINT
DIG=5;

/END

```

## NOTES

1. In Versions 7.16 and earlier, convergence could not be obtained when structured means were incorporated into the model, due to program bugs.

2. The corrections consisted of (a) converting the reported correlations and standard deviations into sums of squares and cross-products, (b) subtracting the sums for the remeasurement subsample from those for the full sample, and (c) recalculating the correlations and standard deviations. These corrections are only approximate because the matrices reported in BHF are actually pairwise-present matrices, so that different correlations may be based on different numbers of cases. Our correction formulas assumed that all correlations were based on either 25,223 or 578 cases.

3. An alternative way of calculating part (b) is to fit the model to both subsamples in such a way that the only constraints are those across subsamples.

4. If the data are missing completely at random, the two subsamples can be regarded as distinct random samples from the same population. Hence all corresponding parameters should be equal. On the other hand, the assumption that the data are missing at random but not observed at random is compatible with any cross-group differences in parameters (Rubin 1976). In other words, the missing-at-random assumption imposes no restrictions on the observed moments.

## REFERENCES

- Allison, P. D. 1987. "Estimation of Linear Models with Incomplete Data." Pp. 71-103 in *Sociological Methodology 1987*, edited by Clifford Clogg. Washington, DC: American Sociological Association.
- Bentler, P. M. 1989. *EQS Structural Equations Program Manual*. Los Angeles: BMDP Statistical Software.
- Bentler, P. M. and M. Berkane. 1985. "Developments in the Elliptical Theory Generalization of Normal Multivariate Analysis." *Proceedings of the Social Statistics Section, American Statistical Association* 291-295.
- Bentler, P. M. and D. G. Weeks. 1980. "Linear Structural Equations with Latent Variables." *Psychometrika* 45:289-308.
- Bielby, W. T., R. M. Hauser, and D. L. Featherman. 1977a. "Response Errors of Nonblack Males in Models of the Stratification Process." *Journal of the American Statistical Association* 72:723-35.
- Bielby, W. T., R. M. Hauser, and D. L. Featherman. 1977b. "Response Errors of Black and Nonblack Males in Models of the Intergenerational Transmission of Social Status." *American Journal of Sociology* 82:1242-88.
- Bowles, S. 1972. "Schooling and Inequality from Generation to Generation." *Journal of Political Economy* 80:S219-51.
- Bowles, S. and V. Nelson. 1974. "The 'Inheritance of IQ' and the Intergenerational Reproduction of Economic Inequality." *Review of Economics and Statistics* 56:39-51.
- Duncan, O. D. 1961. "A Socioeconomic Index for All Occupations." Pp. 109-38 in *Occupations and Social Status*, edited by A. J. Reiss, Jr. New York: Free Press.
- Featherman, D. L. 1980. "Retrospective Longitudinal Research: Methodological Considerations." *Journal of Economics and Business* 32:152-69.
- Hauser, R. M. and W. H. Sewell. 1986. "Simple Models of Education, Occupational Status, and Earnings: Findings from the Wisconsin and Kalamazoo Studies." *Journal of Labor Economics* 4:S83-115.
- Intriligator, M. D. 1978. *Econometric Models, Techniques, and Applications*. Englewood Cliffs, NJ: Prentice-Hall.

- Jöreskog, K. G. 1969. "A General Approach to Confirmatory Maximum Likelihood Factor Analysis." *Psychometrika* 34:183-202.
- Jöreskog, K. G. and D. Sörbom. 1979. *Advances in Factor Analysis and Structural Equation Models*. Cambridge, MA: Abt Associates.
- Jöreskog, K. G. and D. Sörbom. 1988. *LISREL 7: A Guide to the Program and Applications*. Chicago: SPSS.
- Little, R.J.A. and D. B. Rubin. 1987. *Statistical Analysis with Missing Data*. New York: Wiley.
- Muthén, B., D. Kaplan, and M. Hollis. 1987. "On Structural Equation Modeling with Data That Are Not Missing at Random." *Psychometrika* 52:431-62.
- Rubin, D. B. 1976. "Inference and Missing Data." *Biometrika* 63:581-92.
- Shapiro, A. and M. W. Browne. 1987. "Analysis of Covariance Structures Under Elliptical Distributions." *Journal of the American Statistical Association* 82:1092-97.
- Werts, C. E., D. A. Rock, and J. Grandy. 1979. "Confirmatory Factor Analysis Applications: Missing Data Problems and Comparison of Path Models Between Populations." *Multivariate Behavioral Research* 14:199-213.

*Paul D. Allison is a professor of sociology at the University of Pennsylvania. In addition to methodological work, he is currently developing models of cultural evolution, especially as applied to the origin and persistence of altruistic behavior. Recent publications include "Change Scores as Dependent Variables in Regression Analysis," Sociological Methodology 1990, and "Departmental Effects on Scientific Productivity," American Sociological Review (with J. Scott Long, August 1990).*

*Robert M. Hauser is the Vilas Research Professor of Sociology at the University of Wisconsin-Madison. He is a demographer and social statistician with interests in education and social mobility. His recent work on sibling resemblance in schooling (with Raymond Sin-Kwok Wong) appeared in Sociology of Education. He chaired the Panel on Education of the National Research Council's Committee on the Status of Black Americans; its report, "A Common Destiny: Blacks in American Society," was the subject of a recent exchange between Hauser (with Gerald Jaynes and Robin M. Williams, Jr.) and Richard Herrnstein in The Public Interest.*