# Measurement Error in the Current Population Survey: A Nonparametric Look

Christopher R. Bollinger, *Georgia State University*

This article utilizes an exact match file between the 1978 March Current Population Survey and administrative records from the Social Security Administration to analyze errors in the reporting of annual income using nonparametric methodology. The article extends work of Bound and Krueger, and the results confirm many of the findings in Bound and Krueger. Three new findings are of interest: there is higher measurement error in cross-sectional samples than in panels. The negative relationship between measurement error and earnings is driven largely by overreporting among low earners. Median response errors are not related to earnings.

## I. Introduction

Econometricians have understood the problem of measurement error in survey data for many years. Aigner et al. (1984) and Fuller (1987) provide excellent surveys of the current literature. Incorrect responses to a survey question can dramatically bias even simple estimation. The information necessary to correct this bias is attainable through validation data: These are data where survey responses can be compared with an independent and presumably error-free source of the same information. The focus here is on the distribution of the response, conditional on the

true value of earnings and other potential correlates. Rather than specify a parametric or even semiparametric relationship, this research uses nonparametric methodology. The advantage of this methodology is that results are not biased because of specification error. The disadvantage is a loss in efficiency (compared with a known specification) and an inability to extrapolate beyond the support of the observed data. This study finds that reported earnings in the Current Population Survey (CPS) have a measurement error that is systematically related to true earnings and gender.

Many validation studies have confirmed the presence of measurement error in survey data. Rogers and Herzog (1987) find reporting errors in the Study of Michigan Generations Project data. Poterba and Summers (1986) find reporting errors in CPS measures of employment status. Greenberg and Halsey (1983) and Mathiowetz and Duncan (1988) find similar results in the Seattle and Denver Income Maintenance Experiment (SIME/DIME) data. Freeman (1984) and Card (1991) find errors in reporting union status in CPS data. Marquis and Moore (1990) and Bollinger and David (1997a) find misreporting of government assistant program participation in the Survey of Income and Program Participation (SIPP). Mellow and Sider (1983), Duncan and Hill (1985), and Bound et al. (1990) find measurement error in the Panel Survey of Income Dynamics (PSID).

Bound and Krueger (1991) utilized the 1978 Current Population Survey-Social Security Administration Exact Match File (CPS-SSA) data to examine the extent of measurement error in earnings reports. While the work of Bound and Krueger is important and useful, a number of extensions are warranted. The authors focus on determining the effect of measurement error on estimated wage equations. A more general understanding of the relationship between true and reported earnings yields implications for other research. Second, they utilized fully parametric estimation. Nonparametric estimation is not subject to specification bias. Finally, the sample was limited to heads of households in the March 1978 CPS who could be matched to their March 1977 CPS responses. This limits the general applicability of the sample.

This article extends the work of Bound and Krueger (1991). The research focuses on an expanded sample including women who were not heads of households and examines a cross-sectional sample in addition to the more restrictive panel structure. Further, nonparametric methodology is used. The results confirm many of those reported in Bound and Krueger: response error is negatively related to earnings, there is more measurement error among men than women, and measurement error is not related to age, education, and weeks worked. Three new findings are of interest. First, comparing cross-sectional samples to the panel samples examined by Bound and Krueger indicates higher measurement error in

the cross-sectional samples. The more stringent requirements of constructing a panel from CPS apparently leads to selecting individuals who are better reporters. Second, the negative relationship between measurement error and earnings is driven largely by a concentration of overreporting among low earners. Finally, while the average response error is negatively related to earnings, the median response error is zero for all levels of earnings, hence median regression will be more robust to response error when earnings are a left-hand side variable.

Section II of this article discusses the data. Section IIA describes the construction of the data sets and compares the data in this study with the data utilized by Bound and Krueger (1991). Section IIB presents specific comparisons of the data sets. Measurement error for men is compared with that for women, and measurement error in simple cross-sectional samples is compared with measurement error in panel samples.

Section III of this article presents the empirical findings. Section IIIB presents nonparametric estimates of the expectation of reported income conditional on Social Security Administration income and other variables. Section IIIC presents nonparametric estimates of the conditional median function. Section IIID summarizes results of nonparametric estimates of the probability of correctly reporting income.

Section IV concludes the article. A discussion of situations when nonparametric methods may be appropriate and potential solutions for the measurement error problem are presented.

## II. Description of Data

### A. Construction and Choice of Data Sets

In an effort to improve the accuracy of the data collected for the Current Population Survey, the Bureau of the Census with the cooperation of the Social Security Administration (SSA) conducted an "exact match" study: CPS data were matched to SSA data. The resulting data set contains the usual Annual Demographic File variables, and data supplied by the SSA including earnings for 1977, 1976, and previous years. This research utilizes the public use files generated from the 1978 CPS-SSA exact match file. Approximately 50% of the original CPS sample was successfully matched. Also used are data from the 1977 CPS Annual Demographic File (ADF).

Bound and Krueger (1991) examined response error utilizing these data. They constructed a panel sample using the 1978 CPS-SSA data and the 1977 March CPS ADF. The sample focused on heads of households only. Details on construction of the sample are found in Bound and Krueger. The samples used here are constructed identically to those used in Bound and Krueger with two exceptions. First, women who are either the householder or the spouse of the householder are included. Hence

the samples of women here are not limited only to female heads of household. The second major difference is that cross-sectional samples are constructed in addition to panel samples similar to those of Bound and Krueger. The cross-sectional samples are the main focus of this study.

Four samples are considered. Sample 1 selects male heads of household who are present in both the 1977 March ADF and in the 1978 CPS-SSA file. Sample 2 contains females who are either heads of household or wives of the head of the household and are present in both the 1977 March ADF and the 1978 CPS-SSA file. These first two samples are the panel samples. Sample 3 includes male heads of household from the 1978 CPS-SSA file selected identically to sample 1 except not matched to the 1977 data. Sample 4 contains women (heads of household and wives) selected identically to sample 2 but not matched to the 1977 data. The first and second sample are proper subsets of the third and fourth sample (respectively). Table 1 presents the sample size, mean, median, and standard deviation for selected variables in each of the four samples.

A number of potential problems in using earnings as reported to the SSA as "true earnings" exist. The earnings data reported by the SSA are censored at $16,500 for 1977 income and at $15,300 for 1976 income. Less than 3% of the women earn more than the threshold. Approximately 40% of the men earn more than the threshold. The voluntary aspect of participation in the match study leads to potential selection bias. Bound and Krueger (1991) addressed the top coding problem by modeling the response error as a Tobit and using maximum likelihood estimation. The approach here is to use the SSA earnings as a right-hand-side variable. Hence, no selection bias exists, but extrapolation to higher values of earnings is tenuous. These issues are discussed at length in Bound and Krueger.

The voluntary aspect of participation in the match study leads to potential selection bias. One would expect that those individuals who refuse to participate in the match to SSA data are most likely to have knowingly misrepresented their income. Bound and Krueger (1991) examined the potential selection bias. They concluded that this selection was unlikely to result in bias to their parameter estimates. However, it is wise to treat the results presented here and in Bound and Krueger as potentially understating the extent of the problem.

An unresolved issue is the potential of earnings not reported to SSA but reported to CPS. There is no way of determining these "under-the-table" earnings. Under-the-table earnings may explain the finding (discussed below) that low-income men overreport earnings to the CPS. If low-income men have a high rate of participation in second jobs, especially self-employment or in industries not covered by SSA (or where tips and other unreported income occurs), then the results below may not be due to errors in the CPS but rather to errors in the SSA earnings

Table 1
Descriptive Statistics of Four Samples

| Variables | Sample 1: Male Panel | Sample 2: Female Panel | Sample 3: Male Cross Section | Sample 4: Female Cross Section |
|---|---|---|---|---|
| Age: | | | | |
| M | 43.1 | 43.0 | 41.1 | 40.9 |
| Median | 42 | 42 | 39 | 39 |
| SD | 12.19 | 12.95 | 13.0 | 14.2 |
| 1976 CPS income: | | | | |
| M | 14,969 | 6,545 | | |
| Median | 14,085 | 6,166 | . . . | . . . |
| SD | 7,323 | 4,197 | | |
| 1977 CPS income: | | | | |
| M | 16,128 | 7,055 | 15,548 | 7,022 |
| Median | 15,049 | 6,705 | 14,875 | 6,500 |
| SD | 8,004 | 4,271 | 8,272 | 4,512 |
| 1976 SSA income: | | | | |
| M | 12,538 | 6,601 | 11,777 | 6,213 |
| Median | 14,408 | 6,482 | 13,367 | 5,954 |
| SD | 3,562 | 3,649 | 4,125 | 3,823 |
| 1977 SSA income: | | | | |
| M | 13,401 | 7,046 | 12,894 | 6,968 |
| Median | 15,436 | 6,841 | 14,665 | 6,606 |
| SD | 4,003 | 4,015 | 4,340 | 4,126 |
| Education: | | | | |
| M | 11.94 | 11.8 | 12.1 | 11.9 |
| Median | 12 | 12 | 12 | 12 |
| SD | 2.95 | 2.32 | 3.0 | 2.4 |
| White: | | | | |
| M | .93 | .91 | .93 | .86 |
| Median | 1 | 1 | 1 | 1 |
| SD | .26 | .29 | 2.6 | .32 |
| Weeks worked 1976: | | | | |
| M | 49.7 | 45.5 | | |
| Median | 52 | 52 | . . . | . . . |
| SD | 7.1 | 12.7 | | |
| Weeks worked 1977: | | | | |
| M | 49.5 | 45.3 | 48.8 | 44.7 |
| Median | 52 | 52 | 52 | 52 |
| SD | 7.6 | 12.9 | 8.8 | 13.4 |
| SSA77 < 16,500: | | | | |
| M | .56 | .98 | .60 | .97 |
| Median | 1 | 1 | 1 | 1 |
| SD | .50 | .45 | .49 | .17 |
| Age < 70 (years): | | | | |
| M | .99 | .99 | .99 | .98 |
| Median | 1 | 1 | 1 | 1 |
| SD | .10 | .10 | .10 | .14 |
| N | 2,338 | 1,566 | 7,380 | 6,499 |

NOTE.—CPS = Current Population Survey. SSA = Social Security Administration.

data. This issue was examined using both the May 1985 CPS and the May 1991 CPS, which ask supplemental questions concerning second jobs. The detailed results are available from me on request. The analysis revealed that there is no concentration of second-job holders among low-

income males. Although the analysis cannot, with any certainty, establish the validity of the results here, it does not present any evidence that the CPS data would be more accurate than the SSA data for low-income workers.

## B. Comparison of Samples

This section examines differences in the mean and variance of the reporting error between the four data sets. The reporting error is defined as the reported earnings in the CPS minus the SSA reported earnings: $e = CPS - SSE$. Two sets of comparisons are of particular importance. First, the panel data sets (samples 1 and 2) are compared with the cross-sectional data sets (samples 3 and 4). This comparison tests the hypothesis that use of only the panel data sets results in selecting individuals who are better reporters than those in the population as a whole. Second, reporting errors are compared between men and women.

Restricting the male sample to only those who could be matched to the 1977 CPS-ADF results in smaller mean and variance of $e$. The null hypothesis that the mean error is the same between samples 1 and 3 is rejected at all standard confidence intervals, while the null hypothesis that the variances are the same cannot be rejected. A joint test of the means and variances rejects the null hypothesis that both are the same at all standard significance levels.[1] The panel sample of women (sample 2) has a lower mean and variance of the errors compared with the cross-sectional sample (sample 4). The test of the differences in means could not reject the null hypothesis that the means were the same, but the hypothesis that the variances are the same was rejected. The joint hypothesis is only rejected at the 10% significance level.[2]

Two important conclusions can be drawn. First, individuals who can be matched between the surveys are better reporters, hence complete panels are less subject to measurement error than cross-sectional data (see also Bollinger and David 1997*b*). This may be the result of additional checks necessary for constructing a panel data set or because individuals who remain in a panel over the long term are more likely to cooperate and provide accurate data. Second, since most practitioners utilize the CPS as a cross-sectional sample, the work of Bound and Krueger (1991)

[1] The asymptotically standard normal test of difference in mean error between samples 1 and 3 had a value of $-2.8$. The test of differences in variance had a value of $-.85$. The asymptotically $\chi^2$ test of the joint hypothesis had a value of 9.6.

[2] The asymptotically standard normal test for differences in the mean error between samples 2 and 4 had a value of $-.64$, while the test of differences in variance had a value of $-2.26$. The asymptotically $\chi^2$ test of the joint hypothesis had a value of 5.16.

needs extending in this direction. For this reason, the remainder of this article focuses on samples 3 and 4.

Comparing men with women (sample 3 vs. sample 4) soundly rejects the null hypothesis that the errors for men and women have the same mean and variance.[3] The results indicate that men have a higher mean error and a higher variance of the error: men are less accurate reporters than women.

A final aspect worth investigating is the normality of the errors. A simple test of normality is to compare the third central moment with zero (skewness) and the fourth central moment to three times the square of the variance (kurtosis). Tests were performed on the natural log errors in samples 3 and 4. Since both samples contain a large proportion of respondents who report exactly (11.7% of the men and 12.7% of the women), tests were also performed for the samples without these respondents. The results of the tests soundly reject the null hypothesis of normality for the log errors.[4] The errors are found to be asymmetrically distributed (positively) and have thick tails (leptokurtic). Goldberger (1983) examines the bias in Tobit estimation when the assumption of normality is violated. He finds that even minor deviations from normality (symmetric distributions) can lead to rather large bias. The large asymmetry of the error distribution apparent here suggests that analysis utilizing Tobit estimation may lead to a large bias. Using nonparametric methodology rather than the maximum likelihood estimation employed by Bound and Krueger (1991) may lead to different conclusions. In this case, no conclusions are overturned; however, this cannot be determined a priori.

### III. Empirical Findings

#### A. Empirical Methodology

One approach to analyzing the measurement error process is to consider the traditional definition of measurement error: $e = CPS - SSE$. The descriptive statistics above were constructed using this definition. The more general approach taken here can be viewed as a structural analysis. The response to the question "Last year (1977) did [the respondent] receive any money in wages or salary" (question 51a from CPS

---

[3] The asymptotically standard normal test for differences in mean error between samples 3 and 4 has a value of 9.7, while the test for differences in variance has a value of 9.4. The asymptotically $\chi^2$ test of the joint hypothesis has a value of 112.5.

[4] For the men, the third-moment test value is 79, the fourth-moment test value is 403, and the joint test value is 168,157. For the women, the test values are 58, 40, and 166,435, respectively. For the men, without exact reporters, the test values are 69, 331, and 114,342, respectively. For the women, without exact reporters, the test values are 54, 33, and 109,556, respectively.

questionnaire) is structurally determined by the true value of income (*SSE*) and perhaps other factors. Hence the regressions of interest are of the form

$$E\,[CPS|\,SSE,\,X\,] = f(\,SSE,\,X\,) \tag{1}$$

The function $f(\cdot)$ represents any systematic reporting behavior, while the "residual" from these is stochastic noise. Classical measurement error analysis assumes that $f(\,SSE,\,X\,) = SSE$ and focuses on the residual. The approach here focuses on the function $f(\cdot)$.

The interested reader will find Bierens (1987) and Härdle (1990) excellent texts on nonparametric methodology. The estimation in this work utilizes Nayadara-Watson kernel-regression estimators. The estimator, for the simple case of only the *SSE* variable, is

$$\hat{E}[CPS|\,SSE = x] = \frac{\displaystyle\sum_{i=1}^{N} CPS_i K\left(\frac{SSE_i - x}{h}\right)}{\displaystyle\sum_{i=1}^{N} K\left(\frac{SSE_i - x}{h}\right)},$$

where $x$ is the chosen value of *SSE*, $N$ is the sample size, $h$ is the bandwidth, and $K(\cdot)$ is the kernel function. The Epanechnikov function is utilized for the kernel function:

$$K(u) = (1 - u^2)*I[u^2 < 1],$$

where $I[\cdot]$ is the indicator function. As noted in Härdle (1990), the choice of the kernel function does not have dramatic effects on the estimates or the sampling properties of the estimates. Choice of the bandwidth, however, is crucial. The bandwidth is chosen via leave-one-out least squares cross validation.

In nonparametric regression, it is only necessary to assume that the function, $f(\cdot)$, be twice continuously differentiable. The estimator reveals to the researcher any nonlinearities present without requiring the research to incur the "pretest" bias associated with arriving at a functional form through the trial and error of various nonlinear specifications. The cost is slower convergence ($\sqrt{Nh}$) and thus less precise estimates at any given sample size (compared with $\sqrt{N}$ consistent estimators such as ordinary least squares [OLS]).

The nonparametric approach is well suited to problems where the data set is relatively large and the researcher has particular interest in observing the true relationship rather than testing specific parameters of a model.

Thus, the question asked in this work is not simply, "What is the bias in OLS from the measurement error?" but rather, "What is the relationship between the truth and the observed response, on average?" In addition to the potential bias from misspecification of the Tobit model when the normality assumption is violated (see the discussion above), linear models may fail to capture important nonlinearity in the relationship between the survey response and the true income. While these nonlinearities may be "found" by experimentation with functional forms, this approach is ad hoc, leads to pretest bias, and may still fail to "find" the true functional relationship. By utilizing nonparametric estimation, the nonlinear relationship between response and income for men was found.

Nonparametric estimates are also used to test particular specifications of the function $f(\cdot)$. Below, the specification where $f(\cdot)$ is just the identity function and where $f(\cdot)$ is the linear projection are both tested. If $f(\cdot)$ is the identity function, then the response error is an additive white-noise term. The additive white-noise measurement error process is often assumed in treatment of measurement error (see, e.g., Fuller 1987). If this assumption is not borne out in practice, new approaches to addressing measurement error may be necessary. If the true function $f(\cdot)$ is the linear projection function, then the correlation coefficients presented in Bound and Krueger (1991) are an excellent summary of the relationship between the response and the true income. Otherwise, more complicated modeling is necessary.

### B. Estimation of the Conditional Mean Function

Figures 1 and 2 present the basic result for the cross-sectional samples (samples 3 and 4). The figures present the estimation of $E[CPS|SSE]$. In each figure, the solid line represents the estimate of the conditional expectation of the CPS earnings response given the level of SSE. The large dashed line represents the 45° line. The two smaller dashed lines represent the 95% confidence region for the estimated line. The confidence band is derived using asymptotic results (Härdle 1990, p. 116–17). Ninety-five percent of all such bands would contain the true conditional expectation function.

Comparison of the estimates for men and women supports the finding in section IIB that women are better reporters than men. The error is higher for men than women at all income levels, and the response of reported earnings to changes in actual earnings is flatter for men than women.

If measurement error were not related to earnings, the solid line would coincide with the 45° line. This would imply that $E[CPS|SSE] = SSE$. The results in figures 1 and 2 suggest that this is likely not the case for men, but potentially the case for women. The null hypothesis that $E[CPS|SSE] = SSE$ is tested for both samples using a test proposed by
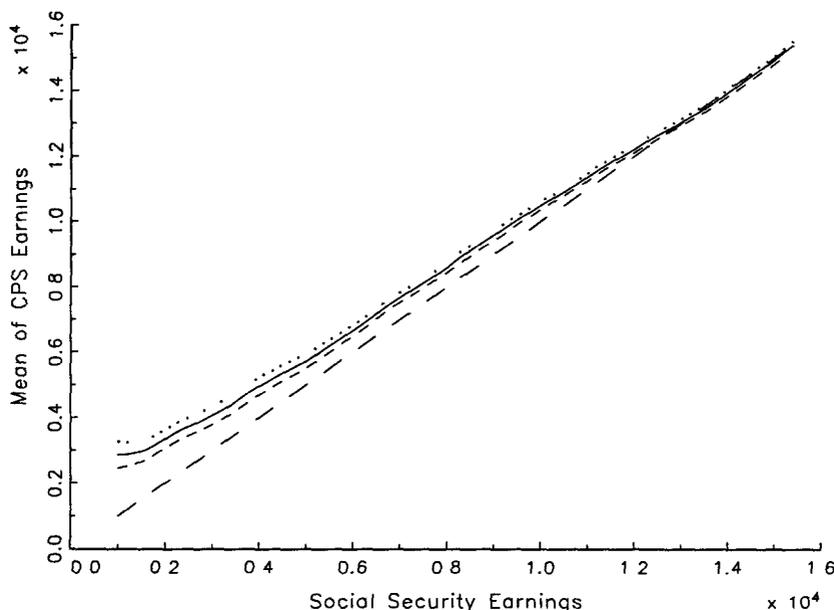
FIG. 1.—CPS earnings conditional on SSA earnings of male head of household. Income in ten thousands of dollars. Key: — = conditional mean, – – = 45° line, - - - = lower 95% confidence bound, and · · · = upper 95% confidence bound.

Lee (1994). The test compares the sum of squared residuals from the parametric specification with the nonparametric estimate. The null hypothesis was rejected for men, but not for women.[5]

Figures 1 and 2 suggest that a linear specification might be appropriate for both men and women. This specification was tested using the test proposed by Lee (1994). The null hypothesis that $E[CPS|SSE] = \alpha + \beta SSE$ was rejected for the sample of men, but not rejected for the sample of women.[6] It is particularly interesting to note that the fit of the OLS line for the men is poorest at low values of SSA earnings. Retesting the hypothesis excluding men with SSA income below $3,000, the null hypothesis that the specification is linear is no longer rejected.[7] These

[5] The value of the test for men was −2.94, and for women −.69. Under the alternative hypothesis, the test must be negative; hence, a one-sided critical value is utilized. The $p$-value for the men is .0016, while for the women it is .2451.

[6] The OLS estimation resulted in a slope of .91 and an intercept of 1,364 for men. For women, the slope is .97, and the intercept is 211. The test value for the men was −2.52, while the test value for the women was −.51. The respective $p$-values were .0059 and .305.

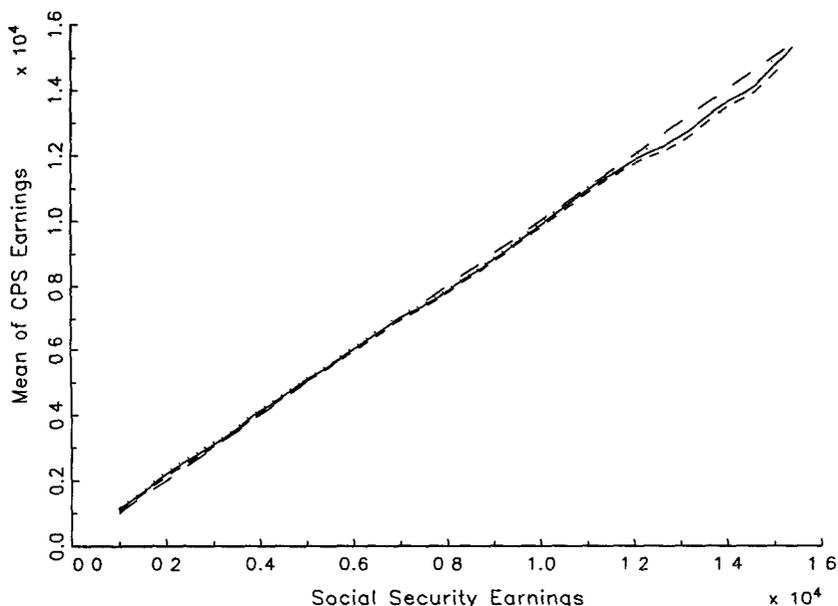[7] The resulting test statistic is −1.15, with a $p$-value of .1251.

FIG. 2.—CPS earnings conditional on SSA earnings of female head and spouse of Head of Household. Income in ten thousands of dollars. Key: — = conditional mean, – – = 45° line, - - - = lower 95% confidence bound, and · · · = upper 95% confidence bound.

results imply that while a linear fit may work well for some range of the data, it would fail, at least for the men, to account for nonlinearity in the lowest ranges of the SSE data. This result highlights the advantage of utilizing nonparametric methodology: a priori it is difficult to know a proper functional form specification.

The finding here that there appears to be a "regression to the mean" of response to the income question conforms with that of Bound and Krueger (1991). Further, the finding that the slope of the response is flatter for men than women also conforms with the results of Bound and Krueger. The nonlinearity for the men at low incomes and establishing the linear relationship for women were not addressed in Bound and Krueger. It should be noted that the methodology of Bound and Krueger is not strictly comparable to the approach here. They are concerned with the question of predicting true income given the observed income. This study focuses on the conditional expectation of observed income given true income. Thus the focus is how true income determines the response of the individual rather than the prediction problem.

In cases where earnings is the dependent variable, the response of men's earnings to variables such as education has a larger downward bias than women's earnings. The differences in average overreporting discussed in
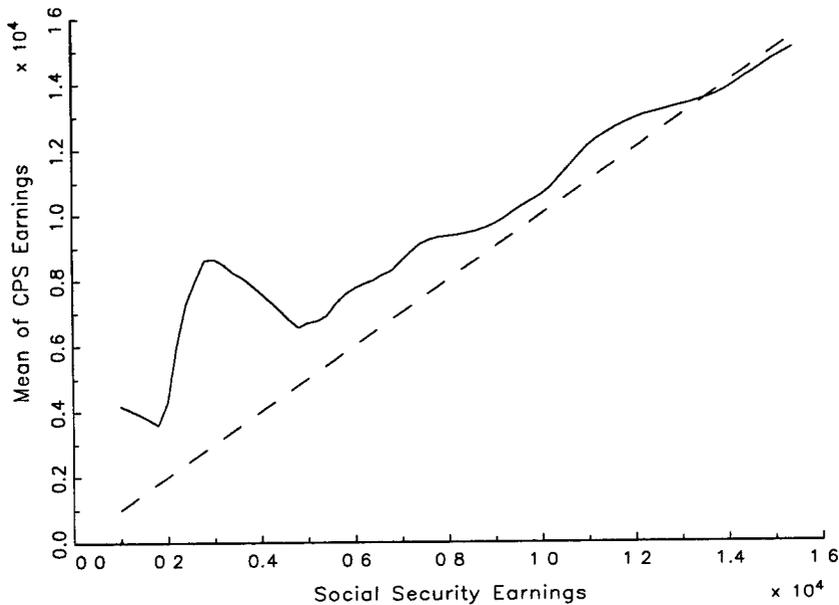
FIG. 3.—CPS earnings conditional on SSA earnings. Age = 39, Education = 12, and Weeks Worked = 50 for male head of household. Income in ten thousands of dollars. Key: — = conditional mean, and – – = 45° line.

section IIB will result in an upward bias for a dummy variable representing male gender. In the case where earnings is a regressor, the response of the dependent variable, such as savings, to earnings also has a larger downward bias for men than women. Comparison across gender must take these results into account.

It might be argued that the systematic bias above represents relationships to other regressors such as age, education, or the number of weeks worked. Bound and Krueger (1991) examined this by considering the regression of the response error on these variables. As noted earlier, they utilized maximum-likelihood estimation of a Tobit model to address this issue. Their results suggest that these variables are not related to the response error. Their approach focused on the relationship between response error and those variables, not controlling for true income. The question addressed here is whether these variables cause a different response structure controlling for true income. However, similar results are found.

Figures 3–6 present the results from the estimation of $E[CPS \mid SSE$, Age, Ed, Weeks Worked]. Figures 3 and 4 present the level curves generated by fixing values of Age, Education, and Weeks Worked near their mean values. For both men and women, it is evident that the relationship
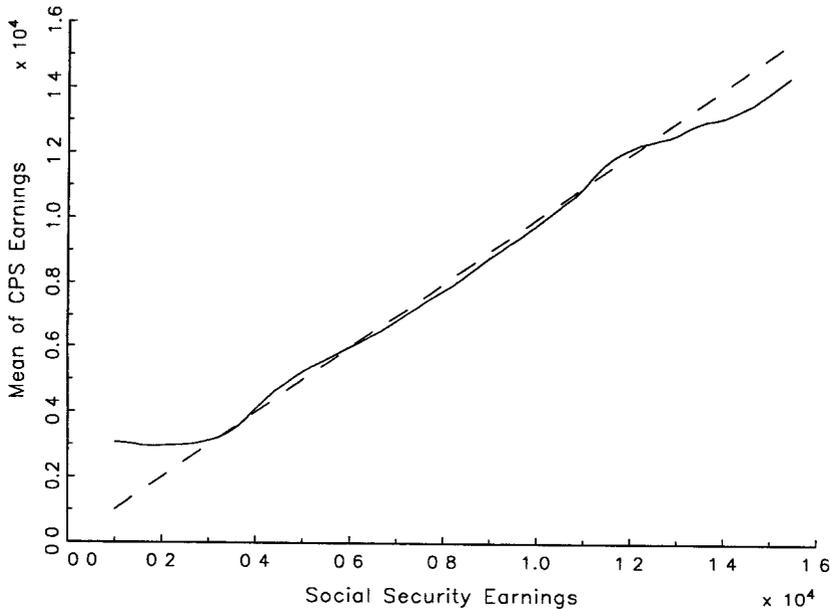
FIG. 4.—CPS earnings conditional on SSA earnings. Age = 39, Education = 12, and Weeks Worked = 50 for female head and spouse of head of household. Income in ten thousands of dollars. Key: — = conditional mean, and - - - = 45° line.

between SSA earnings and CPS earnings is not affected by these other variables. The overreporting for males at low income is strikingly confirmed; there is much less overreporting for women at low income. As with any estimation, the more parameters estimated with the data, the less precise the estimates. The reader should be aware that confidence bands for these estimates are much larger than those in figures 1 and 2.

Figures 5 and 6 present a set of level curves for Education. The Education variable has some interesting peaks for men near 14 years of education (in most cases a trade school or associates degree). There is seemingly large systematic bias here. However, the large confidence bands for these curves prevent any strong conclusions from being drawn. The peaks may simply be an artifact of these data. Similar level curves for the Age and Weeks Worked variables are available from the author. These curves show a very flat response to these variables.

## C. Estimation of the Conditional Quantiles

The median is an interesting measure of central tendency in this case. The median measures what proportion of people are overreporting or underreporting. Quantiles measure the symmetry and spread of the distribution.
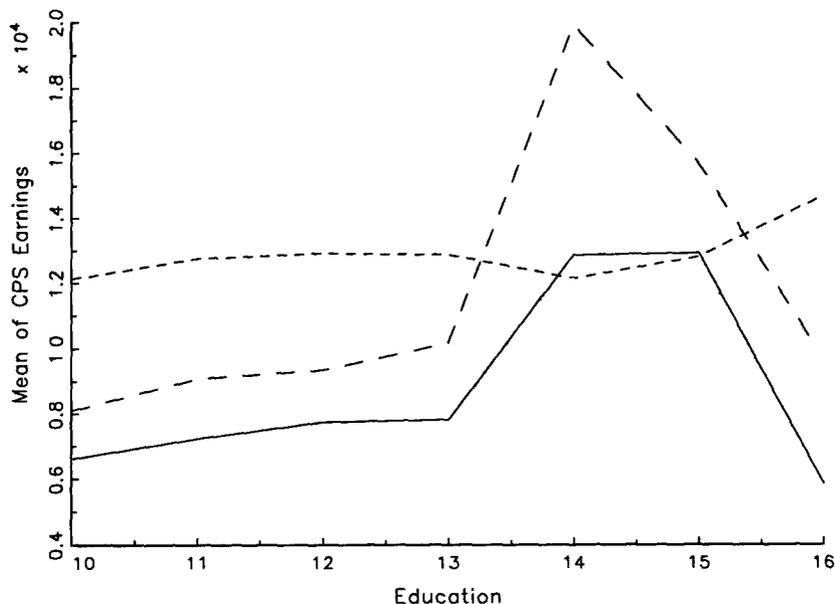
FIG. 5.—CPS earnings conditional on SSA earnings and education. Weeks Worked = 50 and Age = 39 for Male Head of Household. Income in ten thousands of dollars. Key: — = Social Security Earnings = $6,000; – – = Social Security Earnings = $8,000; · · · = Social Security Earnings = $10,000; - - - = Social Security Earnings = $12,000; and · · · = Social Security Earnings = $16,000.

Figures 7 and 8 display the median and the 10%, 25%, 75%, and 90% quantiles for the cross-sectional samples of the men and women. The median is very close to the 45° line. This indicates that overreporting of income is equally likely as underreporting for all income groups. Median regressions for earnings will be more robust to the measurement error problem than mean regressions. Hence, many questions that have been examined using OLS and other mean regression estimators are worth reexamining using median regression methodology.

There is substantial asymmetry in the distribution for men and less so for women. The asymmetry is most pronounced for men at low incomes and is positively skewed (toward overreporting). Most notably, the upper decile for men is very high for low SSA earnings. At higher earnings, the overreporting is less pronounced. Also at higher earnings, the lower decile demonstrates a tendency for large underreporting of earnings. These findings substantiate those of Bound and Krueger (1991) and the findings in the previous subsection, that a "regression to the mean" exists for men. However, from the quantile analysis, it is learned that this result is primarily generated by a subpopulation of individuals (about 10%), rather than a population-wide phenomenon.
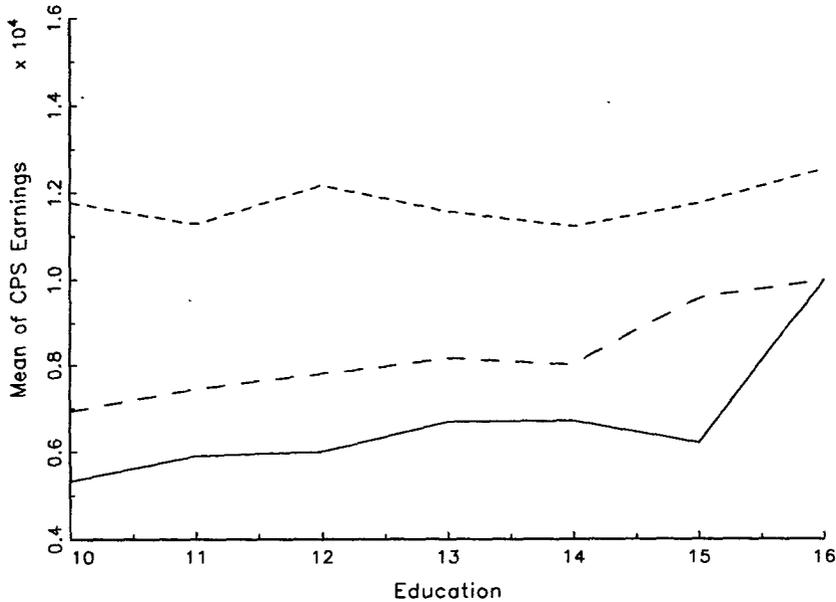
FIG. 6.—CPS earnings conditional on SSA earnings and education. Weeks Worked = 50 and Age = 39 for Female head and spouse of head of household. Income in ten thousands of dollars. Key: — = Social Security Earnings = $6,000; – – = Social Security Earnings = $8,000; · · · = Social Security Earnings = $10,000; - - - = Social Security Earnings = $12,000; and · · · = Social Security Earnings = $16,000.

The results for women are less dramatic but also suggest some regression to the mean phenomena. As with the men, overreporting is equally likely as underreporting. For women with high earnings, an underreporting phenomenon similar to that of the men is found.

The asymmetry and the overreporting of earnings will severely bias simulations of eligibility for government assistance programs. Such simulations will understate the true eligible population. The results of the quantile analysis as well as the conditional mean analysis quantify this problem. Researchers are often in a quandary as to how to handle sample observations where an individual is not income-eligible for a program but reports participation. Marquis and Moore (1990) and Bollinger and David (1997a) show that false positives are unlikely. Hence, it is most plausible that "ineligible participants" are actually eligible participants that are overreporting income.

## D. Estimates of Probability of Correct Report

In samples 3 and 4, 11.7% of the men and 12.7% of the women report their earnings correctly. In fact, 53.9% of the men and 56.2% of the
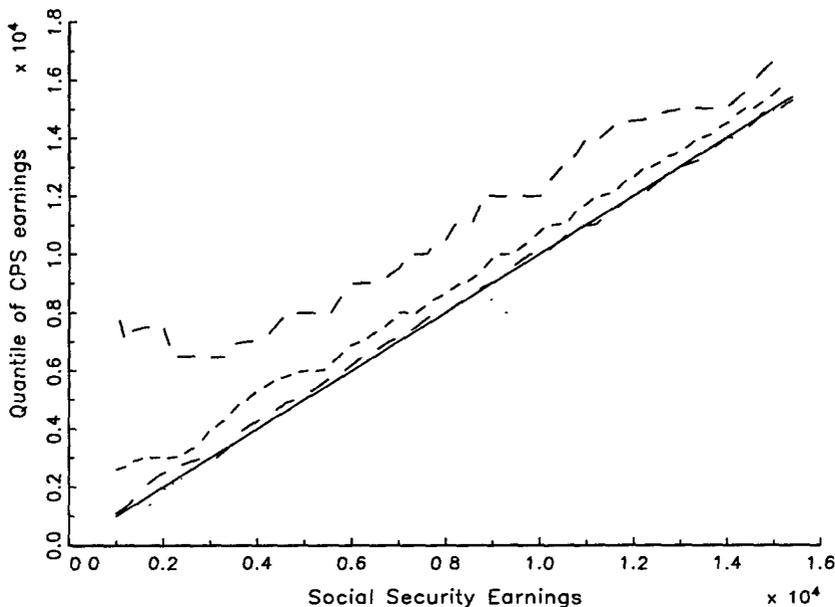
FIG. 7.—CPS earnings conditional on SSA earnings. Median and quantile regressions for male head of household. Income in ten thousands of dollars. Key: — = 45° line, – – = median, · · · = lower quartile, - - - = upper quartile, · · · = lower decile, and - · - = upper decile.

women report their earnings within 5% of the mean income (within $534 for the men and $338 for the women). Nonparametric estimation is applied to estimation of the conditional probability functions in a straightforward way.

The results of estimation of the conditional probability of correctly reporting income levels show no relationship between the probability of correctly reporting income and the level of true income. This observation is also supported when the definition of accurate reporting is broadened to include responses within 5% of the mean income. These results are available from the author by request.

As noted in section IIIC, the results indicating that men, on average, overreport earnings are driven by a small subsample of the population. There is little or no relationship between the probability of giving a correct response and the level of income. The finding that errors are negatively correlated with true income is true only for "gross errors," not for the probability of giving an accurate response.

## IV. Conclusions

Utilization of nonparametric methodology has allowed a deeper understanding of the structure of response error in earnings reports of the CPS.
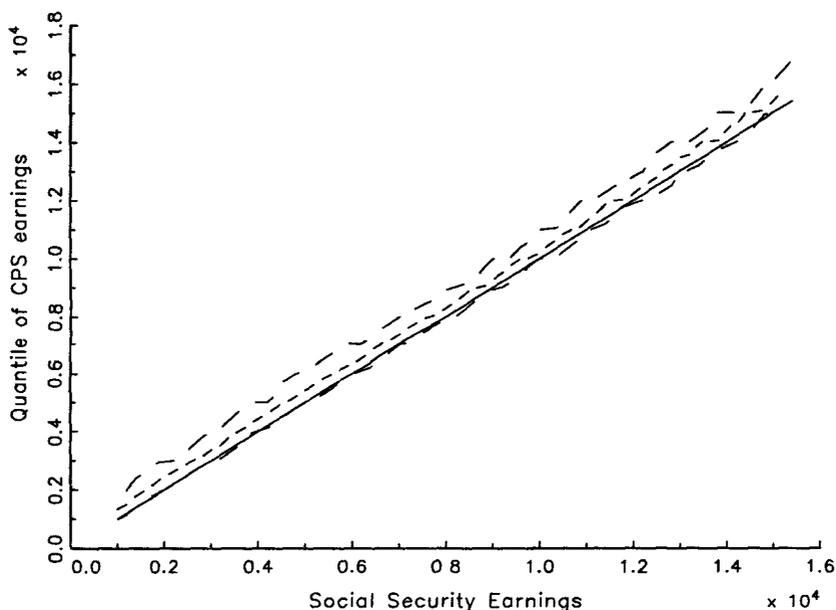
FIG. 8.—CPS earnings conditional on SSA earnings. Median and quartile regressions for female head and spouse of head of household. Income in ten thousands of dollars. Key: — = 45° line, – – = median, · · · = lower quartile, - - - = upper quartile, · · · = lower decile, and - · - = upper decile.

While in some cases, nonparametric results simply confirmed those found using OLS and correlation analysis, the finding that overreporting of income is concentrated in the lower end of the income distribution for men is new. Testing some of the linear specifications demonstrated that for men a nonlinear relationship between reported earnings and true earnings existed but for women the relationship was linear. Further, analysis of nonparametric median and quantile regressions presents a more complete picture of the underlying response distribution. It is clear that high overreporting of income for low-income men is driven by about 10% of the reporters who grossly overreport their income.

The results of Bound and Krueger (1991), in particular that response error is unrelated to variables such as age and education but is related to gender, are confirmed. Although estimation of coefficients on education and age (or constructed experience) are not likely to be biased by the measurement error, comparisons between genders are clearly biased. Response error in income cannot be treated as additive white noise because of its relationship with gender and true earnings.

The first major finding here demonstrates that the additional consistency checks of constructing a panel sample from the CPS result in lower

overall measurement error in income. This may indicate that inability to follow individuals is related to response error. Hence, differences in regression results between cross-sectional and panel samples may be due to differences in response error.

By examining median and quantile regression, two additional conclusions can be made. First, since the median of the response error is not related to income or gender, median wage regressions will be less affected by response error. Second, severe bias is likely to result in construction of samples of individuals who are eligible for means-tested transfer programs. Since the response error is asymmetric, with the largest errors occurring for low-income males, a sample constructed using an income threshold will contain too few males at the very lowest income levels. This may significantly affect predictions.

## References

Aigner, Dennis J.; Hsiao, Cheng; Kapteyn, Arie; and Wansbeek, Tom. "Latent Variable Models in Econometrics." In *Handbook of Econometrics*, vol. 2, edited by Z. Griliches and M. D. Intriligator, pp. 1323–93. New York: Elsevier, 1984.

Bierens, Herman J. "Kernel Estimators of Regression Functions." In *Advances in Econometrics—Fifth World Congress*, vol. 1, edited by Truman F. Bewley, pp. 99–144. New York: Cambridge University Press, 1987.

Bollinger, Christopher R., and David, Martin H. "Modeling Food Stamp Participation in the Presence of Reporting Errors." *Journal of the American Statistical Association 92* (September 1997): 827–35. (*a*)

———. "Using Attrition as a Proxy for Response Error." Working paper. Atlanta: Georgia State University, 1997. (*b*)

Bound, John; Brown, Charles; Duncan, Greg J.; and Rogers, Willard L. "Measurement Error in Cross-Sectional and Longitudinal Labor Market Surveys: Results from Two Validation Studies." *Panel Data and Labor Market Studies*, edited by J. Hartog, G. Ridder, and J. Theeuwes, pp. 1–19. New York: Elsevier, 1990.

Bound, John, and Krueger, Alan B. "The Extent of Measurement Error in Longitudinal Earnings Data: Do Two Wrongs Make a Right?" *Journal of Labor Economics 9* (January 1991): 1–24.

Card, David. "The Effect of Unions on the Distribution of Wages: Redistribution or Relabelling?" Working Paper no. 187. Princeton University, Princeton, NJ: Industrial Relation Section, 1991.

Duncan, Greg J., and Hill, Daniel H. "An Investigation of the Extent and Consequences of Measurement Error in Labor-Economic Survey Data." *Journal of Labor Economics 3* (October 1985): 508–32.

Freeman, Richard B. "Longitudinal Analysis of the Effects of Trade Unions." *Journal of Labor Economics 2* (January 1984): 1–26.

Fuller, Wayne A. *Measurement Error Models.* New York: Wiley, 1987.

Goldberger, Arthur S. "Abnormal Selection Bias." *Studies in Economet-*

*rics, Time Series, and Multivariate Statistics,* edited by Samuel Karlin, Takeshi Amemiya, and Leo A. Goodman, pp. 67–84. New York: Academic Press, 1983.

Greenberg, David, and Halsey, Harlan. "Systematic Misreporting and Effects of Income Maintenance Experiments on Work Effort: Evidence from the Seattle-Denver Experiment." *Journal of Labor Economics 1* (October 1983): 380–407.

Härdle, W. *Applied Nonparametric Regression.* Cambridge, MA: Cambridge University Press, 1990.

Lee, Byung-Joo. "Asymptotic Distribution of the Ullah-Type Specification Test against the Nonparametric Alternative." *Journal of Quantitative Economics* 10 (January 1994): 73–92.

Marquis, Kent, and Moore, Jeffrey. "Measurement Errors in SIPP Program Reports." SIPP Working Paper no. 9008. Washington, DC: U.S. Bureau of the Census/Center for Survey Methods Research, 1990.

Mathiowetz, Nancy A., and Duncan, Greg J. "Out of Work. Out of Mind: Response Error in Retrospective Reports of Unemployment." *Journal of Business and Economic Statistics 6* (April 1988): 221–29.

Mellow, Wesley, and Sider, Hal. "Accuracy of Response in Labor Market Surveys: Evidence and Implications." *Journal of Labor Economics 1* (October 1983): 331–44.

Poterba, James M., and Summers, Lawrence H. "Reporting Errors and Labor Market Dynamics." *Econometrica 6* (November 1986): 221–29.

Rogers, Willard L., and Herzog, A. Regula. "Covariances of Measurement Errors in Survey Responses." *Journal of Official Statistics 3* (October 1987): 403–18.

U.S. Department of Commerce, Bureau of the Census. *Current Population Survey, March 1977: Annual Demographic File* [Machine readable file and technical documentation]. Washington, DC: U.S. Department of the Census (producer). Ann Arbor, MI: Inter-university Consortium for Political and Social Research (distributor).

———. *Current Population Survey, March 1978: Annual Demographic File and Social Security Administration Exact Match File.* [Machine readable file and technical documentation]. Washington, DC: U.S. Department of the Census (producer). Ann Arbor, MI: Inter-university Consortium for Political and Social Research (distributor).

———. *Current Population Survey, May 1985: Work Schedules, Multiple Job Holding, and Premium Pay.* [Machine readable file and technical documentation]. Washington, DC: U.S. Department of the Census (producer). Ann Arbor, MI: Inter-university Consortium for Political and Social Research (distributor).

———. *Current Population Survey, May 1991: Multiple Job Holding and Work Schedules.* [Machine readable file and technical documentation]. Washington, DC: U.S. Department of the Census (producer). Ann Arbor, MI: Inter-university Consortium for Political and Social Research (distributor).