

# Measurement Error and Misclassification: A Comparison of Survey and Administrative Data

Arie Kapteyn, *RAND*

Jelmer Y. Ypma, *RAND and University of Groningen*

We provide both a theoretical and empirical analysis of the relation between administrative and survey data. By distinguishing between different sources of deviations between survey and administrative data we are able to reproduce several stylized facts. We illustrate the implications of different error sources for estimation in (simple) econometric models and find potentially very substantial biases. This article shows the sensitivity of some findings in the literature for the assumption that administrative data represent the truth. In particular, the common finding of substantial mean reversion in survey data largely goes away once we allow for a richer error structure.

## I. Introduction

Microdata are an essential ingredient of research in economics and other social sciences. Such data, where information is available for each micro-unit (individual, firm, etc.) separately, are usually obtained through a survey or from administrative records. Both methods of data collection have their advantages and disadvantages.

This article is part of an NIA-funded project (Comparison of Survey and Register Data: The Swedish Case, R03AG21780) in collaboration with Anders Klevmarken (Uppsala University) and Susann Rohwedder (RAND). We thank Susann Rohwedder, Anders Klevmarken, Fredrik Johansson, and two anonymous referees for their helpful comments. Contact the corresponding author, Arie Kapteyn, at [kapteyn@rand.org](mailto:kapteyn@rand.org).

[*Journal of Labor Economics*, 2007, vol. 25, no. 3]  
© 2007 by The University of Chicago. All rights reserved.  
0734-306X/2007/2503-0004\$10.00

One problem in surveys is nonresponse, both to single questions (item nonresponse) and to the entire questionnaire (unit nonresponse). Another problem is measurement error. Furthermore, surveys are generally costly. In administrative files, such as the U.S. Social Security Administration (SSA), information is typically available for large numbers of individuals. These data are generally assumed to be reliable but may not measure exactly the concept a researcher is interested in.

There is an additional problem with administrative data that may have severe consequences for the models that one would want to estimate based on these data: administrative databases typically need to link data from different sources, introducing the possibility of mismatching, due to imperfect linkage information (e.g., errors in social security numbers [SSNs] that would be used for linking records for a given individual).

In this article we concentrate on sources of measurement error in survey and administrative data by comparing individual survey information and administrative information on the same variables. By doing so, we both replicate a number of earlier studies using new data and propose extensions of models in the literature that may help us better understand the nature of measurement error in both survey data and administrative data.

The number of data sets allowing for validation studies of survey information appears to be quite limited. We discuss three of them here.<sup>1</sup>

A first example of a study comparing survey data and data from administrative sources is the Panel Study of Income Dynamics Validation Study (PSIDVS). In 1983 and 1987 a questionnaire based on the PSID questionnaire, but shorter, was administered to employees of a manufacturing company in the Detroit area to measure their earnings in the preceding years. At the same time, payroll records of the employees were collected from this firm. The data are assumed to be very accurate, since the firm was highly cooperative (Pischke 1995).

Duncan and Hill (1985) use the PSIDVS from 1983, which includes questions about annual earnings in the 2 preceding years (1981 and 1982). They find that means of log earnings in the survey data and the validation data do not differ significantly. However, at the individual level there may still be substantial differences. Assuming the administrative data to contain the “true” values, measurement error can be defined as the difference between survey data and administrative data. They find a reliability ratio between .64 and .84, depending on which year they look at and whether outliers are removed or not. As expected, the reliability ratio is lower for the question with a longer recall period, that is, earnings in 1981 versus earnings in 1982. Pischke (1995) finds, when using data from 1982 and 1986, that administrative data and survey data do not differ much in either

<sup>1</sup> Bound, Brown, and Mathiowetz (2001) provide a more extensive overview of validation studies dealing with earnings measures.

mean or variance. However, measurement error is found to be weakly negatively correlated with “true” log earnings. In cross-section data of 1982 this correlation is stronger than in cross-section data of 1986. When restricting data to respondents who are present in both years, no significant correlation is found between true log earnings and measurement error. The negative correlation is not significant when leaving out hourly workers, that is, when looking only at salaried employees. Rodgers, Brown, and Duncan (1993) and Bound et al. (1994) also report a negative correlation between measurement error and the true value. In contrast to Pischke (1995), they do not distinguish between hourly workers and salaried employees.

Another data set suitable for comparison of survey and administrative data is a match constructed between the Current Population Survey (CPS) and data from the SSA for the years 1976 and 1977 (see, e.g., Bound and Krueger 1991). Once again, the maintained hypothesis for most of their paper is that the administrative data are error free. Using cross-section data, the reliability ratio for log earnings for men is .844 in 1976 and .819 in 1977. For women this ratio is higher, .939 and .924. For men they find large negative correlations between measurement error and true log earnings,  $-.46$  in 1976 and  $-.42$  in 1977. This correlation is small for women. Bollinger (1998) finds that the negative relationship between measurement error and earnings is mainly driven by overreporting among low earners.

When comparing their results obtained with the PSIDVS with the CPS-SSA data, Bound et al. (1994) find that, qualitatively, the results are similar. However, they notice a large difference in the standard deviation of the measurement error (.13 in the PSIDVS data vs. .32 in the CPS-SSA data). The difference seems to lie in the tails of the measurement error distribution, with very large outliers in the CPS-SSA data. Bound et al. (1994) suggest that these very large measurement errors are not necessarily due to misreporting in the survey but rather to errors in the SSA data.

A number of studies have addressed the possibility of measurement error in the earnings data collected in the Survey of Income and Program Participation (SIPP). We briefly discuss two studies relevant to our analysis. Pedace and Bates (2001) use a match between the 1992 SIPP longitudinal file and the Social Security Summary Earnings Records. These authors also assume the administrative data to represent the truth. It appears that respondents with low SSA earnings tend to overreport their earnings, whereas respondents with high earnings underreport. In contrast to the studies mentioned so far, Stinson (2002) does not make the assumption that the administrative data represent the truth. She instead estimates an earnings function allowing for (mutually uncorrelated) measurement error in both the survey and administrative data. She finds that both measures have similar magnitudes of measurement error, with the

error in the administrative data being slightly larger than in the survey data.

The studies cited here reach somewhat contrasting conclusions regarding whether there is negative correlation between the true value and measurement error or not. While this correlation seems evident in the CPS-SSA data and the SIPP data, in the PSIDVS data no such correlation is found for the salaried workers, but it is found when including hourly workers. This “mean reversion” in the measurement error has gradually attained the status of a stylized fact. Kim and Solon (2005) discuss its implications for the modeling of earnings dynamics.

In the remainder of the article we will provide both a theoretical and empirical analysis of the relation between administrative and survey data. We distinguish a number of different sources of measurement error. Measurement error in the administrative data will be due only to mismatching; that is, with a certain probability a value recorded in an administrative file refers to a different observation. Measurement error in survey data can be (1) absent, (2) classical but potentially mean reverting, or (3) the result of contamination. By distinguishing between different sources of deviations between survey and administrative data we will be able to reproduce several stylized facts in the literature. In doing so, we deviate from the almost universal assumption that the administrative data represent the truth.<sup>2</sup> We will illustrate the implications of these error sources for estimation in (simple) econometric models. The analysis is applied to Swedish data that have been collected for a validation study as part of a larger European health and retirement study (SHARE: Survey of Health, Ageing, and Retirement in Europe). Thus this article makes two contributions: (1) it adds to the limited number of empirical validation studies of earnings measurement in surveys, and (2) it shows the sensitivity of some findings in the literature for the assumption that administrative data represent the truth.

In Section II we describe the data. In Section III a relatively straightforward model of different error sources is proposed and their empirical implications are explored, both for the observed relation between survey and administrative data and for some econometric models incorporating survey and/or administrative data. Among other things we address the question when it is preferable to use administrative data and when survey data are to be preferred. Section IV estimates the model of Section III using our Swedish data set. We are able to identify various sources of error in the data and actually flag observations that suffer from different types of error. Section V concludes.

<sup>2</sup> As in most of the papers discussed above. The same assumption is made in more formal analyses of the use of validation samples, including Lee and Sepanski (1995) and Chen, Hong, and Tamer (2005).

## II. Data and Project Description

The Scandinavian countries have 30 years of experience with using administrative data for statistical purposes. The statistical offices in Denmark, Norway, Sweden, and Finland have made important progress in making information from various administrative files compatible and to ensure the link among various sources. The wealth of information contained in these administrative files is considerable. In our empirical work we will use Swedish data.

Every Swede has a unique SSN, which is also available in every administrative record. In principle this allows interviewers to ask respondents for their SSN and permission to link them to the information available in the administrative files. If the respondent agrees this may substantially shorten interviews, because questions, for instance about income, can be skipped. Interviews can then focus on information not available in administrative files, while combining the interview information with the (presumably) more reliable administrative information.

The experiment generating our data is motivated exactly by this consideration. As part of SHARE, a pan-European data collection effort among individuals 50 and over, an experiment was devised in Sweden to assess the usefulness of combining survey and administrative data. The purpose of the experiment was to inform SHARE about the possibilities of combining such data and potentially implement it in other countries than Sweden. Aspects to be investigated include (1) selectivity in survey responses in a number of important domains; (2) reliability of selected survey measures including income and pension entitlements by comparison with administrative data; (3) estimating biases in a number of important empirical relationships (e.g., health and socioeconomic status) when using survey data rather than administrative data; (4) generalizing from these findings to the limitations of international comparisons if administrative data are available in some countries, but not in others. In this article we mainly concentrate on aspect 2 and pay some attention to aspect 3.

### Administrative Data

In our empirical analysis we will use a sample of individuals over age 50 from LINDA (Longitudinal Individual Data for Sweden). LINDA is a registry-based longitudinal database that is representative of the Swedish population since 1960. LINDA has two subsamples. The first subsample is the population sample, representative of the entire population, with a coverage rate of 3.35%. The second subsample is the immigrant sample, covering 19.5% of the immigrant population. The samples are kept both cross-sectionally and longitudinally representative of their respective populations. There are two principal data sources: (1) the Income Registries, available annually since 1968 and (2) the Population Censuses, available

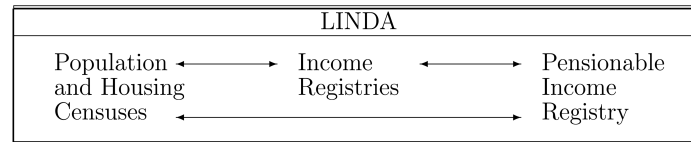


FIG. 1.—Some of the registries used in LINDA, linked via SSNs

every fifth year from 1960 to 1990 (no census has been taken after 1990). Other administrative files have been added during the nineties (see Edin and Frederiksson [2000] for more details).

By drawing from LINDA we have a wealth of information available about our respondents from the administrative data.<sup>3</sup> The information in the different administrative files is linked by the SSN of individuals (see fig. 1). Since the information comes from several administrative files an incorrect SSN in an administrative record may lead to a wrong link.

As Abowd and Vilhuber (2005) noticed while looking at unemployment spells, errors in linking can be a real problem. Biases arise that are the result of errors in period-to-period linking of records. When using a different probabilistic matching algorithm, Abowd and Vilhuber achieve smaller error rates than the overall error rate of 7.8%, the Bureau of Labor Statistics found in an SSN validation project. Errors in database-linking can be caused by similar errors in period-to-period linking. These would be most likely the result of recording a wrong or mistyped SSN.

We will be using mostly demographic variables (e.g., education, age, and gender) and financial variables (earnings, pensions, and taxes). Since the survey contained mostly questions about 2002, we will be using data from that year.

The LINDA data are generally thought to be a reasonably accurate measure of earnings of the population they cover. At the same time, as with the American data discussed above, the administrative files are based on the linking of data from different sources, and, hence, one would suspect that a certain number of errors will occur. How serious this problem is will be a focus of our empirical analysis. We believe that mismatching is likely to be the most important source of error. Conditional

<sup>3</sup> Among others, LINDA includes annual cash earnings, annual taxable benefits, social security sickness compensation (only if not old-age pension), if single family home (owner occupied), if condominium, if secondary house, tax assessed value on house, market value of house, stocks and shares, bank holdings, bonds, mutual funds, mortgages, other loans, schooling by number of years and category, pensions (social security, old age and disability pension, group [occupational] pensions, private pension insurance [annuity]), capital incomes (interest and dividends received, realized capital gains, interest paid), and total tax (income tax on earnings, capital income tax, real estate tax, wealth tax).

**Table 1**  
**Evolution of the Sample**

	Earnings	Pensions	Taxes
Total number of persons in sample	1,431	1,431	1,431
Number of respondents	881	881	881
Number of respondents with positive administrative values	511	492	845
Number of respondents with positive survey values	414	376	495
Both values positive	400	369	487

on correctly identifying an observation, errors in financial variables are less likely. For instance, a variable representing pension benefits will be verified by both the beneficiary and the payee, since both have an interest in avoiding errors. The same is true of other financial variables like earnings. There is no such mechanism that would verify the correct linking of records in the construction of administrative analysis files.

#### Survey Data

In the beginning of every year, just before tax returns need to be filed, people in Sweden receive preprinted tax statements from the tax authorities. Around this time in 2003 a survey was conducted of 1,431 individuals age 50 or over. Out of the 1,431 individuals who were contacted, 881 responded—469 of them women and 412 men. The timing of the survey was chosen so as to optimize the information available to respondents when answering the questions in the survey. The questionnaire contains several questions about household income and expenses, partner income, and assets. Besides the financial questions there are some smaller sections about household composition, health, retirement, and education of the respondent.

#### Descriptive Comparisons of Administrative and Survey Data

Our analysis will concentrate on three monetary variables: earnings, pensions, and taxes. Table 1 shows the evolution of the sample if we move from the gross sample of 1,431 respondents drawn from the administrative files to the sample of respondents who answered at least some questions over the phone. The survey has 881 respondents, a response rate of 61.6%. In view of the age distribution of the sample, it is not surprising that many respondents report zero earnings. The answer to the survey question about taxes could be given either as an amount or as a percentage of income. We do not consider the percentage answers, as the calculated amount of taxes paid would most likely exhibit an error structure different from the other responses. Roughly speaking, the implied error is a result of taking the ratio of two variables (cf. Duncan and Hill 1985). Thus the

**Table 2**  
**Comparison of Population and Sample Characteristics**

Administrative Variable	Original Sample* (1)		Respondents Sample† (2)		Earnings Sample (3)		Pension Sample (4)		Taxes Sample (5)	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD
Female (%)	53.0		53.2		52.5		53.7		54.4	
Age	63.6	9.6	63.5	9.5	57.1	5.3	70.8	7.2	64.3	9.4
Education (%):										
Low	28.1		26.1		20.0		28.2		24.8	
Middle	35.8		37.8		48.3		27.4		39.4	
High	17.5		19.1		31.3		11.9		18.1	
Missing	18.6		17.0		.5		32.5		17.7	
Earnings‡	114	153	122	153	...		...		...	
Earnings > 0	214	149	211	147	246	138	...		...	
Pensions‡	74	129	78	154	...		...		...	
Pensions > 0	128	147	140	184	...		141	74	...	
Taxes‡	59	85	64	100	...		...		...	
Taxes > 0	64	87	67	101	...		...		67	121
N	1,431		881		400		369		487	

NOTE.—Earnings, pension, and taxes are given in 1,000 SEK per year.

\* In the original sample the number of observations with a positive amount is 765 for earnings, 828 for pensions, and 1,328 for taxes.

† In the respondents sample the number of observations with a positive amount is 511 for earnings, 492 for pensions, and 845 for taxes.

‡ 1,430 observations are used, since administrative values are missing for one individual.

number of respondents in the survey with positive survey taxes is much lower than in the administrative data.

Table 2 compares a number of administrative variables across different subsamples. Comparing the sample of respondents (col. 2) to the original (gross) sample (col. 1) shows that the age and gender composition is essentially the same. The respondent sample is slightly better educated. Perhaps related to that, the respondent sample exhibits somewhat higher (administrative) earnings (7%). If we only consider observations with positive earnings the difference disappears: respondents earn 1.5% less on average than the overall mean of the gross sample. Pensions among the respondents are 5% higher on average in the respondent sample than in the gross sample, and 9% higher if we only consider positive pensions. For taxes the difference is 8.5% and 6% if we consider only respondents with positive (administrative) tax payments.

Our empirical analysis will be primarily concerned with a comparison of nonzero observations in both the administrative and survey data. Columns 3–5 present administrative values for the subsamples that have positive values for both the administrative and survey variables. Now age and education vary considerably across subsamples for obvious reasons. Respondents with earnings are typically younger and respondents with pensions are typically older. Education levels are also different, reflecting



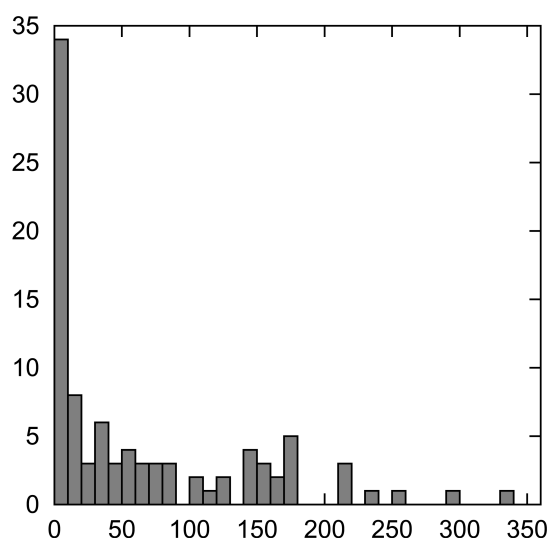


FIG. 2.—Histogram of administrative earnings (in 1,000 SEK/year) for respondents reporting no earnings.

cohort differences in educational achievement. Comparing column 2 with columns 4 and 5 shows essentially no difference in means for pensions and taxes, respectively. A comparison of columns 2 and 3 shows that mean earnings are substantially higher in the latter column than in the former. It is worthwhile therefore to investigate the differences between administrative data and survey data in some more detail.

The survey contained a number of questions on income and income-related variables. A number of measures were taken to improve survey quality and to improve immediate comparison to the administrative variables. According to Hurd et al. (2004), one way to increase the quality of report is by giving respondents the opportunity to report income in a time span consistent with how they receive their income. For example, instead of forcing respondents to provide a yearly amount for earnings and pensions, they were given the possibility to report these amounts either per month or per year.

The question “Did you have any income from employment in 2002?” was answered affirmatively by 435 respondents. When comparing these answers with the administrative files, some discrepancies are found (see table 3). Seventeen respondents claimed to have earnings, when according to the administrative files they have zero earnings. Ninety-three respondents reported to have no earnings, when they should have earnings according to the administrative files. As can be seen from figure 2, two groups can be distinguished. One large group has low earnings (e.g., 34

**Table 3**  
**Correspondence between Administrative and Survey “Answers” on**  
**Relevant Questions**

	Yes		No		DK		RF	
Did you have any income from employment in 2002?								
Adm. > 0	418	81.8%	93	18.2%	0	0%	0	0%
Adm. = 0	17	4.6%	352	95.1%	0	0%	1	.3%
Did you receive any type of old-age pension in 2002, such as . . . ?								
Adm. > 0	430	87.4%	62	12.6%	0	0%	0	0%
Adm. = 0	7	1.8%	381	97.9%	1	.3%	0	0%
	> 0		= 0		DK		RF	
How much did you earn per month in 2002, before taxes?								
Adm. > 0	327		1		6		5	
Adm. = 0	7		1		1		0	
Altogether, about how much did you earn from your main job in 2002, before taxes?								
Adm. > 0	73		0		6		0	
Adm. = 0	7		0		1		0	
How much did you receive in pension payments per month in 2002, before taxes?								
Adm. > 0	353		2		43		6	
Adm. = 0	6		0		0		0	
Altogether, about how much did you receive in old-age pension payments (before taxes) in 2002?								
Adm. > 0	16		1		8		0	
Adm. = 0	1		0		0		0	
Think about the total income you received in 2002 from employment, pensions and taxable benefits. About how much did you pay (will you pay) in income tax on that amount?								
Amount								
Adm. > 0	487		1		12		0	
Adm. = 0	8		8		0		0	
Percentage								
Adm. > 0	230		0		10		0	
Adm. = 0	3		0		0		0	

respondents earn less than 10,000 SEK per year), which are easy to forget or perhaps thought not to be worth mentioning.<sup>4</sup> Since this group is included in the mean earnings in column 2 of table 2, but not in column 3, the mean in column 3 is substantially higher than in column 2. For the other group, respondents with substantial amounts of earnings, the explanation of rounding to zero is less plausible. Errors can be made by both the respondent and the interviewer (e.g., an interviewer may simply type in the wrong code). A different explanation could be mismatching, where a nonearning respondent is incorrectly matched with an individual in the administrative files who has positive earnings.

It is interesting to compare these findings to those of Pedace and Bates (2001) in their analysis of the SIPP. They find that, of those who had earnings according to the administrative data, 5.5% of those surveyed said that they did not have any earnings in the survey. In our survey that percentage is 18.2. However, Pedace and Bates (2001) find that, of those who had no earnings according to the administrative files, 18% reported earnings in the SIPP. In our data that percentage is 4.6. Thus, in percentage terms, the discrepancy between administrative data and survey data appears similar across the two data sets, but in the SIPP zero administrative earnings are reported to be positive in the survey data at about the same rate as positive administrative earnings in Sweden are reported to be zero in the survey. Not too much should be read into this, if only because of the different age compositions of the Swedish and U.S. samples.

For the pension data we find a similar pattern of false responses. The percentage of respondents with positive administrative pensions but stating they do not have any pensions is 12.6. A smaller percentage of respondents (1.8%) report having a pension, while this pension is not found in the administrative data. Four hundred ten respondents gave an answer to the monthly pension question, including 2 zeros, 6 refusals, and 43 don't knows. Only 26 respondents chose to give a yearly amount. For 369 respondents we have both a positive administrative value and a positive survey value.

From now on we concentrate on positive values of the monetary variables. We will not pay attention to sources of selectivity, such as refusals, zero responses, and don't knows. See Johansson and Klevmarken (2006) for a detailed analysis of various sources of nonresponse in the data. We will generally use logarithms of the monetary variables we are considering to achieve distributions that are approximately symmetric.

In table 4 some statistics for the monetary variables of interest are given. The difference between log earnings measured in the survey and the administrative data is on the order of .02 on average. Strikingly, the variance

<sup>4</sup> SEK represents the Swedish kronor. In 2002, 10 SEK equal approximately one U.S. dollar.

**Table 4**  
**Summary of Administrative and Survey Variables**

	<i>N</i>	Survey	Administrative	$s - r$	Reliability*	$\text{corr}(m, r)$
Log earnings	400	12.196 (.821)	12.172 (.961)	.0243 (.656)	.6821	-.5395
Log pension	369	11.649 (.945)	11.694 (.658)	-.0451 (.878)	.3595	-.2699
Log taxes	487	10.869 (.917)	10.786 (.829)	.0828 (.577)	.6734	-.1875

NOTE.—SD is given in parentheses.

\*  $\sigma_s^2/(\sigma_s^2 + \sigma_m^2)$ , where  $r$  is the administrative variable and  $m$  is the difference between survey and administrative variable.

of the survey data is smaller than the variance of the administrative data. Under the assumption of classical measurement error for the survey data and no error in the administrative data, this would not be possible. Looking at log-pensions, we again observe modest, though somewhat larger, differences between the log-means of the administrative values and the survey values. Just as in the earnings data, we see a substantial negative correlation between the difference between survey and administrative data,  $s - r$ , and the administrative data,  $r$ . Finally, mean log-taxes are about .08 higher in the survey data than in the administrative data. The correlation between  $s - r$  and  $r$  is still negative but closer to zero than for earnings and pensions.

Ten people (2.5%) gave an earnings amount exactly equal to the administrative value, while an additional 49 (12.3%) respondents provided amounts within 1,000 SEK of the administrative value. For the pension and taxes data we find, respectively, 40 (13.3%) and 25 (5.1%) answers exactly equal to the administrative value, with an additional 59 (16.0%) and 72 (14.8%) answers within the 1,000 SEK bound.

Figure 3 presents histograms of the difference between survey data and administrative data for earnings, pensions, and taxes. Most of the values are close to zero, but we also see a large positive difference for earnings and large negative differences for pensions and taxes. The large positive and negative values force the scales in figure 3 to be rather coarse. Hence, in figure 4 the same histograms are given but now truncated at a value of  $\pm 6$ . Although the mean difference between survey and administrative earnings is positive (table 4), we see from the truncated histogram that for most observations administrative earnings are larger than survey earnings. If the administrative data are assumed correct, this would indicate that most people underestimate their earnings. The difference in pensions data seems to center around zero, whereas the difference in taxes data is mostly positive.

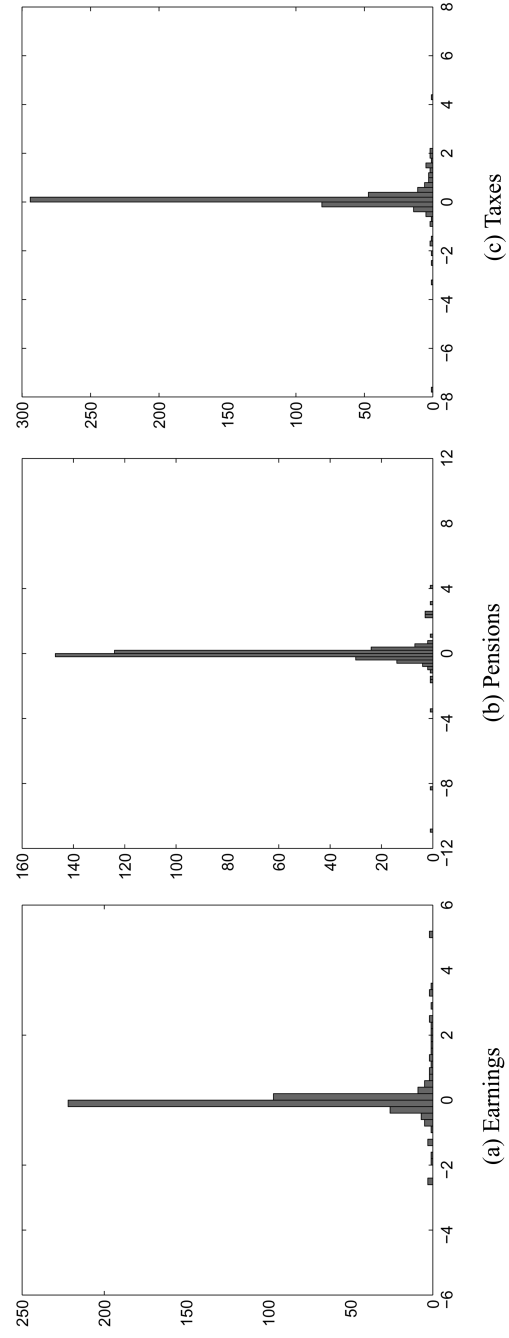


FIG. 3.—Histograms of  $m_i = s_i - r_i$  for earnings, pensions, and taxes

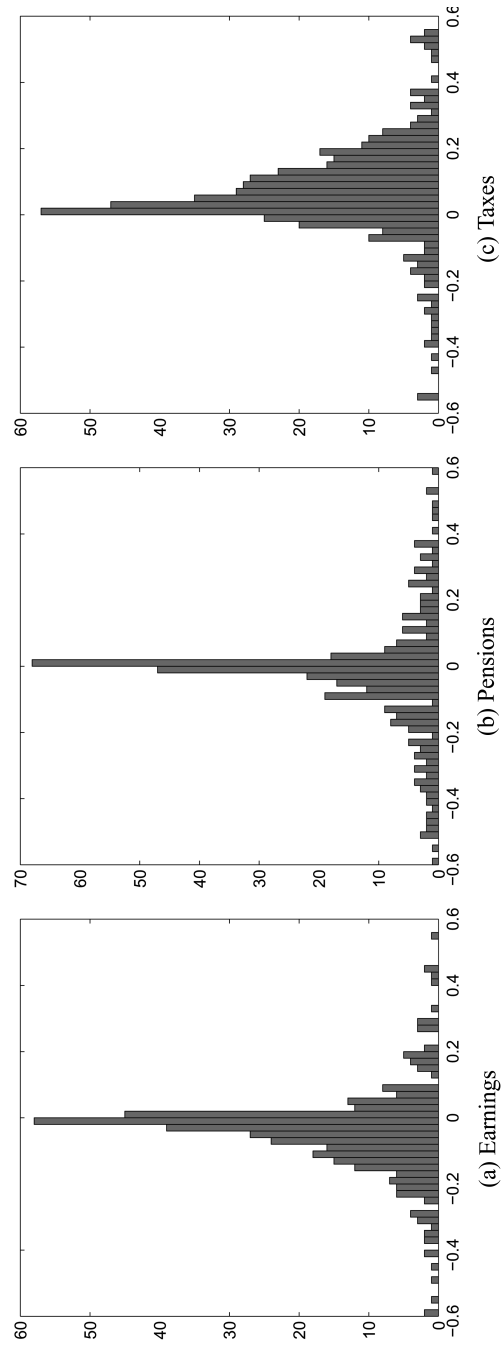


FIG. 4.—Histograms of  $m_i = s_i - r_i$  for earnings, pensions, and taxes, truncated at  $\pm .6$

### III. Modeling Different Errors

Let  $\xi_i$  be the logarithm of the true value of the variable of interest (e.g., log-income) for individual  $i$ . This true value is not measured directly, but instead two sources of data, capturing the same concept, are available. In contrast to the assumptions typically made in the literature, we will assume that both sources of information, administrative and survey data, may contain error. The structure of the error, however, is different per source.

Define four independently and identically distributed and mutually independent normal variables:  $\xi_i \sim N(\mu_\xi, \sigma_\xi^2)$ ,  $\zeta_i \sim N(\mu_\zeta, \sigma_\zeta^2)$ ,  $\eta_i \sim N(\mu_\eta, \sigma_\eta^2)$ , and  $\omega_i \sim N(\mu_\omega, \sigma_\omega^2)$ , where  $i$  indexes the unit of observation. Denote the value elicited in the survey by  $s_i$  and the corresponding administrative value by  $r_i$ . For the derivations in this section the normality assumptions we are making are not necessary, but will be exploited later in maximum likelihood estimation.

Errors in the administrative data are assumed to be due only to mismatching. With probability  $\pi_r$  the observed administrative value,  $r_i$ , is equal to the true value of individual  $i$ ,  $\xi_i$ . In the case of a mismatch, which occurs with probability  $1 - \pi_r$ , the administrative value  $r_i$  corresponds to the true value of someone else.<sup>5</sup> This mismatched value will be denoted by  $\zeta_i$ , where no correlation exists between  $\xi_i$  and  $\zeta_i$ . Note that our sample comes from a subset of the population, containing only individuals of age 50 and older, whereas a mismatch may come from the complete sample available in LINDA. The distributions of  $\xi$  and  $\zeta$  can therefore be different. The observed values  $r_i$  are a mixture of correct matches and mismatches:

$$r_i = \begin{cases} \xi_i & \text{with probability } \pi_r \\ \zeta_i & \text{with probability } (1 - \pi_r) \end{cases}. \quad (1)$$

For the survey data we distinguish three cases. The observed survey value is correct with probability  $\pi_s$ . With probability  $1 - \pi_s$  the survey data contain response error, part of which is mean-reverting. Some of these observations, a proportion  $\pi_\omega$ , are contaminated, modeled by adding an extra error term,  $\omega_i$ :

$$s_i = \begin{cases} \xi_i & \text{with probability } \pi_s \\ \xi_i + \rho(\xi_i - \mu_\xi) + \eta_i & \text{with probability } (1 - \pi_s)(1 - \pi_\omega) \\ \xi_i + \rho(\xi_i - \mu_\xi) + \eta_i + \omega_i & \text{with probability } (1 - \pi_s)\pi_\omega \end{cases}. \quad (2)$$

<sup>5</sup> Although it is attractive to think of mismatch as administrative data being linked to the wrong individual, other cases may be covered as well. For instance in some cases administrative data may be formally correct but measure something that is conceptually different. An example would be an individual who writes off a heavy capital loss in a given year. This may lead to low or even negative taxable income, while for most modeling purposes income would probably be defined differently.

Contamination can, for instance, be the result of erroneously reporting income as annual, whereas the amount is a monthly amount, or vice versa, omitting a second job or working only part of the year. Each of these cases may result in large differences between survey value and true value. A value of  $\rho$  smaller than zero implies mean-reverting response error in the sense of Bound and Krueger (1991).

Next define the difference between survey data and administrative data,  $m_i$ , as

$$m_i \equiv s_i - r_i. \quad (3)$$

Note that the difference between survey data and administrative data no longer equals the response error, since in this model the administrative data may contain error. It is useful to calculate a number of moments of  $r_i$ ,  $s_i$ , and  $m_i$  to gain insight in the behavior of the model and its differences with a model without administrative error. We have (for derivations, see app. A):

$$\mu_r \equiv E(r_i) = \pi_r \mu_\xi + (1 - \pi_r) \mu_\zeta, \quad (4)$$

$$\mu_s \equiv E(s_i) = \mu_\xi + (1 - \pi_s) [\mu_\eta + \pi_\omega \mu_\omega], \quad (5)$$

$$\mu_m \equiv E(s_i) - E(r_i) = (1 - \pi_r) [\mu_\xi - \mu_\zeta] + (1 - \pi_s) [\mu_\eta + \pi_\omega \mu_\omega], \quad (6)$$

$$\sigma_r^2 \equiv \text{Var}(r_i) = \pi_r \sigma_\xi^2 + (1 - \pi_r) \sigma_\zeta^2 + \pi_r (1 - \pi_r) [\mu_\xi - \mu_\zeta]^2, \quad (7)$$

$$\begin{aligned} \sigma_s^2 \equiv \text{Var}(s_i) \\ = \pi_s \sigma_\xi^2 + (1 - \pi_s) [(1 + \rho)^2 \sigma_\xi^2 + \sigma_\eta^2 + \pi_\omega \sigma_\omega^2 + \pi_s [\mu_\eta + \pi_\omega \mu_\omega]^2 \\ + \pi_\omega (1 - \pi_\omega) \mu_\omega^2], \end{aligned} \quad (8)$$

$$\begin{aligned} \sigma_m^2 \equiv \text{Var}(m_i) \\ = [(1 - \pi_r) \pi_s + \pi_r (1 - \pi_s) \rho^2 + (1 - \pi_r) (1 - \pi_s) (1 + \rho)^2] \sigma_\xi^2 \\ + (1 - \pi_r) \sigma_\zeta^2 + (1 - \pi_s) [\sigma_\eta^2 + \pi_\omega \sigma_\omega^2] + \pi_r (1 - \pi_r) [\mu_\xi - \mu_\zeta]^2 \\ + \pi_s (1 - \pi_s) [\mu_\eta + \pi_\omega \mu_\omega]^2 + (1 - \pi_s) \pi_\omega (1 - \pi_\omega) \mu_\omega^2, \end{aligned} \quad (9)$$

$$\begin{aligned} \sigma_{mr} \equiv E[(m_i - \mu_m)(r_i)] \\ = \rho \pi_r (1 - \pi_s) \sigma_\xi^2 - (1 - \pi_r) \sigma_\zeta^2 - \pi_r (1 - \pi_r) [\mu_\xi - \mu_\zeta]^2. \end{aligned} \quad (10)$$

The last expression can be seen as a measure of mean reversion we expect to see in the data under the (possibly incorrect) assumption that the administrative data are measured without error. We note that  $\sigma_{mr}$  is, for negative  $\rho$ , unambiguously negative, implying indeed mean reversion. Observe, however, that  $\sigma_{mr}$  is still negative if  $\rho = 0$ . That is, even without



“true” mean reversion the data will still suggest that mean reversion is present as long as  $\pi_r \neq 1$ , that is, if the administrative data suffer from at least some mismatch. This is not unique to the current set-up and is essentially an example of regression toward the mean. As soon as  $r_i$  suffers from measurement error, we expect  $s_i - r_i$  to be negatively correlated with  $r_i$ . As a second observation we consider

$$\begin{aligned} \sigma_s^2 - \sigma_r^2 &= \sigma_\xi^2[\pi_s - \pi_r + (1 - \pi_s)(1 + \rho)^2] + (1 - \pi_s)[\sigma_\eta^2 + \pi_\omega\sigma_\omega^2] \\ &\quad - (1 - \pi_r)\sigma_\xi^2 - \pi_r(1 - \pi_r)[\mu_\xi - \mu_s]^2. \end{aligned} \quad (11)$$

For  $\pi_r = 1$ , that is, no mismatch in the administrative data, (11) reduces to

$$\sigma_s^2 - \sigma_r^2 = (1 - \pi_s)\{[(1 + \rho)^2 - 1]\sigma_\xi^2 + \sigma_\eta^2 + \pi_\omega\sigma_\omega^2\}, \quad (12)$$

which shows that if the administrative data are assumed to be measured perfectly, the variance of the survey data can only be smaller than the variance of the administrative data if the survey data exhibit mean reversion ( $\rho < 0$ ). As a matter of fact, (12) will be negative if

$$(1 + \rho)^2 < 1 - \frac{\sigma_\eta^2 + \pi_\omega\sigma_\omega^2}{\sigma_\xi^2}. \quad (13)$$

Thus, the bigger the measurement errors in the survey data are assumed to be, the more mean reversion one needs to rationalize the data.

Under the scenario that neither the survey data nor the administrative data are flawless the question arises, which of the two should one use in modeling, and under what circumstances. An alternative, related, question would be: given that survey data are usually more easily available than administrative data, how much do we lose by using survey data rather than administrative data? Below we illustrate the answers to these questions for a very simple linear model. One could also ask, what is the best way of combining survey and administrative data if both are available, but that question is beyond the scope of this article.<sup>6</sup>

#### Implications if $\xi_i$ Is a Dependent Variable

We consider a very simple linear regression model of the form:

$$\xi_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad (14)$$

where we make the conventional assumption that  $\varepsilon_i$  is uncorrelated with  $x_i$  (and with any of the other random variables we have defined so far). If we have only survey data available, we would replace  $\xi_i$  by  $s_i$  on the

<sup>6</sup> Our ML estimation using both administrative and survey data is one answer to that question.

left-hand side of (14). Denote the ensuing estimate of  $\beta_1$  by  $\hat{\beta}_1$ . Then we have

$$p \lim \hat{\beta}_1 = \beta_1 [1 + \rho(1 - \pi_s)]. \quad (15)$$

Thus the estimator is consistent if  $\rho = 0$ , that is, if there is no mean reverting error. The biasing effect of the mean reversion is mitigated by the observations that are exactly correct.

If we replace  $\xi_i$  by  $r_i$  we obtain for the probability limit of the estimator (say  $\tilde{\beta}_1$ ):

$$p \lim \tilde{\beta}_1 = \pi_r \beta_1, \quad (16)$$

so that the percent bias is equal to the percentage of mismatched administrative observations. Comparing (15) to (16) shows that, for  $\rho < 0$ , using survey data will lead to less bias if

$$1 + \rho(1 - \pi_s) > \pi_r. \quad (17)$$

Clearly, this always holds if  $\rho = 0$  and the administrative data are not perfect.

#### Implications if $\xi_i$ Is an Independent Variable

Now consider a model of the form:

$$z_i = \gamma_0 + \gamma_1 \xi_i + \nu_i. \quad (18)$$

Let  $\hat{\gamma}_1$  be the ordinary least squares (OLS) estimator of  $\gamma_1$  if we replace  $\xi_i$  by  $s_i$ . It is straightforward to show that

$$p \lim \hat{\gamma}_1 = \gamma_1 \frac{[1 + \rho(1 - \pi_s)]\sigma_\xi^2}{\pi_s \sigma_\xi^2 + (1 - \pi_s)((1 + \rho)^2 \sigma_\xi^2 + \sigma_\eta^2 + \pi_\omega \sigma_\omega^2 + \pi_s [\mu_\eta + \pi_\omega \mu_\omega]^2 + \pi_\omega (1 - \pi_\omega) \mu_\omega^2)}, \quad (19)$$

where the denominator is equal to (8). Clearly if  $\rho = 0$ ,  $\pi_s = 0$ , and  $\pi_\omega = 0$ , we are back in the case of classical measurement error. In that case (19) reduces to

$$p \lim \hat{\gamma}_1 = \gamma_1 [\sigma_\xi^2 / (\sigma_\xi^2 + \sigma_\eta^2)].$$

For  $\rho = 0$  and  $\pi_s, \pi_\omega \neq 0$  (19) reduces to

$$p \lim \hat{\gamma}_1 = \gamma_1 \frac{\sigma_\xi^2}{\sigma_\xi^2 + (1 - \pi_s)(\sigma_\eta^2 + \pi_\omega \sigma_\omega^2 + \pi_s [\mu_\eta + \pi_\omega \mu_\omega]^2 + \pi_\omega (1 - \pi_\omega) \mu_\omega^2)}. \quad (20)$$

As one would expect, having a proportion  $\pi_s$  of exactly measured variables mitigates the attenuation due to measurement error in the explanatory variable, while on the other hand a proportion  $\pi_\omega$  of contaminated samples

worsens the effect. Having measurement errors or contaminated samples with mean unequal to zero also increases the bias.

Let  $\tilde{\gamma}_1$  be the OLS estimator of  $\gamma_1$  if we replace  $\xi_i$  by  $\eta_i$ . Now we obtain:

$$p \lim \tilde{\gamma}_1 = \gamma_1 \frac{\pi_r \sigma_\xi^2}{\pi_r \sigma_\xi^2 + (1 - \pi_r) \sigma_\zeta^2 + \pi_r (1 - \pi_r) (\mu_\xi - \mu_\zeta)^2}, \quad (21)$$

where the denominator is equal to (7). Clearly  $\tilde{\gamma}_1$  is consistent for  $\pi_r = 1$ , as it should be. A special case of interest is where a mismatch is a drawing from the same distribution, that is,  $\mu_\xi = \mu_\zeta$ ,  $\sigma_\xi^2 = \sigma_\zeta^2$ . In that case (21) reduces to  $p \lim \tilde{\gamma}_1 = \gamma_1 \pi_r$ . This is exactly the same shrinkage as in (16).

Comparison of biases introduced by either using less than perfect survey data or partly mismatched administrative data shows that survey data are to be preferred if

$$\begin{aligned} & \frac{[1 + \rho(1 - \pi_r)] \sigma_\xi^2}{\pi_r \sigma_\xi^2 + (1 - \pi_r) [(1 + \rho) \sigma_\xi^2 + \sigma_\eta^2 + \pi_\omega \sigma_\omega^2 + \pi_r [\mu_\eta + \pi_\omega \mu_\omega]^2 + \pi_\omega (1 - \pi_\omega) \mu_\omega^2]} \\ & > \frac{\pi_r \sigma_\xi^2}{\pi_r \sigma_\xi^2 + (1 - \pi_r) \sigma_\zeta^2 + \pi_r (1 - \pi_r) (\mu_\xi - \mu_\zeta)^2}, \end{aligned} \quad (22)$$

which is not particularly informative. More insight can be gained for the case where  $\rho = 0$  (which is close to the empirically relevant case as we shall see for our data),  $\mu_\omega = \mu_\eta = 0$  and  $\mu_\xi = \mu_\zeta$ ,  $\sigma_\xi^2 = \sigma_\zeta^2$ . In that case (22) reduces to

$$\frac{\sigma_\xi^2}{\sigma_\xi^2 + (1 - \pi_r) [\sigma_\eta^2 + \pi_\omega \sigma_\omega^2]} > \pi_r. \quad (23)$$

Thus, survey data exhibit less bias if the reliability ratio of the survey data is greater than the proportion of perfect administrative data. In the more general case where  $\mu_\xi \neq \mu_\zeta$  and  $\sigma_\xi^2 \geq \sigma_\zeta^2$ , the balance tips a little more in favor of the survey data.

#### IV. Estimation and Results

We extend the model, discussed in the last section, by including covariates. We parameterize  $\mu_\xi$  as a function of individual characteristics. The “true” variable,  $\xi_i$ , is assumed to be dependent on variables such as gender, age, and education. In this case we have

$$\xi_i = x_i \beta + \varepsilon_i, \quad (24)$$

where  $\varepsilon_i$  is Gaussian noise. Note that, when covariates are included, mean-reversion gets a slightly different interpretation. Values do not get adjusted toward the overall mean, but toward the mean within a group of people

with the same values of  $x_i$ . Both the model with covariates, where  $\xi_i \sim N(x_i\beta, \sigma_\xi^2)$ , and the model without covariates, where  $\xi_i \sim N(\mu_\xi, \sigma_\xi^2)$ , are estimated. The covariates used are age, age<sup>2</sup>, and dummies for gender, degree of education, and self-assessed retirement status.<sup>7</sup> These are all survey values in order to avoid the mismatching problems accompanying administrative values. One drawback of this method is possible error in survey data and a lot of don't know values for education (coded as *dkedu*).

For estimation we write the model as a mixture of one univariate normal distribution, when  $r_i = \xi_i$  and  $s_i = \xi_i$ , and five bivariate normal distributions with different means and covariance matrices. Maximum likelihood is used to obtain estimates. For a detailed description of the estimation procedure we refer to appendix B. It is perhaps somewhat remarkable that such a rich error structure can be identified. This is the direct result of the nonnormality of the error structure. It has been known for quite a while that normality is a very unfavorable assumption for identifiability of parameters in linear errors in variables models (see, e.g., Aigner et al. 1984; and Bekker 1986). Kane, Rouse, and Staiger (1999) provide an example of how exploiting nonnormality aids identifiability. Meijer and Ypma (2006) provide a simple proof of identification for the case of a mixture of two normal distributions, of which the current model is a generalization. A different strand of literature examines how bounds on measurement error can be exploited to bound parameter estimates (see, e.g., Klepper and Leamer 1984; and Bekker, Kapteyn, and Wansbeek 1987). In the current context such bounds do not seem necessary, although conceivably this would further narrow the range of plausible parameter estimates.

Appendix C presents all estimation results. For each of the variables of interest, earnings, pensions, and taxes, two tables are presented—one with covariates included and one without. Besides the full model, three other models are estimated, with certain constraints imposed on the full model. These include a model without contamination of the survey data,  $\pi_\omega = 0$ , a model where no mismatching occurs,  $\pi_r = 1$ , and a model where both are left out. This last model can be seen as the model used in most previous studies, with only the addition that survey observations are equal to the truth with positive probability.<sup>8</sup>

<sup>7</sup> The LINDA variable used to define age is *bald*, which is the age of the individual at the end of the tax year. We define  $\text{age} = \text{bald} - 40$  and  $\text{age}^2 = \text{bald} - 40^2/100$ .

<sup>8</sup> Since it is sometimes found that women provide more accurate answers than men (see, e.g., Bound and Krueger 1991), we have also estimated a model where  $\mu_e$  and  $\sigma_e$  are allowed to differ by gender, but the differences are negligible. For instance, for earnings the estimate of  $\sigma_e$  is .100 for men and .104 for women; the *t*-value of the difference is .33. The *t*-value for the difference in  $\mu_e$  is even lower: .10.

**Table 5**  
**Proportional Biases Resulting from Using Administrative or Survey Data**

	Earnings		Pensions		Taxes	
	Administrative	Survey	Administrative	Survey	Administrative	Survey
$\xi_i$ LHS	.959	.989	.981	.904	.935	.991
$\xi_i$ RHS	.491	.701	.719	.363	.789	.735

NOTE.—Cells present proportional asymptotic biases in OLS estimates if we replace the true variable by administrative or survey data; 1 means no bias.

Table C2 presents the estimation results for earnings when no covariates are included. The most striking observation is probably that allowing for mismatches or contaminated samples in the model leads to a dramatic fall in the estimated value of the mean-reversion parameter,  $\rho$ . Only when we do not allow for either contamination or mismatches, do we reproduce the “stylized fact” of substantial mean reversion. Table C1, which includes covariates, leads to qualitatively similar conclusions. The fact that allowing for mismatch may lead to a sharp drop in the estimate of  $\rho$  is consistent with equation (10), which shows that a negative covariance between the “measurement error” and the administrative values can be generated by mismatches even if  $\rho = 0$ . The fact that contamination of the survey data can also rationalize a negative covariance between  $m$  and  $r$  is a little harder to grasp intuitively.

With respect to the pattern of estimated mean reversion, tables C4 and C3 (for pensions) provide a qualitatively similar picture, although mean reversion in the full model is somewhat higher than for earnings. For taxes, mean reversion in the full model and in the models with only mismatches or only contaminated samples is essentially zero (tables C6 and C5).

The estimated percentage of correct survey data,  $\pi_s$ , ranges from 15% in the earnings data to a little over 25% in the pension data. The fraction of contaminated survey values,  $\pi_a(1 - \pi_s)$ , lies between .04 and .13.

A parameter of particular interest is  $\pi_s$ , or rather  $1 - \pi_s$ , the proportion of mismatched administrative values. The estimate of  $1 - \pi_s$  varies from 2% in the pension data to about 8% in the tax data. We can use equations (15), (16), (19), and (21) to assess the biases that would arise in the estimate of a slope parameter if  $\xi_i$  were either a dependent (left-hand side) or an independent (right-hand side) variable in a univariate regression. Table 5 presents the proportional biases for both cases and for the three variables we are considering in this article.<sup>9</sup> When  $\xi_i$  is a dependent (left-hand side) variable, biases are modest, with the administrative data leading to slightly smaller biases than the survey data for pensions, whereas for earnings and taxes survey data yield less bias. Biases are on the order of 5%. When

<sup>9</sup> We use the estimates for the models without covariates.

$\xi_i$  is a right-hand side variable, the picture is dramatically different. Biases are much larger, up to 65%. When using administrative data, we find that the bias is largest for earnings, about 50%. Inspection of formula (21) shows that this is due to the fact that the 4% mismatched data appear to be drawn from a distribution that has a much lower log-mean ( $\mu_\xi = 9.187$ , while  $\mu_\xi = 12.283$ ) and substantially higher dispersion ( $\sigma_\xi = 1.807$ , while  $\sigma_\xi = .717$ ). For the survey data the bias is particularly large in the case of pensions, about 65%. Inspection of formula (19) suggests that the main cause for the big bias lies in the low mean and high variance of the contaminated data ( $\mu_\omega = -1.632$ ,  $\sigma_\omega = 3.801$ ).

A different way to gauge the importance of a proper treatment of the various error sources is to compare different conventional estimation methods and how their results differ from the full model. Tables C7–C9 present estimates of the parameters of economic interest for a number of different estimation methods: ML on the full model; OLS, robust regression, and median regression. The latter three estimation methods are applied twice, once with administrative data as dependent variable and once with survey data as dependent variable.

Considering earnings (table C7) we note that the estimates of the effect of education on earnings may vary by at least a factor of two, depending on the estimation method chosen. Table C7 suggests that running OLS of the administrative variables on the explanatory variables provides estimates close to those of the full model. However when we consider pensions that conclusion changes quite a bit. Robust regression and median regression yield estimates that may be quite far removed from the estimates obtained with ML on the full model. Another noteworthy phenomenon is the wide variation in the estimates of the effect of the semi-retired dummy on pensions (table C8). Estimates vary from significantly positive to significantly negative depending on the estimation method used and the choice of dependent variable. For taxes the different estimation methods appear to provide roughly comparable estimates of parameters of economic interest.

Since the full model is a mixture of a number of different regimes it is of interest to assign observations to different regimes. To do so we have used the fact that the likelihood for each observation is a weighted sum of densities corresponding to the different regimes. We have assigned observations to the regime that produces the highest density for that observation. Figures C1–C3 present the results. For the earnings data there is some suggestion that respondents with low earnings tend to give high survey values, as also found by Bollinger (1998). Of particular interest are the observations that are classified as mismatched administrative variables. Most of these points lie above the 45 degree line, whereas the points classified as contaminated in the survey data lie below the 45 degree line. Naturally, the assignment procedure used here is merely indicative

and probabilistic, and hence no great importance should be attached to the classification of each observation.

It is of interest to see if some of the classifications may be externally validated in some way. For about 70% of the LINDA sample, a third source of earnings information is available from employer records. It consists of two variables, one containing the full-time equivalent monthly earnings and one containing the percentage of full-time equivalent employed. Although there are some problems with this information, we can still compare these values with the survey value and the administrative value. A simple way to look for possible mismatches is to consider observations where survey and employer earnings are fairly close, while differing substantially from the administrative data. For instance, if we select observations for which survey earnings and administrative earnings are at least 50,000 SEK apart, while survey earnings and employer earnings differ by less than 10,000 SEK, we find 12 cases where this condition holds true. For these 12 observations the ratio of survey and administrative earnings varies between .5 and 4.5.

## V. Conclusion

In comparison with most studies in the literature we have allowed for a richer specification of possible error sources in survey data and administrative data. Our results suggest that some conclusions in the literature may be quite sensitive to the assumption that administrative data are flawless. In a sense, the question if administrative data represent the truth, is almost a philosophical question. For instance, the examples of detected mismatches given above do not necessarily imply there is true mismatching going on in the administrative data. Rather, it appears that sometimes the survey data (and the alternative source of administrative data) measure a different concept than the administrative data. Be this as it may, also under the latter interpretation one would be hard pressed to maintain that the difference between survey data and administrative data exhibits strong mean reversion.

Our results also suggest that substantive conclusions may be affected quite a bit by changes in assumptions on the nature of error in survey and administrative data. Application of robust methods, such as median or robust regression, yields results quite far removed from ML on the full model. Thus these methods do not appear to provide a solution for dealing with different sources of error in survey or administrative data.

There are many good reasons for wanting to use administrative data, including sample sizes, cost of surveys, and data quality. However, unless administrative data measure exactly the concept that one is interested in and do so without error, these data are not a panacea. Our illustrative calculations in table 5 suggest that biases resulting from using adminis-

trative data as right-hand side variables may be very substantial. As always, one has to be careful in modeling the sources and nature of errors and take that into account when investigating hypotheses of substantive interest.

## Appendix A

### Derivation of Moments

The expectation of  $r_i$

$$\begin{aligned}\mu_r &\equiv E[r_i] = \pi_r E[r_i | r_i = \xi_i] + (1 - \pi_r) E[r_i | r_i = \zeta_i] \\ &= \pi_r E[\xi_i] + (1 - \pi_r) E[\zeta_i] \\ &= \pi_r \mu_\xi + (1 - \pi_r) \mu_\zeta.\end{aligned}$$

The expectation of  $s_i$

$$\begin{aligned}\mu_s &\equiv E[s_i] \\ &= \pi_s E[s_i | s_i = \xi_i] + (1 - \pi_s)(1 - \pi_\omega) E[s_i | s_i = \xi_i + \rho(\xi_i - \mu_\xi) + \eta_i] \\ &\quad + (1 - \pi_s)\pi_\omega E[s_i | s_i = \xi_i + \rho(\xi_i - \mu_\xi) + \eta_i + \omega_i] \\ &= \pi_s E[\xi_i] + (1 - \pi_s)(1 - \pi_\omega) E[\xi_i + \rho(\xi_i - \mu_\xi) + \eta_i] \\ &\quad + (1 - \pi_s)\pi_\omega E[\xi_i + \rho(\xi_i - \mu_\xi) + \eta_i + \omega_i] \\ &= \pi_s \mu_\xi + (1 - \pi_s)(1 - \pi_\omega)[\mu_\xi + \mu_\eta] + (1 - \pi_s)\pi_\omega[\mu_\xi + \mu_\eta + \mu_\omega] \\ &= \mu_\xi + (1 - \pi_s)[\mu_\eta + \pi_\omega \mu_\omega].\end{aligned}$$

Expectation of  $m_i$

$$\begin{aligned}\mu_m &\equiv E[m_i] = E[s_i - r_i] = E[s_i] - E[r_i] \\ &= (1 - \pi_r)[\mu_\xi - \mu_\zeta] + (1 - \pi_s)[\mu_\eta + \pi_\omega \mu_\omega].\end{aligned}$$

The variance of  $r_i$  is

$$\begin{aligned}\sigma_r^2 &\equiv E[r_i - \mu_r]^2 \\ &= \pi_r E[r_i - \mu_r | r_i = \xi_i]^2 + (1 - \pi_r) E[r_i - \mu_r | r_i = \zeta_i]^2 \\ &= \pi_r E[\xi_i - \pi_r \mu_\xi - (1 - \pi_r) \mu_\zeta]^2 + (1 - \pi_r) E[\zeta_i - \pi_r \mu_\xi - (1 - \pi_r) \mu_\zeta]^2 \\ &= \pi_r E[(\xi_i - \mu_\xi) + (1 - \pi_r)(\mu_\xi - \mu_\zeta)]^2 + (1 - \pi_r) E[(\zeta_i - \mu_\zeta) - \pi_r(\mu_\xi - \mu_\zeta)]^2 \\ &= \pi_r [\sigma_\xi^2 + (1 - \pi_r)^2 (\mu_\xi - \mu_\zeta)^2] + (1 - \pi_r) [\sigma_\zeta^2 + \pi_r^2 (\mu_\xi - \mu_\zeta)^2] \\ &= \pi_r \sigma_\xi^2 + (1 - \pi_r) \sigma_\zeta^2 + \pi_r (1 - \pi_r) (\mu_\xi - \mu_\zeta)^2.\end{aligned}$$



The variance of  $s_i$  is

$$\begin{aligned}
 \sigma_s^2 &= E[s_i - \mu_s]^2 \\
 &= \pi_s E[s_i - \mu_s | s_i = \xi_i]^2 \\
 &\quad + (1 - \pi_s)(1 - \pi_\omega) E[s_i - \mu_s | s_i = \xi_i + \rho(\xi_i - \mu_\xi) + \eta_i]^2 \\
 &\quad + (1 - \pi_s)\pi_\omega E[s_i - \mu_s | s_i = \xi_i + \rho(\xi_i - \mu_\xi) + \eta_i + \omega_i]^2.
 \end{aligned} \tag{A1}$$

First calculate the three variances separately. We have

$$\begin{aligned}
 E[s_i - \mu_s | s_i = \xi_i]^2 &= E[\xi_i - \mu_\xi - (1 - \pi_s)(\mu_\eta + \pi_\omega \mu_\omega)]^2 \\
 &= E[\xi_i - \mu_\xi - (1 - \pi_s)(\mu_\eta + \pi_\omega \mu_\omega)]^2 \\
 &= \sigma_\xi^2 + (1 - \pi_s)^2 (\mu_\eta + \pi_\omega \mu_\omega)^2,
 \end{aligned} \tag{A2}$$

and

$$\begin{aligned}
 E[s_i - \mu_s | s_i = \xi_i + \rho(\xi_i - \mu_\xi) + \eta_i]^2 &= E[\xi_i + \rho(\xi_i - \mu_\xi) + \eta_i - \mu_\xi - (1 - \pi_s)(\mu_\eta + \pi_\omega \mu_\omega)]^2 \\
 &= E[(1 + \rho)(\xi_i - \mu_\xi) + (\eta_i - \mu_\eta) + \pi_s(\mu_\eta + \pi_\omega \mu_\omega) - \pi_\omega \mu_\omega]^2 \\
 &= (1 + \rho)^2 \sigma_\xi^2 + \sigma_\eta^2 + [\pi_s(\mu_\eta + \pi_\omega \mu_\omega) - \pi_\omega \mu_\omega]^2,
 \end{aligned} \tag{A3}$$

and

$$\begin{aligned}
 E[s_i - \mu_s | s_i = \xi_i + \rho(\xi_i - \mu_\xi) + \eta_i + \omega_i]^2 &= E[\xi_i + \rho(\xi_i - \mu_\xi) + \eta_i + \omega_i - \mu_\xi - (1 - \pi_s)(\mu_\eta + \pi_\omega \mu_\omega)]^2 \\
 &= E[(1 + \rho)(\xi_i - \mu_\xi) + (\eta_i - \mu_\eta) + (\omega_i - \mu_\omega) + \pi_s(\mu_\eta + \pi_\omega \mu_\omega) + (1 - \pi_\omega)\mu_\omega]^2 \\
 &= (1 + \rho)^2 \sigma_\xi^2 + \sigma_\eta^2 + \sigma_\omega^2 + [\pi_s(\mu_\eta + \pi_\omega \mu_\omega) + (1 - \pi_\omega)\mu_\omega]^2.
 \end{aligned} \tag{A4}$$

Define

$$\delta = \mu_\xi - \mu_\eta, \quad \text{and} \quad \Delta = \mu_\eta + \pi_\omega \mu_\omega.$$

Substituting (A2), (A3), and (A4) in (A1) we have

$$\begin{aligned}
 \sigma_s^2 &= E[s_i - \mu_s]^2 \\
 &= \pi_s [\sigma_\xi^2 + (1 - \pi_s)^2 \Delta^2] \\
 &\quad + (1 - \pi_s)(1 - \pi_\omega) [(1 + \rho)^2 \sigma_\xi^2 + \sigma_\eta^2 + [\pi_s \Delta - \pi_\omega \mu_\omega]^2] \\
 &\quad + (1 - \pi_s)\pi_\omega [(1 + \rho)^2 \sigma_\xi^2 + \sigma_\eta^2 + \sigma_\omega^2 + [\pi_s \Delta + (1 - \pi_\omega)\mu_\omega]^2] \\
 &= [\pi_s + (1 - \pi_s)(1 + \rho)^2] \sigma_\xi^2 + (1 - \pi_s) [\sigma_\eta^2 + \pi_\omega \sigma_\omega^2] + \Omega,
 \end{aligned} \tag{A5}$$

where  $\Omega$  is

$$\begin{aligned}
 \Omega &= \pi_s(1 - \pi_s)^2\Delta^2 + (1 - \pi_s)(1 - \pi_\omega)[\pi_s\Delta - \pi_\omega\mu_\omega]^2 + (1 - \pi_s)\pi_\omega[\pi_s\Delta + (1 - \pi_\omega)\mu_\omega]^2 \\
 &= \pi_s(1 - \pi_s)^2\Delta^2 + (1 - \pi_s)\pi_s^2\Delta^2 + (1 - \pi_s)(1 - \pi_\omega)\pi_\omega^2\mu_\omega^2 + (1 - \pi_s)\pi_\omega(1 - \pi_\omega)^2\mu_\omega^2 \\
 &\quad - 2(1 - \pi_s)(1 - \pi_\omega)\pi_s\Delta\pi_\omega\mu_\omega + 2(1 - \pi_s)\pi_\omega\pi_s\Delta(1 - \pi_\omega)\mu_\omega \\
 &= [\pi_s(1 - \pi_s)^2 + (1 - \pi_s)\pi_s^2]\Delta^2 + (1 - \pi_s)[(1 - \pi_\omega)\pi_\omega^2 + \pi_\omega(1 - \pi_\omega)^2]\mu_\omega^2 \\
 &= \pi_s(1 - \pi_s)\Delta^2 + (1 - \pi_s)\pi_\omega(1 - \pi_\omega)\mu_\omega^2.
 \end{aligned} \tag{A6}$$

We can now calculate the variance using (A5), (A6), and the definition of  $\Delta$ :

$$\begin{aligned}
 \sigma_s^2 &= [\pi_s + (1 - \pi_s)(1 + \rho)^2]\sigma_\xi^2 + (1 - \pi_s)[\sigma_\eta^2 + \pi_\omega\sigma_\omega^2] \\
 &\quad + \pi_s(1 - \pi_s)(\mu_\eta + \pi_\omega\mu_\omega)^2 + (1 - \pi_s)\pi_\omega(1 - \pi_\omega)\mu_\omega^2.
 \end{aligned}$$

Using the same procedure as above, the covariance between  $r$  and  $s$  is

$$\begin{aligned}
 \sigma_{rs} &\equiv E[r_i - \mu_r][s_i - \mu_s] \\
 &= \pi_r\pi_s(\sigma_\xi^2 - (1 - \pi_r)(1 - \pi_s)\delta\Delta) + \\
 &\quad + \pi_r(1 - \pi_s)(1 - \pi_\omega)((1 + \rho)\sigma_\xi^2 + (1 - \pi_r)\pi_s\delta\Delta - (1 - \pi_r)\delta\pi_\omega\mu_\omega) \\
 &\quad + \pi_r(1 - \pi_s)\pi_\omega((1 + \rho)\sigma_\xi^2 + (1 - \pi_r)\pi_s\delta\Delta + (1 - \pi_r)\delta(1 - \pi_\omega)\mu_\omega) \\
 &\quad + (1 - \pi_r)\pi_s(\pi_r(1 - \pi_s)\delta\Delta) \\
 &\quad + (1 - \pi_r)(1 - \pi_s)(1 - \pi_\omega)(-\pi_r\pi_s\delta\Delta + \pi_r\delta\pi_\omega\mu_\omega) \\
 &\quad + (1 - \pi_r)(1 - \pi_s)\pi_\omega(-\pi_r\pi_s\delta\Delta - \pi_r\delta(1 - \pi_\omega)\mu_\omega) \\
 &= \pi_r\sigma_\xi^2 + \pi_r(1 - \pi_s)\rho\sigma_\xi^2.
 \end{aligned}$$

The variance of  $m_i$  is then easily obtained as

$$\begin{aligned}
 \sigma_m^2 &\equiv E[m_i - \mu_m]^2 = E[(s_i - \mu_s) - (r_i - \mu_r)]^2 = \sigma_s^2 + \sigma_r^2 - 2\sigma_{rs} \\
 &= [\pi_s + (1 - \pi_s)(1 + \rho)^2 + \pi_r - 2\pi_r - 2\pi_r(1 - \pi_s)\rho]\sigma_\xi^2 \\
 &\quad + (1 - \pi_r)\sigma_\eta^2 + (1 - \pi_s)[\sigma_\eta^2 + \pi_\omega\sigma_\omega^2] + \pi_r(1 - \pi_r)[\mu_\xi - \mu_\eta]^2 \\
 &\quad + \pi_s(1 - \pi_s)[\mu_\eta + \pi_\omega\mu_\omega]^2 + (1 - \pi_s)\pi_\omega(1 - \pi_\omega)\mu_\omega^2,
 \end{aligned}$$

which is equivalent to

$$\begin{aligned}
 \sigma_m^2 &= [(1 - \pi_r)\pi_s + \pi_r(1 - \pi_s)\rho^2 + (1 - \pi_s)(1 - \pi_s)(1 + \rho)^2]\sigma_\xi^2 \\
 &\quad + (1 - \pi_r)\sigma_\eta^2 + (1 - \pi_s)[\sigma_\eta^2 + \pi_\omega\sigma_\omega^2] + \pi_r(1 - \pi_r)[\mu_\xi - \mu_\eta]^2 \\
 &\quad + \pi_s(1 - \pi_s)[\mu_\eta + \pi_\omega\mu_\omega]^2 + (1 - \pi_s)\pi_\omega(1 - \pi_\omega)\mu_\omega^2.
 \end{aligned}$$

Finally, the covariance between  $m$  and  $r$  is

$$\begin{aligned}
 \sigma_{mr} &\equiv \sigma_{sr} - \sigma_r^2 \\
 &= \pi_r \sigma_\xi^2 + \pi_r(1 - \pi_r) \rho \sigma_\xi^2 - [\pi_r \sigma_\xi^2 + (1 - \pi_r) \sigma_\xi^2 + \pi_r(1 - \pi_r)(\mu_\xi - \mu_r)^2] \\
 &= \rho \pi_r(1 - \pi_r) \sigma_\xi^2 - (1 - \pi_r) \sigma_\xi^2 - \pi_r(1 - \pi_r)(\mu_\xi - \mu_r)^2.
 \end{aligned}$$

## Appendix B

### Maximum Likelihood

As described in Section IV, we assume that the administrative data are a mixture of two normal distributions, while the survey data are a mixture of three different normal distributions. Since we assume the processes underlying the administrative data and the survey data to be independent, the combined set of observations  $(r_i, s_i)$  follow a mixture of six distributions.

The general shape of the log-likelihood function of a mixture of  $M$  distributions and  $N$  observations is the following (Redner and Walker 1984):

$$l(\theta) = \sum_{i=1}^N \log \left( \sum_{m=1}^M \pi_m f_m(x_i | \theta) \right),$$

where  $\theta$  is the vector of parameters, including mixing proportions  $\pi_m$  and parameters describing the distributions. Redner and Walker (1984) mention some special cases of mixtures, for instance, when some of the observations are labeled. In our case, the observations where  $r_i = s_i$  can be seen as labeled observations. We assume that these observations come from the distribution, where both administrative and survey data are correct, referred to as group 1. Since all of the observations from group 1 can be labeled, we have a completely labeled group.

Let's assume that observations  $i = 1, \dots, n_1$  are observations from the completely labeled group 1, and the other observations,  $i = n_1 + 1, \dots, N$  are a mixture of the remaining five distributions. The following log-likelihood can then be derived

$$l(\theta) = \sum_{i=1}^{n_1} \log(\pi_1 f_1(x_i | \theta)) + \sum_{i=n_1+1}^N \log \left( \sum_{m=2}^5 \pi_m f_m(x_i | \theta) \right),$$

where

$$\begin{aligned}
 \pi_1 &= \pi_r \pi_s \\
 \pi_2 &= \pi_r (1 - \pi_s) (1 - \pi_\omega) \\
 \pi_3 &= \pi_r (1 - \pi_s) \pi_\omega \\
 \pi_4 &= (1 - \pi_r) \pi_s \\
 \pi_5 &= (1 - \pi_r) (1 - \pi_s) (1 - \pi_\omega) \\
 \pi_6 &= (1 - \pi_r) (1 - \pi_s) \pi_\omega.
 \end{aligned}$$

The density function  $f_1$  is the probability density function of a  $N(\mu_\xi, \sigma_\xi^2)$  distribution, since both  $r_i$  and  $s_i$  are equal to  $\xi_i$  in this case. In the five other cases we have the bivariate normal distributions listed below:

$$f_2(r_i, s_i) \sim N \left[ \begin{pmatrix} \mu_\xi \\ \mu_\xi + \mu_\eta \end{pmatrix}, \begin{pmatrix} \sigma_\xi^2 & \frac{(1+\rho)\sigma_\xi^2}{\sigma_\xi \sqrt{(1+\rho)^2 \sigma_\xi^2 + \sigma_\eta^2}} \\ \frac{(1+\rho)\sigma_\xi^2}{\sigma_\xi \sqrt{(1+\rho)^2 \sigma_\xi^2 + \sigma_\eta^2}} & (1+\rho)^2 \sigma_\xi^2 + \sigma_\eta^2 \end{pmatrix} \right], \quad (B1)$$

when  $r_i = \xi_i$  and  $s_i = \xi_i + \rho(\xi_i - \mu_\xi) + \eta_i$ ;

$$f_3(r_i, s_i) \sim N \left[ \begin{pmatrix} \mu_\xi \\ \mu_\xi + \mu_\eta + \mu_\omega \end{pmatrix}, \begin{pmatrix} \sigma_\xi^2 & \frac{(1+\rho)\sigma_\xi^2}{\sigma_\xi \sqrt{(1+\rho)^2 \sigma_\xi^2 + \sigma_\eta^2 + \sigma_\omega^2}} \\ \frac{(1+\rho)\sigma_\xi^2}{\sigma_\xi \sqrt{(1+\rho)^2 \sigma_\xi^2 + \sigma_\eta^2 + \sigma_\omega^2}} & (1+\rho)^2 \sigma_\xi^2 + \sigma_\eta^2 + \sigma_\omega^2 \end{pmatrix} \right], \quad (B2)$$

when  $r_i = \xi_i$  and  $s_i = \xi_i + \rho(\xi_i - \mu_\xi) + \eta_i + \omega_i$ ;

$$f_4(r_i, s_i) \sim N \left[ \begin{pmatrix} \mu_\xi \\ \mu_\xi \end{pmatrix}, \begin{pmatrix} \sigma_\xi^2 & 0 \\ 0 & \sigma_\xi^2 \end{pmatrix} \right], \quad (B3)$$

when  $r_i = \zeta_i$  and  $s_i = \xi_i$ ;

$$f_5(r_i, s_i) \sim N \left[ \begin{pmatrix} \mu_\zeta \\ \mu_\xi + \mu_\eta \end{pmatrix}, \begin{pmatrix} \sigma_\zeta^2 & 0 \\ 0 & (1+\rho)^2 \sigma_\xi^2 + \sigma_\eta^2 \end{pmatrix} \right], \quad (B4)$$

when  $r_i = \zeta_i$  and  $s_i = \xi_i + \rho(\xi_i - \mu_\xi) + \eta_i$ ; and

$$f_6(r_i, s_i) \sim N \left[ \begin{pmatrix} \mu_\zeta \\ \mu_\xi + \mu_\eta + \mu_\omega \end{pmatrix}, \begin{pmatrix} \sigma_\zeta^2 & 0 \\ 0 & (1+\rho)^2 \sigma_\xi^2 + \sigma_\eta^2 + \sigma_\omega^2 \end{pmatrix} \right], \quad (B5)$$

when  $r_i = \zeta_i$  and  $s_i = \xi_i + \rho(\xi_i - \mu_\xi) + \eta_i + \omega_i$ . One final note has to be made about the labeling of group 1. Observations are labeled as a member of group 1 if the difference between  $r_i$  and  $s_i$  is smaller than 1,000 SEK. The proportion  $\pi_s$  then is the proportion of survey observations that differ less than 1,000 SEK from the administrative data. In principle, this broader definition of “equal” observations affects the consistency of our estimates. However, we expect these effects to be minor.

## Appendix C

### Estimation Results

**Table C1**  
**Estimates Using Log Earnings**

	Full Model		No Contamination		No Mismatch		Basic Model	
	Coefficient	SD	Coefficient	SD	Coefficient	SD	Coefficient	SD
Log likelihood	-503.2		-558.8		-623.4		-762.8	
Female	-.268	.052	-.253	.062	-.163	.100	-.225	.071
Age	.087	.024	.095	.031	.136	.025	.111	.034
Age <sup>2</sup>	-.275	.062	-.275	.082	-.421	.059	-.330	.095
Fullret	-1.015	.125	-1.152	.091	-1.253	.127	-1.210	.114
Semiret	-.421	.114	-.273	.112	-.694	.276	-.378	.106
Midedu	.244	.102	.370	.117	.266	.128	.361	.228
Highedu	.345	.101	.438	.122	.364	.148	.435	.211
Dkedu	.568	.104	.694	.119	.603	.120	.684	.197
$\mu_{\xi}$	11.592	.262	11.335	.315	11.092	.301	11.181	.287
$\sigma_{\xi}$	.543	.021	.653	.021	.840	.032	.754	.031
$\mu_{\gamma}$	9.843	.646	11.385	.243	...	...	...	...
$\sigma_{\gamma}$	2.026	.334	1.738	.170	...	...	...	...
$\mu_{\omega}$	-.259	.211	...	...	.442	.234	...	...
$\sigma_{\omega}$	1.313	.152	...	...	1.791	.174	...	...
$\mu_{\eta}$	-.046	.007	-.053	.008	-.040	.008	...	...
$\sigma_{\eta}$	.102	.008	.110	.007	.104	.006	.523	.020
$\pi_{\gamma}$	.948	.015	.864	.020	...	...	...	...
$\pi_{\gamma}$	.155	.019	.170	.020	.148	.018	...	...
$\pi_{\omega}$	.123	.028	...	...	.168	.024	...	...
$\rho$	-.064	.020	-.002	.024	-.072	.021	-.525	.036

**Table C2**  
**Estimates Using Log Earnings**

	Full Model		No Contamination		No Mismatch		Basic Model	
	Coefficient	SD	Coefficient	SD	Coefficient	SD	Coefficient	SD
Log likelihood	-607.5		-646.5		-708.5		-881.2	
$\mu_{\xi}$	12.283	.032	12.246	.041	12.178	.046	12.191	.038
$\sigma_{\xi}$	.717	.021	.843	.026	1.116	.035	.960	.030
$\mu_{\gamma}$	9.187	.691	11.387	.254	...	...	...	...
$\sigma_{\gamma}$	1.807	.424	1.751	.178	...	...	...	...
$\mu_{\omega}$	-.304	.174	...	...	.432	.187	...	...
$\sigma_{\omega}$	1.239	.129	...	...	1.407	.151	...	...
$\mu_{\eta}$	-.048	.007	-.056	.007	-.047	.008	...	...
$\sigma_{\eta}$	.099	.007	.112	.006	.100	.007	.552	.019
$\pi_{\gamma}$	.959	.013	.867	.020	...	...	...	...
$\pi_{\gamma}$	.152	.018	.169	.020	.148	.018	...	...
$\pi_{\omega}$	.156	.028	...	...	.187	.026	...	...
$\rho$	-.013	.014	.022	.012	-.013	.014	-.369	.029

**Table C3**  
Estimates Using Log-Pensions

	Full Model		No Contamination		No Mismatch		Basic Model	
	Coefficient	SD	Coefficient	SD	Coefficient	SD	Coefficient	SD
Log likelihood	-589.5		-755.7		-653.6		-775.8	
Female	-.422	.048	-.423	.106	-.375	.063	-.410	.059
Age	.058	.025	.020	.061	.077	.029	.072	.028
Age <sup>2</sup>	-.085	.041	-.011	.100	-.111	.048	-.100	.047
Fullret	.527	.154	.276	.268	.575	.163	.531	.159
Semiret	.034	.187	-.487	.554	-.219	.253	-.150	.273
Midedu	.224	.055	.167	.124	.201	.075	.165	.069
Highedu	.156	.089	.352	.167	.111	.090	.088	.099
Dkedu	.487	.068	.520	.155	.489	.099	.407	.091
$\mu_{\xi}$	10.379	.429	10.964	.814	9.948	.379	10.059	.390
$\sigma_{\xi}$	.501	.021	1.060	.071	.669	.027	.572	.021
$\mu_{\zeta}$	8.957	.307	11.234	.226	...	...	...	...
$\sigma_{\zeta}$	.764	.235	1.206	.160	...	...	...	...
$\mu_{\omega}$	-1.472	.983	...	...	-.305	.809	...	...
$\sigma_{\omega}$	3.676	.713	...	...	3.707	.610	...	...
$\mu_{\eta}$	-.049	.013	-.037	.022	-.040	.015	...	...
$\sigma_{\eta}$	.212	.012	.187	.012	.209	.012	.838	.031
$\pi_r$	.981	.007	.907	.021	...	...	...	...
$\pi_s$	.267	.023	.295	.025	.263	.023	...	...
$\pi_{\omega}$	.056	.018	...	...	.077	.019	...	...
$\rho$	-.199	.032	-.161	.027	-.186	.031	-.459	.073

**Table C4**  
Estimates Using Log-Pensions

	Full Model		No Contamination		No Mismatch		Basic Model	
	Coefficient	SD	Coefficient	SD	Coefficient	SD	Coefficient	SD
Log likelihood	-650.5		-777.3		-696.4		-830.3	
$\mu_{\xi}$	11.742	.032	11.679	.056	11.693	.038	11.685	.025
$\sigma_{\xi}$	.628	.030	1.123	.033	.769	.033	.657	.023
$\mu_{\zeta}$	9.023	.409	11.256	.222	...	...	...	...
$\sigma_{\zeta}$	.843	.344	1.202	.159	...	...	...	...
$\mu_{\omega}$	-1.632	1.077	...	...	-.302	.773	...	...
$\sigma_{\omega}$	3.801	.775	...	...	3.608	.584	...	...
$\mu_{\eta}$	-.044	.014	-.036	.015	-.038	.015	...	...
$\sigma_{\eta}$	.217	.012	.190	.012	.211	.012	.845	.022
$\pi_r$	.981	.008	.905	.022	...	...	...	...
$\pi_s$	.268	.023	.295	.025	.263	.023	...	...
$\pi_{\omega}$	.050	.017	...	...	.080	.020	...	...
$\rho$	-.131	.023	-.117	.022	-.128	.023	-.361	.067

**Table C5**  
**Estimates Using Log-Taxes**

	Full Model		No Contamination		No Mismatch		Basic Model	
	Coefficient	SD	Coefficient	SD	Coefficient	SD	Coefficient	SD
Log likelihood	-759.2		-818.5		-788.0		-924.4	
Female	-.453	.064	-.412	.080	-.471	.076	-.509	.059
Age	.030	.023	.028	.022	.039	.020	.035	.020
Age <sup>2</sup>	-.070	.040	-.076	.042	-.097	.035	-.087	.036
Fullret	-.409	.129	-.368	.117	-.381	.132	-.320	.100
Semiret	.153	.120	.234	.129	.156	.157	.198	.183
Midedu	.332	.083	.290	.095	.291	.091	.284	.077
Highedu	.381	.103	.406	.110	.389	.107	.394	.099
Dkedu	.648	.110	.664	.099	.648	.109	.641	.097
$\mu_{\xi}$	10.752	.293	10.741	.269	10.647	.258	10.709	.256
$\sigma_{\xi}$	.686	.026	.890	.028	.790	.028	.689	.022
$\mu_{\zeta}$	9.846	.300	10.235	.137	...	...	...	...
$\sigma_{\zeta}$	1.197	.145	1.097	.094	...	...	...	...
$\mu_{\omega}$	-.591	.360	...	...	.025	.188	...	...
$\sigma_{\omega}$	1.847	.267	...	...	1.639	.136	...	...
$\mu_{\eta}$	.095	.007	.097	.007	.101	.005	...	...
$\sigma_{\eta}$	.105	.006	.103	.006	.106	.006	.567	.018
$\pi_r$	.921	.023	.848	.020	...	...	...	...
$\pi_s$	.216	.020	.239	.021	.199	.018	...	...
$\pi_{\omega}$	.094	.028	...	...	.180	.024	...	...
$\rho$	-.043	.013	-.043	.015	-.056	.014	-.192	.037

**Table C6**  
**Estimates Using Log-Taxes**

	Full Model		No Contamination		No Mismatch		Basic Model	
	Coefficient	SD	Coefficient	SD	Coefficient	SD	Coefficient	SD
Log likelihood	-840.0		-876.0		-862.8		-1018.2	
$\mu_{\xi}$	10.839	.089	10.801	.043	10.784	.042	10.809	.034
$\sigma_{\xi}$	.846	.117	1.006	.037	.958	.019	.828	.023
$\mu_{\zeta}$	9.645	.448	10.220	.140	...	...	...	...
$\sigma_{\zeta}$	1.189	.167	1.109	.095	...	...	...	...
$\mu_{\omega}$	-.474	.442	...	...	.051	.174	...	...
$\sigma_{\omega}$	1.640	.471	...	...	1.500	.130	...	...
$\mu_{\eta}$	.092	.010	.088	.010	.095	.007	...	...
$\sigma_{\eta}$	.107	.008	.108	.007	.105	.005	.572	.015
$\pi_r$	.935	.031	.853	.020	...	...	...	...
$\pi_s$	.213	.020	.235	.021	.199	.018	...	...
$\pi_{\omega}$	.111	.045	...	...	.189	.024	...	...
$\rho$	-.011	.021	-.0004	.019	-.015	.012	-.133	.030





Table C8

545

Table C9

## Both Sources

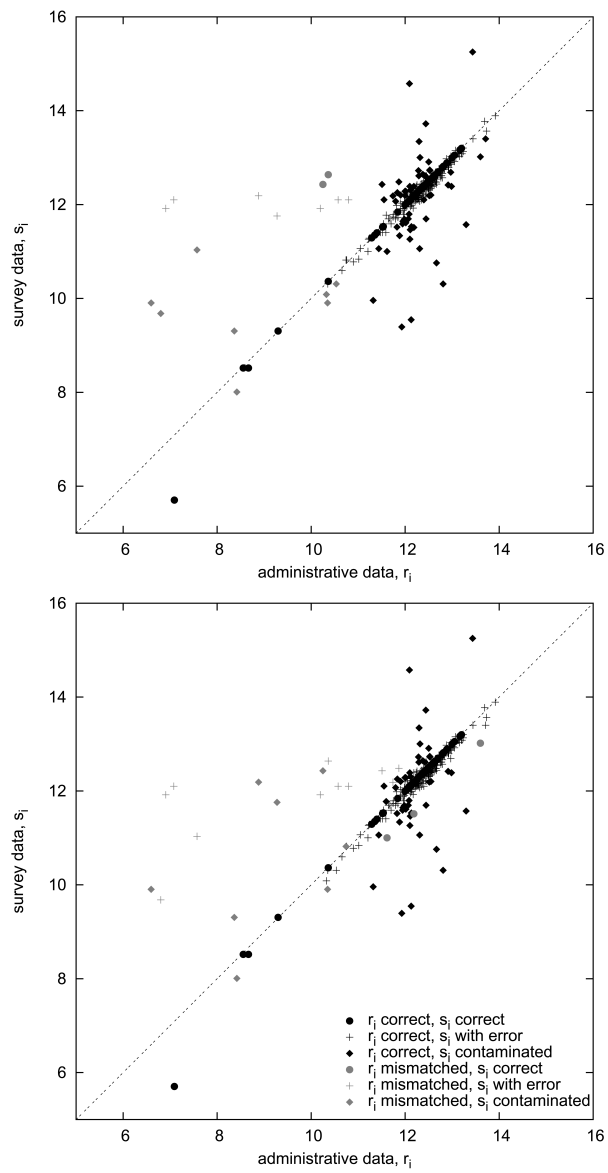


FIG. C1.—Log earnings without (upper) and with (lower) covariates

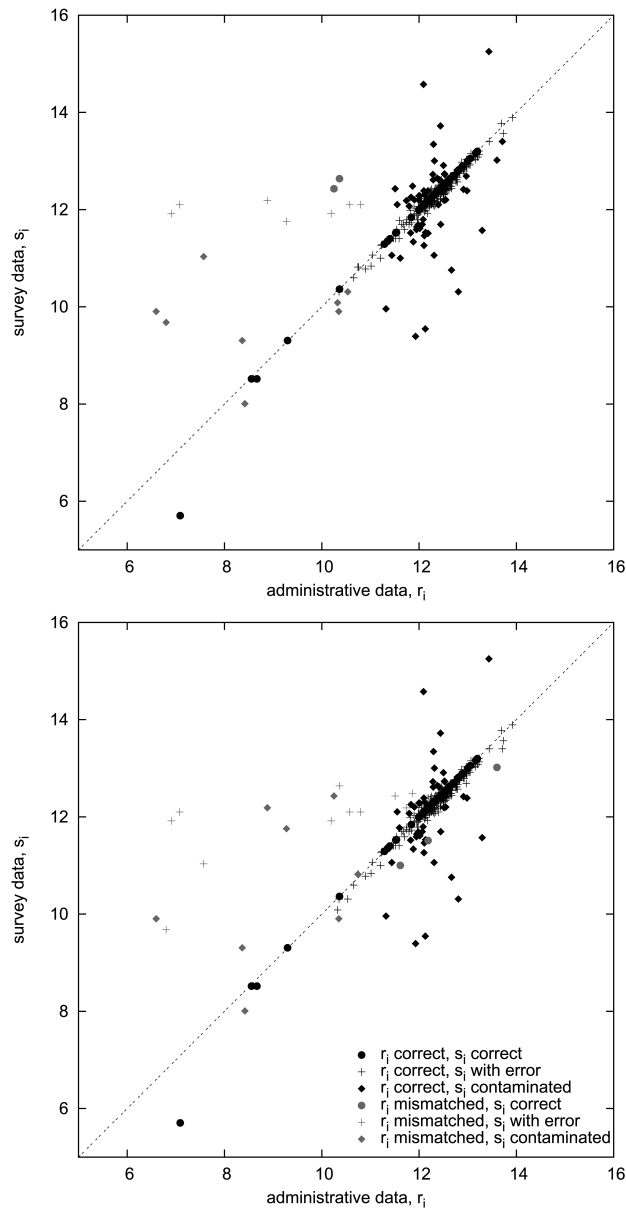


FIG. C2.—Log pensions without (upper) and with (lower) covariates. Two observations, with logarithmic survey values smaller than 3 and logarithmic administrative values between 11 and 11.5, are outside the scale of this figure. Both are classified as correct administrative data and a contaminated survey value.

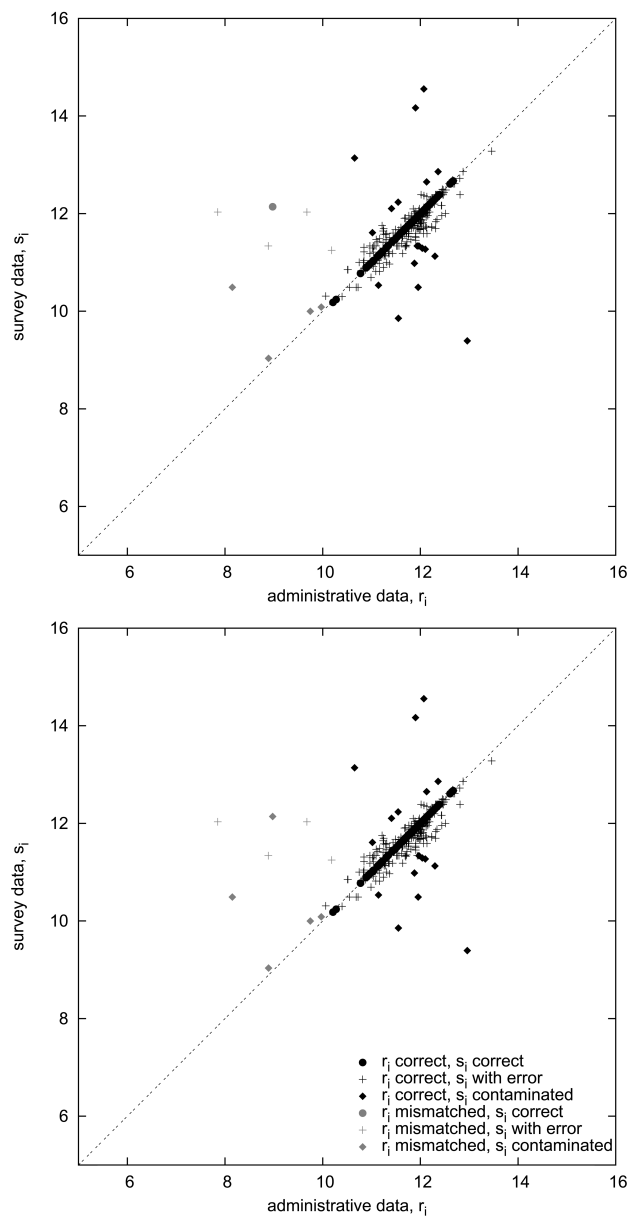


FIG. C3.—Log taxes without (upper) and with (lower) covariates. One observation, with logarithmic survey value smaller than 4 and logarithmic administrative value between 11 and 11.5, is outside the scale of this figure. The survey value is classified as contaminated and the administrative value is classified as correct (with covariates) or as a mismatch (without covariates).

## References

- Abowd, John M., and Lars Vilhuber. 2005. The sensitivity of economic statistics to coding errors in personal identifiers. *Journal of Business and Economic Statistics* 23, no. 2 (April): 133–52.
- Aigner, Dennis, Cheng Hsiao, Arie Kapteyn, and Tom Wansbeek. 1984. Latent variable models in econometrics. In *Handbook of econometrics*, vol. 2, ed. Zvi Griliches and Michael Intriligator, 1321–93. Amsterdam: Elsevier.
- Bekker, Paul A. 1986. Comment on identification in the linear errors in variables model. *Econometrica* 54:215–17.
- Bekker, Paul, Arie Kapteyn, and Tom Wansbeek. 1987. Consistent sets of estimates for regressions with correlated or uncorrelated measurement error in arbitrary subsets of all variables. *Econometrica* 55: 1223–30.
- Bollinger, Christopher R. 1998. Measurement error in the current population survey: A nonparametric look. *Journal of Labor Economics* 16, no. 3 (July): 576–94.
- Bound, John, Charles Brown, Greg J. Duncan, and Willard L. Rodgers. 1994. Evidence on the validity of cross-sectional and longitudinal labor market data. *Journal of Labor Economics* 12, no. 3 (July): 345–68.
- Bound, J., C. Brown, and N. Mathiowetz. 2001. Measurement error in survey data. In *Handbook of labor economics*, vol. 5, ed. J. Heckman and E. Leamer, 3707–3843. Amsterdam: Elsevier.
- Bound, John, and Alan B. Krueger. 1991. The extent of measurement error in longitudinal earnings data: Do two wrongs make a right? *Journal of Labor Economics* 9, no. 1 (January): 1–24.
- Chen, Xiaohong, Han Hong, and Elie Tamer. 2005. Measurement error models with auxiliary data. *Review of Economic Studies* 72, no. 2 (April): 343–66.
- Duncan, Greg J., and Daniel H. Hill. 1985. An investigation of the extent and consequences of measurement error in labor-economic survey data. *Journal of Labor Economics* 3, no. 4 (October): 508–32.
- Edin, Per-Anders, and Peter Frederiksson. 2000. LINDA—Longitudinal INdividual DAta for Sweden. Working Paper no. 2000:19, Department of Economics, Uppsala University (November).
- Hurd, Michael, F. Thomas Juster, and James P. Smith. 2004. Enhancing the quality of data on income: Recent developments in survey methodology. *Labor and Demography* 0412001, EconWPA (December), <http://ideas.repec.org/p/wpa/wuwpla/0412001.html>.
- Johansson, Fredrik, and Anders Klevmarken. 2006. Explaining the size and nature of response behavior in a survey on health status and economic standard. Working Paper no. 2006:2, Department of Economics, Uppsala University (January).

- Kane, Thomas, Cecilia Rouse, and Douglas Staiger. 1999. Estimating returns to schooling when schooling is misreported. NBER Working Paper no. 7235, National Bureau of Economic Research, Cambridge, MA.
- Kim, Bonggeun, and Gary Solon. 2005. Implications of mean-reverting measurement error for longitudinal studies of wages and employment. *Review of Economics and Statistics* 87, no. 1 (December): 193–96.
- Klepper, Steven, and Edward Leamer. 1984. Consistent sets of estimates for regressions with errors in all variables. *Econometrica* 52:163–84.
- Lee, Lung-Fei, and Jungsywan H. Sepanski. 1995. Estimation of linear and nonlinear errors-in-variables models using validation data. *Journal of the American Statistical Association* 90, no. 429 (March): 130–40.
- Meijer, Erik, and Jelmer Ypma. 2006. A simple identification proof for a mixture of two univariate normal distributions. Working paper, University of Groningen (March).
- Pedace, Roberto, and Nancy Bates. 2001. Using administrative records to assess earnings reporting error in the survey of income and program participation. *Journal of Economic and Social Measurement* 26, nos. 3–4:173–92.
- Pischke, Jörn-Steffen. 1995. Measurement error and earnings dynamics: Some estimates from the PSID validation study. *Journal of Business and Economic Statistics* 13, no. 3 (July): 305–14.
- Redner, Richard A., and Homer F. Walker. 1984. Mixture densities, maximum likelihood and the EM algorithm. *SIAM Review* 26, no. 2 (April): 195–239.
- Rodgers, Willard L., Charles Brown, and Greg J. Duncan. 1993. Errors in survey reports of earnings, hours worked, and hourly wages. *Journal of the American Statistical Association* 88, no. 424 (December): 1208–18.
- Stinson, Martha. 2002. Estimating measurement error in SIPP annual job earnings: A comparison of census survey and SSA administrative data. U.S. Census Bureau Technical Report TP-2002-24, Suitland, MD (September).