

MEASUREMENT ERROR IN THE BANK OF ITALY'S SURVEY OF HOUSEHOLD INCOME AND WEALTH

BY CLAUDIA BIANCOTTI,* GIOVANNI D'ALESSIO AND ANDREA NERI

Bank of Italy, Economic and Financial Statistics Department

This paper is aimed at evaluating the incidence of measurement error in the Bank of Italy's Survey of Household Income and Wealth (SHIW). In the case of time-invariant variables, we assess the degree of inconsistency of answers given by panel households in subsequent survey waves. For quantities that vary with time, we estimate the incidence of measurement error by decomposing observed variability into true dynamics and error-induced noise. We apply the Heise model or the latent Markov model, depending on whether the data are continuous or categorical. We also present regression models that explain the error-generating process. Our results are relevant to researchers who use SHIW data for economic analysis, but also to data producers involved in similar income and wealth surveys. The methods we describe and test can be employed in a number of contexts to gain better understanding of data-related problems and plans for survey improvement.

1. INTRODUCTION

Estimates that are based on sample surveys are subject to a number of possible errors. A first source of inaccuracy is implicit in the nature of the inferential process that yields population parameters on the basis of a set of sampled units; a type of disturbance known as *sampling error* arises, whose incidence can be evaluated precisely if we are aware of some features of the sample (e.g. size, design) and of the population (e.g. distribution of the variable we are interested in).¹

Other causes of imprecision involve the process of measurement and estimation; the resulting mistakes are known as *non-sampling errors*, the sometimes unavoidable costs of transforming a theoretical scheme into an actual survey. Broadly speaking, the literature on these topics focuses on the problems relating to the following aspects: (a) sample composition, as a consequence of incomplete sampling frames (non-coverage) or failure to participate in the survey on the part of some sampled units (non-response); and (b) discrepancies between recorded data and "true" data, originating from response error or oversights in the processing phase prior to estimation.

The effects induced by some types of error on the estimated values of aggregates in the Bank of Italy's Survey of Household Income and Wealth (SHIW from here on) have been studied in the past. For example, sampling errors are normally published along with survey results (Banca d'Italia, 2002); the consequences of

Note: We would like to thank Luigi Cannari, Ivan Faiella, Grazia Marchese, Luigi Federico Signorini and two anonymous referees for their useful insights. The opinions expressed in this paper are those of the authors and should not be attributed to the Bank of Italy.

*Correspondence to: Claudia Biancotti, Bank of Italy, Economic and Financial Statistics Department, Via Nazionale 91, I-00184 Rome, Italy (claudia.biancotti@bancaditalia.it).

¹For an overview of survey methodology, including an in-depth discussion of sampling error, see Kish (1995).

non-response on the most important estimates have been assessed;² and efforts have been made to evaluate the magnitude of under-reporting of assets and income.³

These analyses notwithstanding, several areas of the data quality territory remain relatively uncharted, especially in relation to response error. The questionnaire is not a neutral instrument: the order and wording of questions and the available response options influence the answers, especially (but not only) where opinions, expectations and other subjective items are concerned. Interviewer behavior is also very important: there are a number of ways of asking the same question in a face-to-face setting, and each can induce a different psychological reaction, ultimately affecting the answer. Further problems can arise from the respondent's cognitive processes: hypothetical questions require some abstract reasoning, retrospective ones need an effort to recall events of the past.⁴ Moreover, people may not actually know the exact answer to the questions they are asked, especially in cases (such as the SHIW) where response by proxy is allowed. Following Groves and Couper (1998), general aspects such as motivation of the respondent and willingness to give time and effort for a survey should also be assumed to influence data quality. Finally, the use of a Computer-Assisted Personal Interviewing electronic interface (CAPI)—complete with range controls, consistency assessment, and outlier detection tools—instead of a printed form can influence the answers.

The possible causes of response error, as summarized above, appear to be too numerous to be tackled in a single paper. We will therefore concentrate on the impact that certain features of fieldwork operations, of the interviewers and of the respondents have on data quality.

The statistical analyses that follow are based entirely on survey data; hence, they may not meet the strict randomization criteria needed for controlled experiments. Some caveats apply to the conclusions: they are only as reliable as the models used to eliminate possible sources of noise. This will be discussed later.

As a further warning, note that this paper is a first exploration of a subset of measurement error issues in the SHIW: it gives some elements to evaluate the magnitude of imprecision in the data collection process, and on the possible reasons why it exists. It only mentions in passing the consequences exerted by

²Cannari and D'Alessio (1992) analyze the behavior of panel households and find that non-response is a common trait in large cities and in Northern Italy. The participation rate decreases with income, and increases with household size. D'Alessio and Faiella (2002) confirm that well-off households and those headed by educated individuals are harder to interview, while households residing in Central Italy and headed by an individual in the central age groups are more likely to participate in the survey.

³The value of housing, which accounts for most of real wealth, appears to be underestimated by 20 percent; the figure is higher when referred to non-primary (vacation etc.) housing only. Financial assets are also exposed to under-reporting (Cannari *et al.*, 1990; Cannari and D'Alessio, 1993) as well as income deriving from self-employment and from capital (Cannari and Violi, 1995; Brandolini, 1999).

⁴For example, a hypothetical question entailing a certain effort on the part of the respondent is asked to homeowners in the SHIW (Banca d'Italia, 2002): "Assuming you wanted to rent this dwelling, what monthly rent do you think could be charged? Do not include condominium charges, heating or other sundry expenses." Memory problems could arise in questions such as this one, directed to pensioners: "Recall when you began to receive your pension. What percentage of your last wage payment (monthly average earnings, if self-employed) was your first pension payment?"

errors on the estimation of mean values; it does not touch upon their impact on regression coefficients and other statistics,⁵ and it does not dwell on the techniques that can be employed to obtain robust estimators.⁶ A vast literature exists on these topics, with papers often devoted to the impact of a single type of error on a single type of estimate. Our work has a different goal: we set out to outline a road map for SHIW users, so that they know which variables must be handled with particular care. The choice of error correction techniques to be employed, if any, is left to the users themselves, seeing how it depends heavily on both context and preference. The insights we draw from the SHIW might also benefit researchers working with similar surveys, which are likely to be affected, at least partly, by the same issues. Finally, we believe that our paper can be useful for data producers looking for a framework to assess information quality in a standardized way.

The paper is structured as follows. Section 2 briefly describes the SHIW, with a special focus on what is relevant for data quality. Section 3 proposes a methodology for evaluating the degree of reliability of collected data. Section 4 presents some descriptive statistics on measurement error for the main SHIW aggregates (income, wealth, consumption) and their individual components. Section 5 puts forward models that try to explain the inconsistency in answers provided over the years by panel households. Section 6 concludes.

2. THE BANK OF ITALY'S SURVEY OF HOUSEHOLD INCOME AND WEALTH

Since 1962, the Bank of Italy conducted a survey on household budgets, examining economic behavior at the micro level.

In the recent waves, the sample size has been of about 8,000 households. The design is two-stage: 300 municipalities, stratified by region and population, are selected at random; and then households are drawn from the municipal registers.

Starting from the 1989 wave, a part of the sample (now roughly 50 percent) is made up of households with previous experience of SHIW participation, the so-called panel households. This structure permits the study of dynamic phenomena such as income, wealth and employment mobility.

The questionnaire always addresses the following topics: demographic structure of the household, educational and occupational features of each member, individual income, household wealth and consumption, housing. Variable monographic sections are added on the basis of specific needs.

⁵Even a completely random error, although devoid of consequences on the estimation on mean values and population totals, and affecting variance in a way that can be corrected by modifying the sample size, distorts a number of statistics such as quantiles or linear regression coefficients (on this point see Carroll *et al.*, 2006 and Wansbeek and Meijer, 2000), in ways that have to be studied by way of non-trivial models such as the ones presented by Biemer and Trewin (1997). Quite often the error is not completely random: we will see in subsequent paragraphs that it can depend on fieldwork, interviewer and/or respondent features. In this case, complex models are necessary to predict the effect of the error on the estimate; a separate paper should be devoted to this problem and possible remedies for each type of estimate.

⁶See Huber (1981).

The survey is materially implemented by a specialized company hired by the Bank of Italy. Those who agree to participate are interviewed personally in their homes, often with the CAPI interface.

Most aspects of fieldwork are documented: information is collected on interviewer features, CAPI use, the presence of household members other than the head⁷ during the interview, the date and time of the interview, and its length.

The typical interviewer is female, a professional (interviewing is her main job), slightly over 40, with a high school degree. In the Northern regions, which are the richest, interviewer turnover is higher: the mean value of experience in the job is 9 years, while in the South it is 11.5 years. These figures, together with the higher incidence of interviewers who are not professionals (36.9 vs. 25.3 percent) and who hold junior high school degrees only (16.5 vs. 6.1 percent), seem to reflect the difficulty of finding work encountered by high school leavers in the South. On average, 67 percent of interviews are carried out with the CAPI method, which is largely present in the areas where most interviewers are professionals. On one hand, the company in charge of the survey might be more inclined to give computers to stable employees than to short-term ones; on the other hand, people who do not interview full-time might not want to put effort into learning how CAPI is used.

Most heads of household are interviewed personally,⁸ so are 27.9 percent of their spouses or live-in partners, while the rest of the family members are normally not present during the interview; the rate of personal response decreases with the number indicating the position in the household, which is declared by the household head during the interview.

On average, an interview takes 55 minutes, with a standard deviation of 19 minutes. Interview length is explained by socio-demographic features, such as the number of members and of earners (single-person households take only 46 minutes, five-member households take an hour, large families over 70 minutes), and by income levels (46 minutes for households earning less than 10,000 Euros per year, 64 minutes for those over 40,000 Euros).⁹ Operational choices also influence the amount of time needed to complete the interview: paper questionnaires take 58 minutes on average, the CAPI method 54 minutes.

The distribution of the responding households per interviewer shows a certain variability: the mean is 33, and 75 percent of the cases fall between 8 and 60. The asymmetry is justified by the fact that during the last weeks of fieldwork the best interviewers get “recovery” assignments in order to boost the response rate.

⁷The head of the household is defined as the person responsible for the household’s economic decisions.

⁸The presence of the head of household is a necessary condition for the interview. The few cases of absence correspond to exceptional situations, such as the death of the household head between the end of the reference year and the day of the interview. The presence of the other members is recorded only on income-related annexes to the main questionnaire; the information is therefore unavailable for non-earners, 39.4 percent of the sample. As a consequence, some of the estimates in Table 3 are biased downwards.

⁹The questionnaire is structured in such a way that each household member, each job held, each real estate asset requires a separate form.

3. THE EVALUATION OF DATA RELIABILITY: METHODOLOGY

3.1. Time-Invariant Phenomena

Let X be a continuous variable measured with an additive error:

$$(1) \quad Y = X + e.$$

The measure Y differs from the true value X by a random component with the following properties:

$$(2) \quad E(e) = 0; E(X, e) = \sigma_{X,e} = 0; E(e, e) = \sigma_e^2.$$

This type of disturbance is called *homoskedastic, uncorrelated error*. Under these assumptions, the variance of Y is a biased estimator of the variance of X , since:

$$(3) \quad \sigma_Y^2 = \sigma_X^2 + \sigma_e^2 = \sigma_X^2 / \lambda^2 \text{ where } \lambda = \frac{\sigma_X}{\sigma_Y}.$$

The coefficient λ is known as the *reliability index*;¹⁰ it expresses the share of variability in Y that originates from the true phenomenon X (Lord and Novick, 1968).¹¹

This index can be interpreted in several ways, taking into account the impact that measurement errors as described in (2) can have on different statistics. For example, the expected value of the measurement Y is an unbiased estimator of the mean of X : $E(Y) = E(X) + E(e) = E(X)$. Still, the presence of an error induces a higher estimator variance:

$$(4) \quad \sigma_{E(Y)}^2 = \sigma_{E(X)}^2 + \sigma_e^2 / n = \sigma_{E(X)}^2 / \lambda^2.$$

¹⁰Following, among others, Hand *et al.* (2001), "A precise measurement procedure is one that has small variability . . . [A]n accurate measurement procedure, in contrast, not only possesses small variability, but also yields results close to what we think of as the true value. . . . The reliability of a measurement procedure is the same as its precision. The former term is typically used in the social sciences whereas the latter is used in the physical sciences." A reliability index evaluates the degree to which an instrument, in our case the SHIW questionnaire, yields results that portray reality consistently; it does not indicate the instrument's truthfulness. We want to see what additional distance between data collected in different waves is introduced by measurement error, possibly net of actual changes in the quantities studied; a reliability index does not assess the distance between *collected* data and *true* data. Moreover, a precise measurement is not necessarily accurate, as shown by the case of correct and consistent recording of false information; reliability indexes are not able to spot the presence of phenomena such as systematic under-reporting.

¹¹In the rest of the paper we will use the reliability index as a descriptive parameter for the specific sample we are dealing with, not as an estimate of the corresponding population parameter. We only consider the variability introduced by measurement error, ignoring the variability connected to the sampling process.

From a sampling point of view, (4) implies that λ allows us to determine the “effective” sample size $n^* = \lambda^2 n$, i.e. the size that would yield the same variance of the sample mean if there was no measurement error.¹²

Turning to correlation analysis, we can say that if measurement error on X is assumed to be uncorrelated with X and with another variable Z , measured free of error, then $\rho_{Y,Z} = \lambda_Y \rho_{X,Z}$. The correlation coefficient between X and Z is attenuated with intensity proportional to the reliability index of Y , the measure of X . If we had a measure W of the variable Z , affected by measurement error uncorrelated with Z , X and the error component on X , then we would have $\rho_{Y,W} = \lambda_Y \lambda_W \rho_{X,Z}$. Under these conditions, the coefficient of univariate linear regressions of the type $Z = \alpha + \beta_X X$ estimated on the basis of the observed variables Y and W are related to the true coefficients by way of $\beta_Y = \lambda_Y^2 \beta_X$.

We can estimate λ if we have two measurements of the same phenomenon.¹³ Let Y_1 and Y_2 be such measurements, with additive errors:

$$(5) \quad Y_1 = X + e_1; \quad Y_2 = X + e_2.$$

Let the assumptions in (2) hold, supplemented by

$$(6) \quad E(e_1, e_2) = \sigma_{e_1, e_2} = 0; \quad E(X_t, e_{t'}) = \sigma_{X_t, e_{t'}} = 0 \forall t, t'.$$

Under these conditions, the correlation coefficient between the two measurements Y_1 and Y_2 equals the square of the reliability index:

$$(7) \quad \rho_{y_1, y_2} = \sigma_{y_1, y_2} / \sigma_{y_1} \sigma_{y_2} = \sigma_x^2 / (\sigma_x^2 + \sigma_e^2) = \sigma_x^2 / \sigma_y^2 = \lambda^2.$$

Panel households are interviewed every two years and sometimes they are asked questions they have answered in preceding waves. For these variables, if they are time-invariant, a quantification of measurement error can be obtained by applying (7) to the two answers.

In case of a multiplicative measurement error $Y = Xe$, i.e. when the observed variable is distorted by an error that is randomly proportional to the value, under the hypotheses $E(e) = 1$; $E(X, e) = \sigma_{X, e} = 0$; $E(e, e) = \sigma_e^2$, and if the additional conditions (6) hold, both the meaning and the estimation of the reliability index do not change. In fact, the correlation coefficient between the two measurements Y_1 and Y_2 still equals the reliability index: $\rho_{y_1, y_2} = \lambda^2$. More in general, relaxing the assumption $E(e) = 1$, i.e. when the measurement implies a systematic bias of the mean (under-reporting or over-reporting), the correlation coefficient ρ_{y_1, y_2} becomes the square of the ratio of the coefficient of variation of the measured variable to the

¹²We might also say, equivalently, that $(1 - \lambda^2)/\lambda^2$ is the additional cost introduced by measurement error; for example, a reliability index λ of 0.8 implies a rise in survey costs of 56 percent; if there were no error, estimates with the same precision could be obtained with a sample smaller by 36 percent $(1 - \lambda^2)$.

¹³What was stated in relation to regression coefficients implies that we may also estimate the reliability of a variable based on a single wave; i.e. under the conditions necessary for the use of instrumental variables (IV), and in the presence of univariate relations that are sufficiently significant and exhaustive, the reliability index can be calculated as the ratio between the OLS regression coefficient and the IV regression coefficient.

coefficient of variation of the true variable: $\rho_{y_1,y_2} = \lambda^2 \mu_y / \mu_x$. The index can thus be interpreted as the reliability of the measure, once the bias is taken into account.¹⁴

Also, as we will show later, measurement error might have a stronger incidence for some types of respondents: for example, when a question requires a large cognitive effort, seniors and those with low educational qualifications might be more inclined to give wrong answers. Should this be the case, the homoskedasticity assumption in the model would be violated, but there is no reason to assume that the errors would be correlated to the true values, or that different measurements would be correlated to each other. Analogously to the case of proportional error discussed above, our reliability measure would retain its descriptive power (across groups or individuals), and it could still be estimated based on the correlation coefficient between measurements, although the estimate would no longer be efficient.

If we are dealing with categorical variables, the above models are no longer adequate and need to be revised. Let X be a categorical variable (with K categories) and Y its measurement. An index of reliability for categorical features measured twice (Y_1 and Y_2) on the same set of n units is the fraction of units that are classified consistently: $\lambda^* = \text{tr}(F)/n = \sum_i f_{ii}/n$ where F is the cross tabulation of Y_1 and Y_2 whose generic element is f_{ij} . The index λ^* , however, does not take into account the fact that consistent answers could be partly random: if the two measures Y_1 and Y_2 are independent random variables, the expected share of consistent unit is $\sum_i f_i f_i / n$. A version of the reliability index which controls for this effect (see Biemer and Trewin, 1997) can be obtained by normalizing the share of observed matching cases with respect to their expected incidence if the two measurements of Y_1 and Y_2 were independent:

$$(8) \quad \lambda^{**} = (\lambda^* - \sum_i f_i f_i / n) / (1 - \sum_i f_i f_i / n).$$

Both the indexes λ^* and λ^{**} can also be adopted to assess the reliability of the categories of the qualitative variables; in fact, you can compute them on the dummy variables opposing each category to the others. This can help in understanding where the main classification problems are.

3.2. Time-Varying Phenomena

The indexes discussed so far allow us to derive a measure of response errors on variables that are measured twice and independently. In the SHIW context, this is the case of some phenomena that do not vary with time; since there is a two-year interval between interviews, the risk of contamination between different waves is very low (respondents most probably do not remember what they said).¹⁵

¹⁴Assuming $Y = Xe$ and $e = b + u$, where b is a constant and u the error term $E(u) = 0$ independent both from X and from Y_1 and Y_2 , it can be easily derived that the covariance between the two measures Y_1 and Y_2 $\sigma_{y_1,y_2} = b\sigma_x^2$. The term $b = \mu_y / \mu_x$ represents the average share of X reported in the measures Y_1 and Y_2 . The correlation coefficient between the two measures Y_1 and Y_2 is thus $\rho_{y_1,y_2} = b\lambda^2 = \sigma_x^2 \mu_y / \sigma_y^2 \mu_x$. If $b = 1$ (absence of systematic bias), then $\rho_{y_1,y_2} = \lambda^2$.

¹⁵Respondents probably remember that they participated in the survey two years before, and this can have an influence on their motivation or on their attitude toward giving information perceived as sensitive. In turn, this impacts on some types of error.

The analysis of measurement errors on the large majority of collected variables, especially the most interesting ones such as income and wealth, requires more sophisticated instruments. If a quantity varies with time, it is necessary to distinguish actual change from movements induced by wrong measurement.

The reliability of data on time-varying quantities can be assessed with the Heise (1969) method: provided we have at least three separate measurements of a variable on the same panel units (e.g. answers to the same question in three survey waves), under mild regularity conditions we are able to separate real dynamics from measurement error.

Let X_1, X_2, X_3 be the true values of the variable X during periods 1, 2 and 3; Y_1, Y_2, Y_3 are the corresponding measurements, for which the following equation applies:

$$(9) \quad Y_t = X_t + e_t \forall t.$$

In addition to this, let X_1, X_2 and X_3 be pairwise related through independent, first-order autoregressive models, which do not need to be stationary:

$$(10) \quad X_1 = \delta_1$$

$$(11) \quad X_2 = \beta_{21}X_1 + \delta_2$$

$$(12) \quad X_3 = \beta_{32}X_2 + \delta_3$$

where $\beta_{t+1,t}$ is the autoregressive coefficient and δ_t is the process innovation. Innovations are uncorrelated pairwise.

Assuming that the level of reliability of a given variable does not vary with time, the correlation coefficient between the observed values Y_t and Y_{t+1} can be written as:

$$(13) \quad \rho_{Y_t, Y_{t+1}} = \lambda_{Y_t} \lambda_{Y_{t+1}} \rho_{X_t, X_{t+1}} = \lambda_Y^2 \rho_{X_t, X_{t+1}}.$$

The ratio between the coefficient observed and the one that would be observed in the absence of measurement error is therefore always smaller than 1 and equal to λ_Y^2 :

$$(14) \quad \frac{\rho_{Y_t, Y_{t+1}}}{\rho_{X_t, X_{t+1}}} = \lambda_Y^2.$$

Since the true values are related by way of independent, first-order autoregressive processes we can say that for each t the following holds:

$$(15) \quad \frac{\rho_{X_{t-1}, X_t} \rho_{X_t, X_{t+1}}}{\rho_{X_{t-1}, X_{t+1}}} = 1.$$

Substituting (14) into (15), the estimate of reliability can be written as:

$$(16) \quad \lambda_Y = \sqrt{\frac{\rho_{Y_{t-1}, Y_t} \rho_{Y_t, Y_{t+1}}}{\rho_{Y_{t-1}, Y_{t+1}}}} \text{ and more generally } \lambda_Y = \sqrt[n]{\frac{\prod_{s=t}^{t+n} \rho_{Y_s, Y_{s+1}}}{\rho_{Y_t, Y_{t+n}}}}$$

The main idea of the method is: if measurement errors are independent of time and of the underlying variable, the absolute value of the estimated autocorrelation coefficients turns out to be lower than what we would get if the observed value did not include measurement error. Assuming that the true values in the three periods X_1 , X_2 and X_3 are related via first-order autoregressive models, the method proposes an estimate of measurement reliability by comparing the product of one-step correlations $\rho_{Y_1, Y_2} \rho_{Y_2, Y_3}$ with the two-step correlation ρ_{Y_1, Y_3} . If no measurement error existed, the quantity $\rho_{Y_1, Y_2} \rho_{Y_2, Y_3}$ would be equal to ρ_{X_1, X_3} ; but measurement error actually impacts on the estimate with incidence proportional to the square of ρ_{Y_1, Y_3} . It is therefore possible to obtain an indicator of measurement reliability by separating the part that the model attributes to the actual variation of the underlying quantity.

For time-invariant variables, with $\rho_{X_1, X_2} = \rho_{X_2, X_3} = \rho_{X_1, X_3} = 1$, (16) yields $\lambda_Y = \sqrt{\rho_{Y_t, Y_{t+1}}}$.

As noted, the index is based on the hypothesis that two independent first-order autoregressive models are a good approximation of the data-generating process. If this assumption does not hold, i.e. if a direct effect of X_1 on X_3 exists, the specification remains as described in (10), (11) and (12), but we have $\beta_{31} = \beta_{21}\beta_{32} + \beta_{31}^*$, where β_{31}^* is the regression coefficient relating X_1 and X_3 in the model including X_2 ; (15) becomes:

$$(17) \quad \frac{\rho_{X_{t-1}, X_t} \rho_{X_t, X_{t+1}}}{\rho_{X_{t-1}, X_{t+1}}} + \frac{\rho_{X_{t-1}, X_{t+1}}^* \sqrt{(1 - \rho_{X_{t-1}, X_t}^2)(1 - \rho_{X_t, X_{t+1}}^2)}}{\rho_{X_{t-1}, X_{t+1}}} = 1$$

This equation allows us to draw some general conclusions on what happens when the assumption of independent AR(1) processes is violated.

If we denote by ξ the term $\frac{\rho_{X_{t-1}, X_{t+1}}^* \sqrt{(1 - \rho_{X_{t-1}, X_t}^2)(1 - \rho_{X_t, X_{t+1}}^2)}}{\rho_{X_{t-1}, X_{t+1}}}$, we can write

$$(18) \quad \lambda_{YAR(1)}^2 = \frac{\rho_{Y_{t-1}, Y_t} \rho_{Y_t, Y_{t+1}}}{\rho_{Y_{t-1}, Y_{t+1}}} = (1 - \xi) \lambda_{YAR(2)}^2$$

The Heise index measured under the AR(1) hypothesis is a distorted estimate of the reliability value that we would have if we took an AR(2) process into account. Since the partial and the simple correlation coefficients are usually positive or null if dealing with strongly persistent phenomena such as income and wealth, we can say that usually $0 < (1 - \xi) < 1$; applying the base Heise method if

the underlying data structure is AR(2) yields reliability estimates that are biased downwards.¹⁶

The Heise model can also be applied in case of multiplicative errors. Under the same hypotheses on error terms showed that in the case of time-invariant phenomena the equation (16) still holds.

The analysis of measurement error for time varying categorical variables relies on the latent Markov model (LMM).¹⁷ The LMM model mostly relies on the same assumptions of the Heise model: the state of the variable at time t only depends on the state at $t - 1$; for identification and simplicity of the results, it is typically assumed that the error component is time-homogeneous: $P(Y_t = y_t | X_t = x_t) = P(Y_{t-1} = y_t | X_{t-1} = x_t)$ for $2 \leq t \leq T$. If no further constraints are imposed, one needs at least three time points to identify the model.

Measurement error can be captured through a latent class formulation by assuming that each observation of the states (manifest variable) corresponds to a latent variable which measures the true distribution over the states. The transition structure for the latent variables has the form of a first-order Markov chain. Moreover, each occasion-specific observed variable (Y_t) depends only on the corresponding latent variable (X_t); in other words, manifest variables are assumed to be independent conditional on the latent ones. As a consequence, the covariation actually observed among manifest variables is due to each manifest variable's relationship to the latent variable. Contrary to what happens in the case of continuous variables, no assumption of is made about the relation between true variable and the error term.

Suppose a single categorical variable of interest X (with K measured levels) is measured at T occasions, and that Y_t denotes the response at occasion t , $1 \leq t \leq T$, and y_t is a particular level of Y_t . Let X_t denote an occasion-specific true latent variable with C latent levels and x_t a particular level at time t . The corresponding LMM has the form:

$$(19) \quad P(Y = y) = \sum_{x_1=1}^C P(X_1 = x_1) \prod_{t=2}^T P(X_t = x_t | X_{t-1} = x_{t-1}) \prod_{t=1}^T P(Y_t = y_t | X_t = x_t).$$

The LMM model consists of two parts. The first part describes the measurement of true systematic change. It is summarized by a *transition matrix* containing the estimated true transition probabilities $P(X_t = x_t | X_{t-1} = x_{t-1})$. The second part describes the measurement of spurious change resulting from measurement error and other types of randomness in the behavior of individuals. It is represented by a *response probability matrix* containing the conditional probabilities of manifest variables having value y given that the latent one has value x at time t : $P(Y_t = y_t | X_t = x_t)$. The closer the response probability matrix is to an identity

¹⁶It is not easy to derive an unbiased estimator for the AR(2) case; it is not possible to obtain the solution by substitution, since the observed $\rho_{t-1,t}^*$ includes the very measurement error that we want to isolate. One possible solution is the correction of the Heise index by estimating the ξ component with instrumental variables.

¹⁷The latent Markov model was introduced by Wiggins (1955); it is also referred to as a latent transition or hidden Markov model (see Wiggins, 1973; Langeheine and Van de Pol, 1994; Vermunt, 1997). Also refer to MacDonald and Zucchini (1997).

TABLE 1
RELIABILITY OF SEX, YEAR AND PLACE OF BIRTH OF RESPONDENTS (PERCENTAGES)

Waves	Sex		Year of Birth ⁽¹⁾	Place of Birth ⁽¹⁾
	Reliability Index λ^*	Adjusted Reliability Index λ^{**}	Reliability Index λ^*	Reliability Index λ^*
1989–1991	98.2	97.2	95.5	98.4
1991–1993	98.1	97.1	96.9	98.4
1993–1995	99.7	98.7	98.8	98.2
1995–1998	99.7	98.7	98.8	97.5
1998–2000	99.9	98.9	98.3	97.3

Notes: ⁽¹⁾Due to the distribution of answers among a large number of categories, year and place of birth have small expected random consistencies: the adjusted reliability index λ^{**} is thus approximately equal to the unadjusted reliability index λ^* .

matrix, the smaller is the measurement error of the variable: the probabilities along the main diagonal can thus be interpreted as measures of reliability.

4. THE EVALUATION OF DATA RELIABILITY: SOME EVIDENCE

4.1. Reliability of Time-Invariant Phenomena

Socio-demographic features that are either time-invariant (such as sex or year of birth) or subject to small changes only (such as educational qualification) are repeatedly measured on panel units. The study of discrepancies in reported values can shed some light on measurement error in the survey.

A number of inconsistencies emerge, even for the simplest questions. For example, 1.8 percent of respondents declared a different gender in 1989 and 1991; this percentage is stable if we compare the 1991 and 1993 waves, and it decreases in subsequent years, down to approximately 0.3 percent in recent times. The analysis of individual cases shows that 3 out of 4 times the error concerns young children, of whom no features other than the basic demographics are recorded in the survey. The tendency for the misclassification rate to diminish with time is explained by the fact that from 1993 a greater effort was made to avoid discrepancies; the introduction of CAPI in 1998 fortified the attempt with automatic consistency controls. Birth dates of respondents also show a small number of misalignments, again decreasing with time. The province of birth varies in 2 percent of cases; a slight increase in misclassifications in recent years is probably due to the introduction of new provinces (Table 1).

Another feature that can be analyzed in order to gain insight on the reason why discrepancies arise is the type of high school diploma that respondents hold. Even if it only concerns a part of the panel sample (1,969 high school graduates), it is time-invariant: any reported difference can be safely labeled as an error. If we compare the 1998 and 2000 waves, we find that about 25 percent of the respondents report two different high school degrees. The transition matrix (Table 2) shows that almost 40 percent of inconsistencies arise between different types of trade schools, professional and technical. The *Technical school* category reveals the lowest reliability index $\lambda^* = 81.6$ percent; however, once the margins are taken into

TABLE 2
RELIABILITY OF TYPE OF HIGH SCHOOL DEGREE, 1998–2000 (PERCENTAGES)

2000							
1998	A	B	C	D	E	F	Total
A. School for professional studies	3.3	4.7	0.9	0.1	0.4	0.3	9.8
B. Technical school	5.3	41.0	1.5	0.3	1.2	0.7	50.0
C. High schools specialized in classical, scientific or language studies	0.4	1.9	16.1	0.2	0.7	0.1	19.5
D. Art schools and institutes	0.2	0.3	0.6	2.0	0.1	0.0	3.2
E. Teacher training school	0.6	1.1	1.0	0.0	11.6	0.0	14.3
F. Other	0.8	1.3	0.3	0.1	0.6	0.2	3.2
Total	10.5	50.4	20.4	2.7	14.6	1.4	100.0
Reliability index λ^* (consistent answers)	86.3	81.6	92.3	98.1	94.3	95.8	74.3
Adjusted reliability index λ^{**}	24.9	63.2	75.9	66.8	76.9	6.9	61.8

TABLE 3
LOCATION OF DWELLING OF RESIDENCE, 1998–2000 (PERCENTAGES)

2000							
1998	A	B	C	D	E	F	Total
A. Isolated area, countryside	3.4	0.9	1.4	0.5	0.2	0.0	6.4
B. Hamlet	0.8	3.0	2.0	0.5	0.1	0.0	6.3
C. Town outskirts	1.6	1.8	15.2	9.3	1.3	0.2	29.5
D. Between outskirts and town center	0.3	0.6	7.8	15.5	6.0	0.1	30.2
E. Town center	0.2	0.6	2.1	6.5	17.5	0.1	27.1
F. Other	0.1	0.0	0.2	0.2	0.1	0.0	0.6
Total	6.3	6.9	28.8	32.5	25.2	0.3	100.0
Reliability index λ^* (consistent answers)	94.1	92.8	72.1	68.3	82.7	99.1	54.6
Adjusted reliability index λ^{**}	50.4	41.6	32.5	26.4	55.2	0.0	38.7

account, the residual category *Other* ($\lambda^{**} = 6.9$) and the *School for professional studies* are the most unreliable items ($\lambda^{**} = 24.9$).

Response errors may be more frequent if the question itself is ambiguous or if the response options can be interpreted variously. For example, the answers given in 1998 and in 2000 on the location of the household's dwelling of residence (city center, suburbs, between the center and the suburbs etc.; Table 3) match only in 54.6 percent of the cases, probably because the classes are not precisely defined.¹⁸ The overall adjusted reliability index is $\lambda^{**} = 38.7$; the indexes referred to the single items reveal that the two extreme categories, *Isolated area, countryside* and *Town center*, are those with the highest adjusted reliability; the lowest indexes are found for the intermediate option *Between outskirts and town center* and for the residual category *Other*, which—considering its very low frequency—appears absolutely unreliable.

In the case of time-invariant continuous variables, the reliability index is based on the computation of the linear correlation coefficient of the answers given

¹⁸The comparison has been carried out only on households that did not move between 1998 and 2000. The "true" class of a dwelling is not necessarily time-invariant; changes, if any, should nevertheless only affect a small minority.

in two different waves ($\rho = \lambda^2$). For example, for the two measures of the floor area of the dwelling of residence, $\rho = 0.65$, and the reliability index is $\lambda = 0.80$. Note that the data only concern respondents who did not move or incur extraordinary renovation expenses between the two survey waves.

The reliability is lower for the construction year of the dwelling ($\lambda = 0.74$); in 73 percent of the cases, the spread is less than five years, but sometimes it is much greater, probably reflecting response difficulties for houses that have been heavily renovated.¹⁹

Other information that is affected by inconsistencies is the starting year of the respondent's working life. The usual recall problems are aggravated by a degree of ambiguity in the question: it is not clear whether occasional jobs or training periods should be included or not. Of 5,117 individuals who answered the question both in 1998 and 2000, 46.5 percent gave answers that do not match ($\lambda = 0.8$).

4.2. *Reliability of Time-Varying Quantities*

Table 4 presents the Heise reliability index for the main variables collected in the 1995, 1998 and 2000 waves of the SHIW.^{20,21} On a macro-aggregate level, net income and net wealth (Heise index: 0.82) seem to be more reliable than consumption (0.69).²²

The income components that show the highest reliability are pensions and wages; both Heise indexes are around 0.95. Fringe benefits such as the right to drive a company car, on the contrary, are not recorded as precisely (0.41): probably it is not easy to express their monetary value. Data on self-employment and capital income are collected with less precision (the Heise indexes are, respectively, 0.74 and 0.72). Serious problems arise with information on depreciation (0.48) and distributed profits (0.35). Expenditure on food seems to show greater reliability (0.80) than consumption as a whole.

The Heise indexes for wealth items are quite heterogeneous. While real estate is surveyed quite well (0.80, with 0.96 for primary housing), valuables do not perform as satisfactorily (0.47); it might be hard to state the value of objects that are not currently on the market, especially when the price of acquisition is also unknown because they were inherited or received as gifts.

¹⁹The existence of renovations is, unfortunately, documented for the year 2000 only; a correct comparison would require data for 1999 too.

²⁰The results presented in this section were obtained from the micro data of the historical archive, which includes imputed values. This implies that the reliability measure is referred to both collection and preliminary processing of information.

²¹The ranking of Heise indexes does not change even if we use, where necessary, the IV correction proposed for AR(2) processes. The direct application of IV methods for univariate regressions, when reasonably applicable (for example, when regressing consumption on income), also yields results that are aligned with Table 10.

²²In order to identify the variables for which the assumptions are more likely to be violated, Heise suggests the comparison of $\rho_{41} \rho_{32}$ and $\rho_{31} \rho_{42}$, which can be calculated if we have four waves; if the AR(1) models are a good approximation of reality, the two quantities should be very close. In the SHIW, they very often are; significant differences exist for valuables and, to a lesser extent, family-owned businesses. Where income components are concerned, the largest discrepancy emerges for distributed profits.

TABLE 4
HEISE RELIABILITY INDEX FOR THE MAIN SURVEY AGGREGATES, 1995–1998–2000

Aggregate	Heise Index	Aggregate	Heise Index
<i>Income</i>		<i>Consumption and savings</i>	
Net disposable income	0.82	Consumption	0.69
Payroll income	0.94	Non-durables	0.69
Net wages and salaries	0.95	Expenditure on food	0.80
Fringe benefits	0.41	Durables	0.27
Pensions and net transfers	0.94	Savings	0.61
Pensions and arrears	0.95		
Other transfers	0.76	<i>Other aggregates</i>	
Net income from self-employment	0.74	Stock of durables	0.43
Income from self-employment	0.79	Means of transport	0.89
Depreciation	0.48	Furniture	0.23
Distributed profits	0.35	Cash	0.57
Net income from capital	0.72		
Income from buildings	0.67	<i>Dwelling of residence</i>	
Income from financial assets	0.72	Owners	
		Surface area	0.84
		Value	0.84
<i>Wealth</i>		Construction year	0.78
Net wealth	0.82	Year of acquisition	0.83
Real wealth	0.79	Imputed rent	0.74
Real estate	0.86	Non-owners	
Dwelling of residence	0.90	Surface area	0.73
Family-owned businesses	0.56	Value	0.82
Valuables	0.47	Construction year	0.83
Financial wealth	0.68	Years of residence	0.96
Deposits	0.38	Rental rate	0.96
Government securities	0.74		
Other securities	0.64		
Debts	0.54		

The index for financial assets as a macro-aggregate is 0.68. Government securities appear to be measured better than deposits and other securities (respectively 0.74 vs. 0.38 and 0.64).²³ Government bonds are perceived as not exposed to market fluctuations, since most holders do not sell them before their maturity date; in contrast to shares and mutual funds, respondents normally declare the face value of the bond, which is easy to remember. Deposits are measured with lower precision because their high degree of liquidity may induce memory problems.

The measurement of debts appears to be quite unreliable (0.54). This applies to consumer durables as well (0.43), probably because the category encompasses many different types of goods, each of which induces different recall difficulties. The value of means of transport is an exception to this tendency (0.89), since information on the market value of used cars is widely available and known.

Finally, the value of primary housing is more reliable for the households that own than for those that rent; conversely, actual rental rates are measured with greater precision than imputed ones.

²³A high value of the reliability index does not exclude problems such as the bias deriving from under-reporting; the Heise coefficient does not change if households systematically withdraw information on a part of their assets.

Reliability indexes calculated for different sets of three waves (1989–1991–1993; 1991–1993–1995; 1993–1995–1998) are close to the ones presented in Table 4. Similar results are also obtained by estimating Heise indexes on the basis of Spearman rank correlation coefficients, which are not as strongly affected by the presence of outliers as are the Pearson coefficients. Reasonably, the indexes for relatively unreliable quantities are also less stable: the only exception is consumer durables, always showing very low indexes. Probably the AR(1) model is not a satisfactory formalization of the data generating process for durables, since they are bought or renovated irregularly.

4.3. *Reliability of Time-Varying Categorical Variables*

One of the most common uses of economic panel data is the analysis of transitions over time among states of categorical data (i.e. occupational status) or classes of quantitative variables (i.e. income quartiles). In what follows we are going to analyze the latter case, since it is an extensively studied topic in the social mobility literature.

It is well known that measurement error can bias the analysis of transitions. If respondents report data with errors, one will find units moving up and down even if their true state is unchanged; the observed transition probabilities (Table 5) are therefore likely to overestimate the mobility among income classes.²⁴

If we apply the LMM method to the measurements on three different occasions (1995, 1998 and 2000 waves) we can obtain the estimated response probabilities (assumed to be constants over time), i.e. the probability of each unit belonging to a class of being classified in each class (Table 5).²⁵ This matrix shows that the misclassification probabilities are lower for the extreme income classes. For instance, households in the fourth quartile have a probability of about 90 percent to be correctly classified while the risk of misclassification is significantly higher (about 20 percent) for the central classes. Once the problem of measurement error is taken into account, the estimated latent transition matrices show in both periods a significantly lower level of mobility than that observed, specifically in the second and third quartile.

5. EXPLANATORY MODELS FOR MEASUREMENT ERROR

The following paragraphs present models that aim to explain errors and inconsistencies in survey data on the basis of fieldwork, interviewer and

²⁴The variable of interest is the household total disposable income (income from payroll employment, from self-employment, from transfers and property income) net of income tax and social security contributions. At each point in time total disposable income is classified in four categories based on the quartiles of income distribution. The weights used in the analysis refer to 1995. The weights for the panel sample have been post-stratified in order to reproduce the main characteristics of the population at 1995 (age, town size and geographical area).

²⁵The fit of the model is satisfactory $\chi^2 = 32.7$ (p-value = 0.2) and $L^2 = 34.7$ (p-value = 0.15) with 27 df.

TABLE 5
HOUSEHOLDS' TRANSITION AMONG INCOME CLASSES: 1995–2000 (ROW PERCENTAGES)

Observed Transition Probabilities					
1998					
1995	First Quartile	Second Quartile	Third Quartile	Fourth Quartile	Total
First quartile	71,4	19,1	7,9	1,6	100,0
Second quartile	20,0	51,7	20,5	7,8	100,0
Third quartile	7,1	22,3	49,5	21,0	100,0
Fourth quartile	1,5	6,8	22,1	69,6	100,0
2000					
1998	First Quartile	Second Quartile	Third Quartile	Fourth Quartile	Total
First quartile	72,3	21,8	4,0	2,0	100,0
Second quartile	20,5	52,6	22,2	4,7	100,0
Third quartile	5,5	18,0	52,7	23,8	100,0
Fourth quartile	1,9	7,4	21,3	69,4	100,0
Estimated Response Probabilities(*)					
Observed Class					
Latent Class	First Quartile	Second Quartile	Third Quartile	Fourth Quartile	Total
First quartile	84,8	12,5	2,2	0,6	100,0
Second quartile	10,3	77,8	9,0	2,9	100,0
Third quartile	1,1	7,6	79,5	11,9	100,0
Fourth quartile	0,1	0,3	10,6	89,1	100,0
<i>Notes:</i>					
(*) Response probabilities are assumed to be time-invariant.					
Reliability index $\lambda^* = 82.7$.					
Adjusted reliability index $\lambda^{**} = 76.9$.					
Estimated Latent Transition Probabilities					
1998					
1995	First Quartile	Second Quartile	Third Quartile	Fourth Quartile	Total
First quartile	93,4	2,4	4,2	0,0	100,0
Second quartile	4,4	76,4	14,0	5,2	100,0
Third quartile	2,7	19,5	72,2	5,6	100,0
Fourth quartile	0,0	2,2	9,4	88,4	100,0
2000					
1998	First Quartile	Second Quartile	Third Quartile	Fourth Quartile	Total
First quartile	94,3	5,8	0,0	0,0	100,0
Second quartile	5,2	78,2	16,6	0,0	100,0
Third quartile	1,7	11,5	75,0	11,8	100,0
Fourth quartile	0,0	3,8	8,8	87,4	100,0

respondent features. Now, rather than seeking to quantify the incidence of measurement error, we want to find the reasons behind it.²⁶

Section 5.1 illustrates a model for the analysis of coding mistakes (e.g. wrong order of magnitude of a quantity); studying this problem is especially useful for an assessment of the interviewer's role in the determination of data quality, since the interviewer alone is responsible for erroneous coding.

During the preliminary editing phase, which precedes the production of statistics, the data undergo a number of quality controls. In many occasions, these controls lead to verification procedures which involve the examination of paper questionnaires (if available) or discussions with the interviewers. It would be theoretically possible to contact respondents again in order to remedy inconsistencies; but, as this is costly and time-consuming for the company in charge of the survey and the respondents alike, the actual strategy used is often different. If the discrepancies can be solved beyond reasonable doubt by looking at other sections of the questionnaire, the data are modified accordingly. This is typically the case, for instance, with real estate values expressed in millions of lire instead of thousands of lire; by comparing the declared worth with surface area and asset type it is easy to make the necessary correction. On the other hand, when the editing required is not so clear and the (presumed) inconsistency appears serious, households are re-contacted; if this is not possible, the answers are left as they are.

This approach reflects an interest in caution; anomalies are edited out of the database only if they can be safely assumed to be errors, and inconsistencies are rectified only when the values are certainly wrong, and they can be univocally replaced with correct ones. Such caution avoids forced "normalization" of micro data, i.e. replacement of information describing uncommon but true situations with numbers that portray standard occurrences. Researchers using survey information are left with the responsibility of deciding how to treat anomalies, based on the specific features of their analysis.

It seems evident from what has been said so far that a study of the preliminary editing process can shed light on measurement error issues, although there are known limits to the insight that can be obtained from such an exercise. The frequency of editing actions remains an imprecise indicator of the incidence of measurement error on the survey. As stated above, these actions are carried out only when an item can be safely considered wrong; some problems are therefore left undocumented. Since interviewer mistakes are easier to catch than mistakes by respondents, a study of the preliminary editing phase is more helpful in relating interviewer features to errors than in explaining why households represent their economic situation incorrectly.

Section 5.2 presents inconsistencies in panel data for some socio-demographic variables and for income, discussing the features that often accompany them. Differently from the analysis of the editing process, to do this we need two or more

²⁶A large part of the literature (e.g. Fabbris, 1989) claims that each interviewer induces an idiosyncratic distortion in answers, but the average bias is assumed to be null. If interviews were assigned casually, it would be possible to estimate the loss in precision caused by interviewers or by specific fieldwork features. We cannot assume casual allocation of assignments for the SHIW, because there is a strong correlation between the area in which an interviewer operates and respondent features. Moreover, this approach does not shed light on which features of data collectors actually affect the response variance; this is the reason why we prefer to study this problem with regression models.

waves. The 1998 and 2000 surveys have been selected, and discrepancies have been related to socio-demographic coordinates of respondents, interviewer features and fieldwork details.²⁷

It must be noted that information on interviewers is available for 2000 only. Since inconsistencies are generated by errors in either of the two waves, this specific lack in the data needs to be briefly discussed.

In a regression model, the omission of significant variables—such as interviewer features in 1998—introduces a bias in the estimated coefficients; the extent of the bias depends on the correlation between the variables omitted and those included. In this specific case, since the two waves were carried out by different companies, it is reasonable to assume that these correlations are equal to zero, controlling for localization.²⁸ The marginal effects of each variable included in the model are therefore estimated without bias, even though they do not attain minimum variance as they would if the model were specified exactly.

5.1. *Explanatory Models: Role of the Interviewer*

Keeping in mind the limits set out above, we now present some results concerning the preliminary editing process for the 2000 wave.

Two types of error are assessed. Firstly, we look at mistakes in the units of measurement; in some cases, answers are quite obviously given in millions of lire despite being requested in thousands of lire. Secondly, time-span errors are studied; in some cases, monthly incomes are declared when annual ones are requested and vice versa.

Issues related to units of measurement are mainly observed in the value of housing capital; time-span errors emerge on the income of employees and pensioners. These quantities are especially suitable for a preliminary study; the question on the value of the dwelling of residence is asked of every household in the sample, and a good share of the respondents are employees and pensioners, who are also not as likely to under-report as the self-employed, because they normally receive their income net of tax. The focus on these two types of income decreases the representativity of our analysis, but it also allows us to identify the errors in the data with ease and to explain them with simple models.

Table 6 gives some descriptive statistics on the editing actions taken on the three variables cited; a first striking fact is that error concentration is high. If the interviewers are ordered by number of errors found in their work, the last quartile appears responsible for a share of editing actions that ranges between 78.6 and 88.6 percent. The Gini coefficient for the number of errors ranges between 0.34 and 0.48.

The correlation coefficient between the percentage of errors and the number of interviews is always negative; the company in charge of the survey probably operates some form of control, giving more assignments to the better

²⁷We could not find significant effects of interviewer and fieldwork features on average income. It is therefore possible to explain answer variability on the basis of such features as they were in each wave, without having to use their variations as additional controls.

²⁸As illustrated above, the distribution of interviewer features is conditioned by localization. This induces positive correlation between features in different periods, but it is possible to eliminate its effects by including a geographical dummy in the regressions.

TABLE 6
 ERRORS ON SOME SURVEY VARIABLES, 2000 (UNITS, PERCENTAGES)

	Value of Dwelling of Residence	Income (employees)	Income (pensioners)
Errors	238	152	534
Records	8,001	6,553	6,175
Error rate	3.0	2.3	8.6
Interviewers with no errors	59.8	69.8	63.1
Interviewers with one error	22.0	15.3	14.0
Interviewers with more than one error	18.2	14.9	22.9
Error concentration	0.34	0.37	0.48
Share of errors accruing to interviewers with the worst performances (25 percent with highest number of errors)	78.6	88.7	86.8
Correlation between error rate and number of interviews	-0.12	-0.13	-0.003

interviewers.²⁹ This evidence also suggests that the skill of each interviewer is indeed observable; there is no *a priori* knowledge of how many corrections will be made on each individual interview, but the distribution of assignments seems to be consistent with *ex post* measures of precision.³⁰

Table 7 presents the results of a logistic regression run on edited and unedited records concerning the value of primary housing, the income of employees, and the income of pensioners. Only the first 60 interviews carried out by each interviewer are considered (interval of one standard error around the mean number of assignments), in order to limit the effect that last-minute interviews—typically assigned to top interviewers, hence not really representative—have on our estimates. Therefore, 6,467 questionnaires are considered (80.8 percent of the sample), corresponding to 722 edited records out of 924 (78.2 percent of the total).

The dependent variable is a dummy, set at 1 if the record has been corrected, 0 otherwise; the vector of independent variables encompasses the main features of interviewers and some fieldwork details, such as the use of CAPI and the length of the interview.

The respondent might have a role in determining mistakes of the type we are now studying, but this role appears to be junior high: mismatch in the measurement units, failure to clarify the time horizon relative to each question, lack of consistency checks are signals of incorrect interviewer behavior.

The results of our logistic regression are sufficiently stable. The probability of recording data that will subsequently require correction seems to be influenced by both interviewer and fieldwork features.

²⁹The interpretation of this information is not straightforward. While 59.8 percent of interviewers commit no mistakes, just 50.6 percent of the questionnaires require no editing actions. This indicates that the best interviewers typically carried out a number of assignments that are below the general average. But the number of questionnaires with a single error is 25.5 percent of the total (22.0 percent of the interviewers); 10.4 percent of interviews have to be edited twice (9.0 percent of interviewers).

³⁰Some reflections on efficient methods of interviewer selection and supervision can be found in Fowler (1991) and, with specific reference to new interviewing technologies, in Nicholls *et al.* (1997).

TABLE 7
PROBABILITY OF WRONGLY RECORDING AT LEAST ONE ANSWER, 2000 (LOGIT ESTIMATE)

	Coefficient ^o
Intercept	36,3734**
North	-0.0449
Center	-0.1954
Municipality: up to 20,000 residents	0.3212**
Municipality: between 20,000 and 40,000 residents	0.1976
Municipality: between 40,000 and 500,000 residents	0.2900
Paper questionnaire	-1.2401***
Interview length	-0.0403***
Interview length squared	0.0003***
Interviewer assessment of the general psychological climate during the interview ^{oo}	-0.1159***
Progressive number of the interview in the interviewer's portfolio	-0.0195***
Interviewer: previous SHIW waves	-0.0454
Interviewer: birth year	-0.0182**
Interviewer: male	-0.2056
Interviewer: junior high school degree	-0.4819
Interviewer: high school degree	-0.4948**
Interviewer: resident in a province different from the respondent	0.4083**
Interviewer: response rate	0.5137
Interviewer: non-professional	0.8381**
Non-professional interviewer: response rate	-1.2495**

Notes:

***Significant at the 1% level; **significant at the 5% level; *significant at the 10% level.

^oSignificance levels take into account intraclass correlation coefficients for each interviewer.

^{oo}Score expressed out of ten.

Predictably, professional interviewers are less inclined to go wrong; also, as the fieldwork progresses experience is gained and the frequency of errors diminishes. The non-professional interviewers who manage to obtain high response rates are also more precise.

The impact of interview length on accuracy is nonlinear: if too little or too much time is taken, data quality suffers. Normally, about one hour is needed to collect information. If the interviewers are excessively fast, they are probably not paying enough attention. Remarkable slowness might be a signal of interaction problems; moreover, fatigue can affect both interviewers and respondents if the interview exceeds average length.

Complementary results are obtained when studying the impact of the psychological climate in which the interview is carried out, as assessed by the interviewer; if tension arises, for example in the cases where some effort is needed to overcome reticence, more errors appear.

Where demographic traits of the interviewers are concerned, it seems that young males with high school degrees are less likely to be mistaken.

The risk of errors increases if the data collector is operating in a province different from his own. This may depend on a number of factors. First, "away" interviews tend to be shorter; interviewers are in a hurry and concentrate less. The psychological climate also tends to worsen.

On average, errors are less frequent when a paper questionnaire is used. It must be pointed out that the CAPI program is very efficient in monitoring the

interview flow; consistency controls, on the other hand, are limited, because it is not possible to exclude most unusual answers *a priori*. The software therefore only asks the user to confirm recorded anomalies, and it is only as accurate as the interviewers. Moreover, the CAPI interface proposes a sequence of screens and it is not convenient to go back to the previous ones to verify internal consistency; a paper questionnaire allows simultaneous vision of all answers, which is an advantage for quality control. The problem worsens if interviewers are not adequately trained using CAPI, as noted by Couper *et al.* (1997). This effect seems to offset the positive impact of automated verification mechanisms.

The superior quality of data collected via paper questionnaire does not hold for all variables; in particular, as we will see shortly, the CAPI method performs better on the items for which specific controls are implemented.

5.2. *Explanatory Models: Role of the Interviewer and Role of the Respondent*

In this section we present the results of some logistic regressions with the purpose of better understanding the characteristics of both respondents and interviewers most frequently associated to inconsistencies.

The first analysis concerns 9,473 individuals who declared their educational qualifications both in 1998 and 2000 (Table 8). The dependent variable is a dummy that has value 1 if an inconsistency is reported: a qualification that is higher in 1998 than in 2000, and a qualification that is higher by more than one level in 2000 than in 1998. The cases of one-level rise (e.g. primary/junior high school) were not labeled as erroneous. It is useful to point out that the estimates discussed here do not refer (as the previous ones did) to the probability of wrongly recording an answer; they concern the probability of finding inconsistent answers. The controls refer to what was stated by respondents in 1998, even if there is no reason to believe that in the presence of a discrepancy the true answer is the one given in any specific year. The replacement of 1998 values with 2000 values does not produce significant differences in the results. This implies a substantial symmetry in the probability of finding inconsistencies with respect to the wave chosen as portraying the “correct” educational qualification. In other words, the probability of observing discrepancies for a given class is approximately the same if the class is studied in 1998 or 2000.

In this case, like others that will be discussed below, respondent and fieldwork features are likely to be useful in explaining the presence of discrepancies; on the contrary, interviewers do not seem to play a crucial role.

Male respondents appear to be more consistent in their answers; elderly people seem to be more inclined to report two different qualifications. This can be explained with recall problems (Pearson *et al.*, 1992). Those who were born before 1955 also experienced an educational system different from the current one: they could choose between junior high school and apprenticeship, and possibly they have trouble reconciling their experience with one of the response options, which refer to the present organization of schools.

This particular circumstance helps to explain the concentration of inconsistencies on intermediate qualifications; those who had no formal education or completed primary school only tend to confirm in 2000 what they stated in 1998,

TABLE 8
PROBABILITY OF FINDING INCONSISTENT ANSWERS ON SOME PHENOMENA, 2000 (LOGIT ESTIMATES)

	Educational Qualifications Coefficient ^o	Type of High School Degree Coefficient ^o	First Year of Working Life Coefficient
Intercept	47.2256**	16.3527	24.3073***
Respondent: male	-0.2512*	-0.2450*	-0.1843***
Respondent: no formal education	-0.5377*	-	0.3370***
Respondent: primary school degree	-0.5637***	-	0.0559
Respondent: junior high school degree	0.1367	-	-0.2839***
Respondent: high school degree	0.6588***	-	-0.0869
Respondent: school for professional studies	-	1.1456***	-
Respondent: technical school	-	-1.0692***	-
Respondent: high school specialized in classical, scientific or language studies	-	-1.5465***	-
Respondent: art schools and institutes	-	-0.3923	-
Respondent: teacher training school	-	-1.2405***	-
Respondent: employee	-	-	-0.1304**
Respondent: self-employed	-	-	0.0618
Respondent: pensioner, former employee	-	-	-0.0315
Respondent: number of jobs held	-	-	-0.0075
Respondent: birth year	-0.0182***	-0.0001	-0.0091***
North	-0.1276	0.2935	-0.6519***
Center	-0.4366	0.4309**	-0.3123***
Municipality: up to 20,000 residents	-0.3626	0.0086	0.1777**
Municipality: between 20,000 and 40,000 residents	-0.3188	-0.0823	0.1969**
Municipality: between 40,000 and 500,000 residents	-0.2916	0.1400	0.6758***
Paper questionnaire	0.3465*	0.0953	0.3201***
Interview length	0.0005	0.0026	0.0014***
Personal interview ⁺	-0.2473**	-0.2370***	-0.1126***
Interview by proxy ⁺	0.2694*	-0.2253	-0.0308
Interviewer assessment of the general psychological climate during the interview ^{oo}	-0.1755*	-0.1644	0.1192***
Progressive number of the interview in the interviewer's portfolio	-0.0070	-0.0055	0.0025
Interviewer: previous SHIW waves	-0.0160	-0.0693**	-0.0080
Interviewer: birth year	-0.0671	-0.0074	-0.0029
Interviewer: male	-0.0070	-0.0138	-0.0270
Interviewer: junior high school degree	0.4248	-0.6487**	-0.3900***
Interviewer: high school degree	-0.0604	-0.1859	-0.2093***
Interviewer: resident in a province different from the respondent	0.2183	0.2922*	0.0742
Interviewer: response rate	-0.0359	0.1753	0.5041**
Interviewer: non-professional	-0.0319	0.0556	0.7525**
Non-professional interviewer: response rate	-	-	-0.8970**

Notes:

***Significant at the 1% level; **significant at the 5% level; *significant at the 10% level.

⁺The base class is "unknown."

^oSignificance levels take into account intraclass correlation coefficients for each interviewer.

^{oo}Score expressed out of ten.

while those who chose "junior high school" or "high school" are more exposed to confusion (the effect is statistically significant for high school only). Finally, even if the question clearly refers to the highest *attained* qualification, drop-outs at various levels may be uncertain in describing their situation.

If we look at fieldwork features and interviewer-respondent interaction, personal interviews are less exposed to discrepancies than the ones conducted by proxy. A good psychological climate reduces the chances of error; paper questionnaires are worse than CAPI, since the software actually points out inconsistencies for this particular question.

Previously outlined quality-improving factors, such as the experience gained in the course of a particular wave, as measured by the number of interviews already carried out, or the residence of the interviewer and the respondent in the same province, do not seem to be relevant in this case. The same goes for interviewer features, such as previous involvement in the SHIW: most coefficients have the expected sign, but they are not statistically significant.

The second logit regression was ran on the basis of answers provided by high school graduates interviewed in both 1998 and 2000 on their type of degree. The analysis confirms that the discrepancies depend heavily on the degree itself;³¹ those who reported graduation from a trade school in 1998 were more likely than the rest to change their answer in 2000. Degrees with a higher level of specificity seem to induce less confusion. Males turn out to be more consistent again; the remaining features do not seem to be significant. The general psychological climate is, again, correlated with a smaller probability of error. The interviewers who have a long record of SHIW waves, reside in the same province as the households surveyed, and hold a junior high school degree, are less inclined to record inconsistencies.

A further logit regression was run on inconsistencies in the reported year in which respondents started working; the analysis confirms the role of respondent features in determining the probability of discrepancies.

Once again, answers provided by males are more stable. The elderly face the usual recall problems; employees are better than the self-employed at remembering when they started working, probably because the concept itself is more formalized for them, and the initial date is more frequently recalled for reasons connected to wages, promotions and pensions.

Educational qualifications have an interesting effect on inconsistencies: those who are at the bottom and at the top of the qualification ladder are more exposed to errors. Where the unschooled are concerned, this can be explained by the usual difficulties in understanding the questions and interacting with the interviewers. University graduates, on the other hand, might have worked part-time while students, and they might be undecided as to whether they should consider these (often occasional) jobs as part of their working life or not.

It is also worth noting that in this case the geographical covariates, normally not significant, reveal more inconsistencies in the South and in small towns, possibly because the informal economy is more important there.

The interviewer-respondent interaction is again significant; interviews by proxy are less precise than personal ones; professional interviewers and the non-professional ones who obtain higher response rates are less exposed to error. A good interpersonal relationship between the interviewer and the family, as signaled

³¹A substantial symmetry in the distribution of inconsistencies exists, as stated for educational qualifications.

by the psychological climate in which the interview is carried out, shows positive effects. Excessive length of the assignment produces adverse consequences.

The CAPI technique reduces the probability of inconsistencies. Even if there is no specific control for this question, the initial working age is cross-examined with other variables, such as the year of birth, the number of years in which the respondent has paid pension contributions, and the year of retirement.

The effect of interviewer experience with the SHIW and residence in the same province as the respondent have the expected sign, but they are not statistically significant. Socio-demographic features of the interviewer, on the other hand, do not seem to explain inconsistencies, with the exception of educational qualifications: interviewers with junior high school degrees seem to perform better.

5.3. *Income Inconsistencies*

The investigation into the causes of discrepancies in incomes reported for 1998 and 2000 has been carried out with two different models.³² The first relates, via a linear regression, the absolute value of differences between the two waves to a number of controls that should catch “true” variations, and to the usual interviewer and fieldwork features. The coefficients yield a measure of the impact of each observed factor on the observed variability. The second is based on a logistic regression that models the probability of observing discrepancies greater than a fixed limit. In particular, assuming that the mistakes made by interviewers mostly appear in the tails of the distribution of differences, we create an indicator variable signaling percentage variations of income below the 5th and above the 95th percentile.³³

Given that the results are robust, we only present the logistic regression since its outcome is less exposed to the influence of outliers. Note that here, in contrast to our earlier procedure, income data is studied after the preliminary editing, and it is hence already devoid of blatant inconsistencies.

The experiment analyzed the differences between incomes reported by employees and pensioners in 1998 and 2000, corresponding to 3,244 households. As expected, a large part of the variability is related to socio-demographic features such as changes in the number of earners, gender, type of occupation, educational qualification, and area of residence (Table 9). Operational conditions and interviewer features do not seem to impact significantly on the discrepancies, except for the interviewer’s experience and the general psychological climate in which the interview is carried out, both of which show the expected sign.

Variability in reported incomes therefore seems to depend on causes external to the survey process itself; the very tendency towards under-reporting could be a further cause of additional variance, if the underestimation is not systematic.

It is worth noting that the results in Table 9 do not imply that the value of measurement error is correlated with individual characteristics. In fact, the phenomenon for which we model probabilities is symmetric: both negative and posi-

³²Since no information on interviewers is available for the 1998 wave, results have to be interpreted on the assumption of independence between interviewer features in the two periods.

³³In order to eliminate the effect of changes in household composition, only the households with the same roster have been studied.

TABLE 9
PROBABILITY OF FINDING EXTREME VARIATIONS ON INCOME (LOGIT ESTIMATE)

	Coefficient ^o
Intercept	-30.8911
Respondent: male	-0.3287
Respondent: birth year	0.0117
Respondent: number of earners in the household	0.1602
Respondent: new earners (employees) in the household	1.0382***
Respondent: new earners (self-employed) in the household	1.3093***
Respondent: new earners (transfers) in the household	0.8781***
Respondent: household wealth below general median	1.0952***
North	0.3960**
Center	-0.1302
Municipality: up to 20,000 residents	-0.3370
Municipality: between 20,000 and 40,000 residents	-0.3141
Municipality: between 40,000 and 500,000 residents	-0.0822
Respondent: no formal education	-0.9134
Respondent: primary school degree	-0.1442
Respondent: junior high school degree	-0.2048
Respondent: high school degree	0.2445
Respondent: employee	-0.6571***
Respondent: self-employed	0.9803***
Respondent: new head of household	0.3150
Paper questionnaire	0.0711
Interview length	0.0015
Interviewer assessment of the general psychological climate during the interview ^{oo}	-0.1230**
Progressive number of the interview in the interviewer's portfolio	-0.0032
Interviewer: previous SHIW waves	-0.1164***
Interviewer: birth year	0.0028
Interviewer: male	0.1594
Interviewer: junior high school degree	0.0395
Interviewer: high school degree	0.2877
Interviewer: resident in a province different from the respondent	0.1274
Interviewer: non-professional	0.2247
Interviewer: response rate	-0.1329

Notes:

***Significant at the 1% level; **significant at the 5% level; *significant at the 10% level.

^oSignificance levels take into account intraclass correlation coefficients for each interviewer.

^{oo}Score expressed out of ten.

tive extreme income variations are taken into account. Nothing can be said about the average effect resulting from those variations and its association with the covariates in the model. The only reasonable implication is that for given groups of individuals, such as the self-employed and people living in Northern areas, a greater variability in the error term may be expected.

6. CONCLUSIONS

We analyzed measurement errors affecting the most important variables of the Bank of Italy's Survey of Household Income and Wealth (SHIW). For time-invariant quantities, we evaluated the consistency of the answers given by panel units in various waves; for time-variant ones, such as income or wealth, we used the Heise (1969) model, which under mild regularity conditions can separate the actual change in a variable from measurement error on the basis of three or more

subsequent waves. We also examined the role played by fieldwork, interviewer and respondent features. Along with idiosyncratic elements referred to specific questions, there are a number of common explanatory factors for discrepancies.

The main results can be summarized as follows:

- Inconsistencies arise for all questions, even those that are neither ambiguous nor difficult to understand: in this case, discrepancies amount to 2 or 3 percent of the total. The number of errors decreases with time as greater attention is paid to avoiding them. Three fourths of the inconsistencies concern young children, surveyed only with respect to basic demographics.
- The number of inconsistencies increases when the question concerns information that might not be available to all family members, or is perceived as sensitive by the respondent, such as the type of high school degree or the level of educational qualification.
- When a question involves memory (e.g. the age at which a respondent started working) or when it does not specify how to treat certain situations (e.g. apprenticeship or occasional jobs), inconsistencies are the result of objective difficulties in determining the correct answer.
- Errors are more frequent when the response options are not precise enough (e.g. “center” or “suburbs” in the question about the location of one’s primary residence).
- The Heise reliability index, which measures the level of precision of the data (but does not catch systematic under-reporting), is higher for income and wealth (0.82) than for consumption (0.69).
- With regard to income, the data for employees and pensioners are the most reliable (around 0.95). Fringe benefits, on the contrary, are quite problematic, showing a Heise index of 0.41. Income from self-employment and capital are collected less precisely (respectively, 0.74 and 0.72).
- The consumption component that performs best is food (0.80).
- The Heise index for real estate wealth is 0.86; primary residence performs even better (0.90). Other wealth components, such as valuables (0.47), are more exposed to error, since it is not easy to evaluate items that are not currently on the market.
- Personal interviews contain fewer discrepancies than those conducted by proxy.
- When the CAPI software includes specific consistency controls, this helps to avoid discrepancies; when such controls are not present, the paper questionnaire is more precise. This is probably because it allows simultaneous view of the answers, whereas the electronic interface requires switching back and forth between screens. The problem is worse when the interviewer is not adequately trained in the use of the program.
- Interview length has an impact on accuracy; if it is too short or too long, data quality worsens. About one hour appears to be needed to complete the questionnaire. If the interview is much shorter, the interviewer probably did not pay enough attention; if it is much longer, fatigue may set in. Moreover, long interviews probably reflect a difficult interviewer–respondent interaction.

- Professional interviewers (and non-professionals who are better at obtaining high response rates) have better results in terms of data quality; previous experience with the SHIW has similar results.
- Experience gathered during a wave, as measured by the number of interviews already carried out, improves accuracy. The last assignments of each interviewer, in fact, are on average significantly better than the rest.
- The risk of errors increases when the interviewer works in a province other than that of residence. This can depend on several factors. Controlling for the number of family members and income earners, “away” interviews are shorter than “home” interviews; possibly, time constraints intervene.

These results are useful on three levels. Firstly, they allow the large number of researchers who use SHIW micro-data to properly take data quality into account when conducting their studies. This may extend to users of similar surveys, which are likely to be affected, at least partly, by the same issues. Secondly, our results can help data producers involved in surveys on income and wealth to implement quality-improving tools; the difficulties we described related to gathering and evaluating information are not specific to the SHIW. Knowing how to quantify problems with the data and how to identify their causes is essential in order to achieve improvement in survey procedures; for example, being aware of the existence and of the magnitude interviewer or questionnaire effects on specific items can lead to cost-effective quality-improving changes in interviewer training or questionnaire design. Finally, the conclusions we draw hopefully serve both as a reminder for data producers that quality-related information is important, and as a blueprint for quality reporting.

REFERENCES

- Banca, d'Italia, “Italian Household Budgets in 2000,” *Supplements to the Statistical Bulletin*, edited by G. D'Alessio, I. Faiella, new series, year XII, No. 6, January 2002.
- Biemer, P. and D. Trewin, “A Review of Measurement Error Effects on the Analysis of Survey Data,” in L. E. Lyberg, P. Biemer, M. Collins, E. D. De Leeuw, C. Dippo, N. Schwarz and D. Trewin (eds), *Survey Measurement and Process Quality*, Wiley, New York, 603–33, 1997.
- Brandolini, A., “The Distribution of Personal Income in Post-War Italy: Source Description, Data Quality, and the Time Pattern of Income Inequality,” *Giornale degli Economisti e Annali di Economia*, 58(2), 183–239, 1999.
- Cannari, L. and G. D'Alessio, “Mancate interviste e distorsione degli stimatori,” Banca d'Italia, Temi di discussione, No. 172, June 1992.
- , “Non-Reporting and Under-Reporting Behaviour in the Bank of Italy's Survey of Household Income and Wealth,” *Bulletin of the International Statistical Institute—Proceedings of the 49th ISI Session*, 55(3), 395–412, 1993.
- Cannari, L. and R. Violi, “Reporting Behaviour in the Bank of Italy's Survey of Italian Household Income and Wealth,” *Research on Economic Inequality*, Vol. VI, JAI Press, 117–30, 1995.
- Cannari, L., G. D'Alessio, G. Raimondi, and A. I. Rinaldi, “Le attività finanziarie delle famiglie italiane,” Banca d'Italia, Temi di discussione, No. 136, July 1990.
- Carroll, R. J., D. Ruppert, and L. A. Stefanski, *Measurement Error in Nonlinear Models: A Modern Perspective*, Chapman and Hall, London, 2006.
- Couper, M. P., S. E. Hansen, and S. A. Sadosky, “Evaluating Interviewer Performance in a CAPI Survey,” in L. Lyberg *et al.* (eds), *Survey Measurement and Process Quality*, Wiley, New York, 267–85, 1997.
- D'Alessio, G. and I. Faiella, “Non-response Behaviour in the Bank of Italy's Survey of Household Income and Wealth,” Banca d'Italia, Temi di discussione, No. 462, December 2002.
- Fabbris, L., *L'indagine campionaria*, La Nuova Italia, Firenze, 1989.

- Fowler, F. J., "Reducing Interviewer-Related Error Through Interviewer Training, Supervision, and Other Means," in P. Biemer, R. Groves, L. Lyberg, N. Mathiowetz, and S. Sudman (eds), *Measurement Error in Surveys*, Wiley, New York, 259–78, 1991.
- Groves, R. M. and M. P. Couper, *Nonresponse in Household Interview Surveys*, Wiley, New York, 269–93, 1998.
- Hand, D., H. Mannila, and P. Smyth, *Principles of Data Mining*, MIT Press, Cambridge, MA, 2001.
- Heise, D., "Separating Reliability and Stability in Test-Retest Correlation," *American Sociological Review*, 34(1), 93–101, 1969.
- Huber, P. J., *Robust Statistics*, Wiley, New York, 1981.
- Kish, L., *Survey Sampling*, Wiley, New York, 1995.
- Langeheine, R. and F. Van de Pol, "Discrete-time Mixed Markov Latent Class Models," in A. Dale and R. B. Davies (eds), *Analyzing Social and Political Change: A Casebook of Methods*, Sage Publications, London, 171–97, 1994.
- Lord, F. M. and M. R. Novick, *Statistical Theories of Mental Test Scores*, Addison-Wesley, Reading, MA, 1968.
- MacDonald, I. L. and W. Zucchini, *Hidden Markov Models and Other Types of Models for Discrete-valued Time Series*, Chapman & Hall, London, 1997.
- Nicholls, W., R. Baker, and J. Martin, "The Effect of New Data Collection Technologies on Survey Data," in L. Lyberg *et al.* (eds), *Survey Measurement and Process Quality*, Wiley, New York, 221–44, 1997.
- Pearson, R. W., M. Ross, and R. M. Dawes, "Personal Recall and the Limits of Retrospective Questions in Surveys," in J. M. Tanur (ed.), *Questions about Questions*, Russell Sage, 65–94, 1992.
- Vermunt, J. K., *Log-linear Models for Event Histories*, Sage, Thousand Oaks, CA, 1997.
- Wansbeek, T. and E. Meijer, *Measurement Error and Latent Variables in Econometrics*, Elsevier, Amsterdam, 2000.
- Wiggins, L. M., "Mathematical Methods for the Analysis of Multi-Way Panels," Unpublished doctoral dissertation, Columbia University, New York, 1955.
- , *Panel Analysis*, Elsevier, Amsterdam, 1973.