# Southern Africa Labour and Development Research Unit



## Genuine Fakes: The prevalence and implications of fieldworker fraud in a large South African survey

*by*
*Arden Finn and Vimal Ranchhod*

Working Papers can be downloaded in Adobe Acrobat format from www.saldru.uct.ac.za.
Printed copies of Working Papers are available for R15.00 each plus vat and postage charges.

Orders may be directed to:
The Administrative Officer, SALDRU, University of Cape Town, Private Bag, Rondebosch, 7701,
Tel: (021) 650 5696, Fax: (021) 650 5697, Email: brenda.adams@uct.ac.za

# Genuine Fakes: The prevalence and implications of fieldworker fraud in a large South African survey[*]

Arden Finn[†], Saldru, UCT
Vimal Ranchhod[‡], Saldru, UCT

**Abstract:** How prevalent is data fabrication by fieldworkers in South African surveys? Does this substantially affect the validity of subsequent empirical analyses? We document how we diagnosed such misbehaviour in the longitudinal National Income Dynamics Study. We found that the existence of fabrication was non-trivial, and affected about 7% of the sample. Since the fabrication was detected while fieldwork was still on-going, the relevant interviews were re-conducted and the fabricated data were replaced with authentic data. We thus have an observed counterfactual that allows us to measure how problematic such fabrication would have been, had it remained undetected. We implement a number of estimators using the dataset that includes the fabricated interviews, and compare these with the corresponding estimates that include the corrected data instead. We find that the fabrication would not have substantially affected our univariate and cross-sectional estimates, but would have led us to reach substantially different findings when implementing panel estimators. We conclude with a policy discussion for survey organizations that might have similar concerns.

# 1 Introduction

For anyone involved in the running of a survey, issues of data quality are of critical importance. Surveys can cost millions of dollars, require years of planning by large teams of people and need considerable levels of sustained effort. All of these resources are allocated for the sole purpose of producing high-quality data. All empirical findings, in turn, are premised on the assumption that the data being used are of a reasonable quality. This caveat applies to vast literatures in economics, sociology, demography and political science, amongst others. Indeed, it is so ubiquitous that it hardly ever gets stated explicitly.

In this paper, we investigate one aspect of the data production process that might lead us to question the quality of survey data. Most survey organizations, either directly or indirectly, employ fieldworkers to conduct their surveys. The fieldworkers, though, might not have the same objectives as the survey organization. These principal-agent problems might result in fieldworkers 'cheating'.[1]

Fieldworkers may engage in cheating behaviour for a variety of reasons. First, fieldworkers may be reluctant to ask sensitive questions about topics related to income, wealth or sexual behaviour. Second, some sections are very long and fieldworkers might want to leave them out in order to save time. Third, the characteristics of the primary sampling unit (PSU) may play a role. If the PSU is in an area that is considered dangerous (which is not uncommon in the South African context) or is very far away, then fieldworkers may end up cheating rather than visiting the PSU. Fourth, fieldworkers might be remunerated according to the number of successful interviews that they have completed. In the case of refusals, or the case where it is easier to fabricate an interview, this would incentivize cheating. Finally, the penalties for cheating may be small. If the survey company is unable or unwilling to monitor the behaviour of the fieldworkers, then the expected payoff to cheating might exceed the expected costs for some workers.

There are also different ways in which cheating behaviour could manifest itself. First, and most problematic, fieldworkers could fabricate entire interviews. In later waves of longitudinal studies, there is usually some pre-population of the questionnaires based on data from previous waves. This often includes a list of household members from the roster and their demographic characteristics. Fieldworkers can view this information and use it to form the basis of their fabrication. Fieldworkers could also cheat by leaving

---

[1]We use the word 'cheating' although in some cases a better word might be 'negligence'. The former implies intent whereas the latter could arise out of ignorance, incompetence or misunderstandings, and we cannot always separate between the two. In either case, fieldworkers do something that they ought not to have done which results in a deterioration of the aggregate quality of the data produced.

out sections of interviews. For example, in wave 1 of the National Income Dynamics Study (NIDS) questionnaire, the labour market section is substantial and has a total of 89 questions.[2] However, a respondent who is 'not economically active' will only answer seven simple 'yes/no' questions. Fieldworkers could save time by setting respondents' labour market statuses to 'not economically active' when they are, in fact, working or looking for work. A different way to save time would be to leave out certain people in the household. This would be easy to implement in a cross-sectional study. In a longitudinal study, the fieldworker might ignore new members in the household, such as babies or in-migrants, or exaggerate the number of people from the previous wave who have died. Our research findings presented in this paper are primarily concerned with the most problematic type of cheating listed, namely data fabrication of entire interviews.

The remainder of this paper is structured as follows. In section 2, we argue that the incidence of fieldworker cheating is a common problem in the implementation of large household surveys in several countries, and particularly in South Africa. In section 3 we turn our focus to the first two waves of the NIDS dataset and evaluate a number of methods that we considered to detect fieldworker cheating.[3] The two most successful methods, Benford's law and anthropometric diagnostics are dealt with in greater detail than the others. Section 4 describes the NIDS team's operational response to the discovery of cheating. Section 5 analyses what the consequences for future research would have been, had the cheating not been detected and corrected for. Section 6 offers recommendations for future fieldwork operations and provides some concluding remarks.

---

[2]This includes 10 sub-questions.

[3]This work was done by the authors while wave 2 of NIDS was still in the field. At the time, both authors were employed in the NIDS office.

## 2 The Prevalence of Fieldworker Cheating in South Africa

The phenomenon of fieldworkers making up data is a global and persistent problem with a sizeable literature dedicated to documenting and detecting it. A number of studies use data from major surveys in the United States to detect whether data fabrication had occurred. Schreiner et al. (1988) use Census Bureau Studies data from 1982 to 1987 to highlight the importance of reinterviewing respondents as a means of fraud detection. In their study, 83% of suspected falsifications turned out to indeed be a result of cheating. Most of the cheating was detected through reinterviews, although some were picked up because of anomalies in the data. In addition, most of the cheating involved total, rather than partial fabrication of individual-level data. The authors find that falsification rates range from 0.4% to 6.5% depending on which one of the Census Bureau surveys is used. Finally, they note that interviewers who had served for longer periods of time are significantly less likely to be data fabricators. Li et al. (2011) make the point that the Census Bureau's reinterview strategy for detecting falsification can be improved upon. The conventional reinterview methods detect falsification in less than 0.1% of the data. The authors use data from the Current Population Survey to try to design an alternative sampling method that should underlie the reinterview process. Using a combination of real data and simulations, they conclude that alternative sampling methods could find up to 20% more fabricated interviews that the current system. Murphy et al. (2004) use data anomalies from the National Survey on Drug Use and Health to identify suspicious fieldworker behaviour. In particular, they flag relatively short or long interview durations as possible signs of falsification and show how taking these durations into account adds to the power of the fraud detection process.

Outside of the US, Schäfer et al. (2004) use data from the German Socio-economic Panel (SOEP) to test the reliability of two methods of fraud detection. Data fabrication was low in all waves and all samples of the SOEP, never exceeding 2.4% of all cases, with the overall share of faked data at about 0.5% (Schräpler and Wagner, 2005). The authors use the fabricated data that was removed from the publicly-released version of the SOEP and find that using Benford's law as the basis for detecting suspicious data correctly identifies all cases of fabrication. In addition, they exploit the fact that cheating interviewers tend to have less variability in their responses over all questions and all interviews than non-cheaters. Interviewer-level tests for surprisingly low variance also correctly identified all of the cases of cheating. All confirmed cheating interviewers were middle-aged and male, and the effect of education on the probability of cheating was not statistically significant.

There are a number of other studies that use characteristics in the data themselves as a means of identifying fabrication. These include, among others, Bredl et al. (2008) in an unspecified non-OECD country, and Porras and English (2004), Cho et al. (2003) and Swanson et al. (2003) in the US. A broad review on much of the literature related to the detection of fabricated data can be found in Birnbaum (2012) who charts the methods used in twelve different datasets in the developed and developing world.

We now narrow our focus somewhat and turn our attention to illustrative cases of fieldworker cheating in the South African context.

## 2.1 KwaZulu-Natal Income Dynamics Study (KIDS)

KIDS is a household level panel dataset that was conducted in 1993, 1998 and 2004. It revisited a subset of the households located in the KwaZulu-Natal province of South Africa that were included in the original SALDRU/PSLSD 1993 survey. Follow-up fieldwork in May of 2001 suggested that there may have been cheating by fieldworkers in some clusters. Subsequent investigations revealed that the fabrication was limited to two clusters and these were permanently removed from the sample.[4]

## 2.2 Survey on time and risk preferences

Between 2010 and 2011, researchers from the University of Cape Town conducted a survey on time and risk preferences in the three major metropolitan regions of South Africa,[5] with a budget of about 300 000 US dollars.[6] They had a sample size of about 300 respondents and visited each of them six times at three monthly intervals. The survey included a background questionnaire as well as two experimental modules. In the experimental modules, respondents were asked to choose between various alternatives in order to ascertain their appetite for risk and their discount factors. In order to obtain truthful responses, all choices were incentivized to have some probability of entailing an actual cash payout.

In the time preferences component, respondents answered 40 questions. They then rolled a 10-sided die, and if it landed on 0, they would get paid for one of their 40 responses. The relevant question was selected by rolling a 10-sided die and a 4-sided die simultaneously. In the risk preferences component, respondents were also asked 40

---

[4]See May et al. (2007) with an earlier version available for download at
`http://www.datafirst.uct.ac.za/catalogue3/index.php/catalog/286`.

[5]These are Johannesburg/Pretoria, Cape Town and Durban.

[6]Information on the details of this study was obtained through interviews with Andre Hofmeyr. At the time, he was a researcher actively involved in the survey.

questions, one of which would yield a cash payout with certainty. The relevant question was also selected by means of rolling a 10-sided die and a 4-sided die simultaneously. The payouts varied by question and by the choices made by the respondents in that question. The fieldworker would then pay the amount of the winnings in cash to the respondents.

After the data were collected, researchers found a suspiciously high rate of interviewees getting paid out for the time preferences component. The *ex ante* probability of this occurring was 10% but the respondents were 'winning' 25% of the time overall. Moreover, respondents were observed to have a disproportionately high probability of having 'randomly selected' questions with relatively higher cash payouts in both the risk preferences and time preferences component of the study. Further investigation indicated that these anomalies were driven by data from a subset of fieldworkers who almost always paid out the maximum amounts permissible. People involved in the study believe that some fieldworkers colluded with respondents so as maximize the actual disbursements, which they could then share. The problem was identified only after the 4th wave of data had been obtained, with the 5th wave already in the field, and both the time preferences and risk preferences components of the study had to be abandoned.

## 2.3 Cape Area Panel Study: Wave 5

The Cape Area Panel Study (CAPS) is a longitudinal study of young adults in the Cape Town metropolitan area. Wave 1 was conducted in 2002 with a sample of about 4 800 young adults aged 14 to 22. In the fifth wave of CAPS, conducted in 2009, part of the interview included a finger-prick test for HIV status which was administered by the fieldworker.[7] The *ex ante* expectation was that about 30% of the women interviewed would be HIV positive. For most fieldworkers, the proportion of HIV-positive female respondents was around this figure, but after a certain date, one fieldworker returned HIV-positive results for every respondent.

It took a considerable amount of the time for the lab results on the blood to be returned to the operational headquarters. Thus, by the time that this was discovered, the fieldworker in question had already been paid and had left the survey. Investigations discovered that the fieldworker in question had not, in fact, taken blood samples from respondents, but had obtained blood samples from some other source. The result was

---

[7]The information on CAPS was obtained through interviews with Dr. Brendan Maughan-Brown, co-ordinator of the fifth wave of CAPS. More information can be obtained in Lam et al. (2012) which can be downloaded at
`http://www.datafirst.uct.ac.za/catalogue3/index.php/catalog/266`.

that all data collected by this fieldworker after a certain date was deleted and did not form part of the fifth wave.

A common method of monitoring fieldworker behaviour is to phone respondents in the weeks or months after the interview in order to verify that they were in fact interviewed. One of the fieldworkers obtained the list of verification questions and set up a system in which her sister-in-law pretended to be a respondent each time she was called by the survey company. This suggests that fieldworkers who do cheat can use quite sophisticated methods to avoid detection. This fieldworker's cheating was only discovered after the conclusion of fieldwork, and all relevant data was deleted from the study.

Overall, a total of 8 fieldworkers had engaged in some form of cheating, out of an average of about 40 fieldworkers over the course of the fieldwork. A total of 289 fraudulent interviews were deleted from the public release version, which represented about 9% of the expected sample at the start of wave 5.

## 2.4 Time Use Study

In 2000, StatsSA, the official statistical agency of South Africa, conducted a national time use study over three different months in order to investigate how South Africans spend their day. The total sample size was approximately 14 000. Household members were eligible to participate if they were aged 10 or older. Fieldworkers were asked to fill in a household roster in descending order of age, and if more than 2 household members were eligible, to select two household members sequentially using a sampling grid.[8] If the fieldworker reached the end of the grid, that is at the $11^{th}$ such household, then she was instructed to 'loop' back to the start - that is, to treat the $11^{th}$ household with 4 eligible people as if it were the $1^{st}$ household with 4 eligible people that she had encountered (Statistics South Africa, 2000).

The sampling grid yields an asymptotic distribution of the frequency with which household members of a particular age-rank ought to have been selected on aggregate, conditional on the number of eligible persons. For example, in households with three eligible persons, we would expect to find that person 1 was selected 50% of the time, person 2 was selected 70% of the time and person 3 was selected 80% of the time.[9] In the dataset, however, we find that persons 1, 2 and 3 in households with three eligible

---

[8] A copy of the sampling grid is included in Appendix A as Table 11. To illustrate how it works, suppose that a fieldworker came to her first house with four eligible members. She should then select persons 2 and 4, i.e. the $2^{nd}$ and $4^{th}$ oldest members of the household. In the second such household, she should select persons 1 and 3, etc.

[9] The total adds up to 200% since two household members were selected.

persons were in fact selected 81%, 81% and 38% of the time, respectively. A Chi-squared test rejects the null hypothesis that the realised distribution corresponds to the expected distribution at any reasonable level of significance.[10] We repeated the analysis for households with 4, 5 and 6 eligible members respectively, and convincingly rejected the null hypothesis of equivalence of distributions in each case. We interpret this as evidence that fieldworkers did not, in fact, follow instructions about whom to select in households where they had some degree of choice.

An alternative explanation to fieldworker cheating is that the asymptotic distribution is not the correct distribution to use, since we would only expect it to be realised if each fieldworker had 'many' households with more than two eligible persons. Nonetheless, at least in the case with three eligible persons, there are more than 1 000 such households in total. If each fieldworker encountered 10 or more such households then the asymptotic distribution would provide a reasonable approximation of the correct expected distribution. Moreover, the magnitude of the divergence is so great that, unless the expected distribution that we use is grossly incorrect, we would continue to reject the null hypothesis that the realised distribution corresponds to the 'true' expected distribution.

Note that in this case we are not claiming evidence of data fabrication. The 'cheating' that we observe here is of a very different nature to those previously documented. What explains such cheating? We conjecture that fieldworkers violated the sampling instructions due to some combination of the availability of respondents as well as variations in the time that it would take for different respondents to complete the questionnaire.

Two empirical observations support this conjecture. First, in households where the fieldworker had some choice, 53.7% of those eligible were female, but 56% of those selected were female. This difference is small in terms of percentage points but since it applies to just over 9 000 observations it is statistically significant. Moreover, there is nothing in the sampling grid that suggests a clear gender bias in terms of who ought to be selected. A more likely explanation is that females are more likely to be available for an interview, since they are much less likely to be employed in South Africa.[11] Second, we observe that in households with 3 or more eligible persons, 51% of those eligible are younger than 21 years of age. Of these potential respondents, only 35% were selected to fill in a questionnaire. Teenagers are probably less likely to be available due to being in school. Additionally, they might be less willing to participate in interviews in the first place, and it might take longer for them to complete a questionnaire.

---

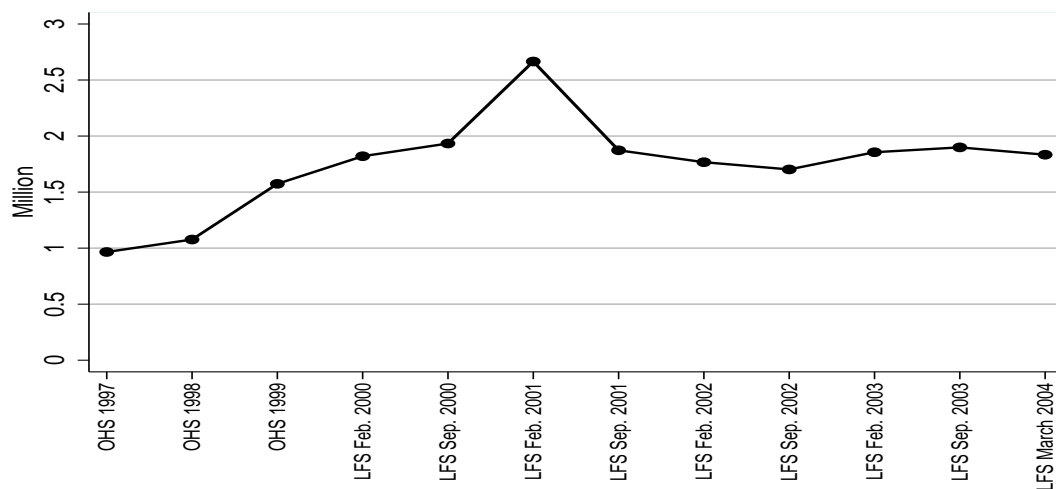[10]The results of this analysis are presented in Table 12 in Appendix A.

[11]In StatsSA's September 2000 Labour Force Survey (LFS), amongst prime working aged adults aged 21 to 59, the employment rates for men and women were 61.6% and 45.5% respectively. This was a nationally representative survey with a sample size of about 50 000 working-age adults.

If our conjecture is correct, then the violation of the sampling framework has potentially serious implications for analyses. Obtaining a disproportionate number of unemployed or 'not economically active' people in the sample will bias measures of aggregate time use, and conventional sampling weights, even if adjusted for non-response, will not correct this bias.

## 2.5 Labour Force Survey 2001

Devey et al. (2006) contains an interesting figure, reproduced below as Figure 1, that charts the number of people in South Africa classified as 'informally employed'. The authors use data from October Household Surveys (OHS) and Labour Force Surveys (LFS) from 1997 to 2004. The most striking feature of the data is the spike in February 2001, with a jump of almost 750 000 informal workers, before falling by approximately the same number to September 2001. It is implausible that such a spike should be present in nationally representative data with a consistent survey instrument. The authors point out that in the February 2001 LFS, interviewers were offered additional incentives to interview informal workers in the households they visited. Although there is certainly no established claim for cheating having taken place whereby interviewers fabricated data on the informal economy, it is nevertheless very suspicious to see the spike when finding informal sector workers was incentivized, and an immediate and equivalent reversal when the incentive was removed.

Figure 1: Informal Employment (in Millions) from OHSs and LFSs 1997 – 2004



Source: Reproduced from Devey et al. (2006).

9

In summation, we have provided evidence of fieldworker cheating in four substantial South African surveys and highlighted potential cheating in a fifth. These surveys have been both cross-sectional and longitudinal, and span a time period from 1993 to 2011. The implementing agencies include local and international fieldwork companies, as well as organisations that employed and managed fieldworkers directly. The cheating manifests in various forms, including outright fabrication of entire interviews, falsification of responses to particular questions and not following the sampling instructions. In some cases the cheating did not affect the overall legitimacy of the study, whereas in one case an entire component of the study needed to be abandoned. The research areas that have potentially been affected include time use, health, risk preferences, labour market status, poverty, inequality and inter-generational mobility. Thus the incidence of fieldworker cheating is widespread and its effects are potentially far reaching.

In the next section, we discuss how we attempted to identify possible cheating in wave 2 of NIDS.

# 3 Fieldworker Cheating in Wave 2 of the National Income Dynamics Study

The second wave of NIDS is the main focus of this paper. NIDS is a nationally representative longitudinal study that collects data from respondents on many socio-economic topics including education, labour market participation, fertility, mortality, migration, income, expenditure and anthropometric measures. The survey starts with a household roster that documents all people who are resident in the household at the time of the interview. From the information captured in the roster, respondents are classified as either children (aged 0–14) or adults (aged 15 and above). Fieldworkers are instructed to then administer either a 'child questionnaire' or an 'adult questionnaire' for each household resident. In cases where the respondent refused or was not available, fieldworkers were asked to try to get a knowledgeable person in the household to fill in a 'proxy questionnaire' on behalf of the respondent who could not be interviewed.

The first wave, which took place in 2008, had a sample size of 7 301 households and about 17 000 people completed the adult questionnaire. In that wave, fieldworkers used paper-based questionnaires and entered responses by hand. The completed questionnaires were then sent from the fieldwork company to a data capturing company and, by the time the full dataset was received by the survey operations team, fieldwork had already been completed. The primary data quality control procedures thus occurred after the fieldwork had been completed in wave 1.

The second wave of NIDS was conducted over 2010 and 2011. Fieldworkers used a Computer Assisted Personal Interview (CAPI) system, whereby fieldworkers filled in responses on a hand-held computer. Data from completed questionnaires were then uploaded to a server on a daily basis. One of the advantages of having data come in 'live' was that a verification process was undertaken while the fieldworkers were still in the field. This allowed for corrections to be made as part of the ongoing fieldwork operations, so that suspicious data could be verified or replaced, rather than deleted.

Our objective for the verification process was to create a measure that could rank fieldworkers by decreasing levels of suspiciousness. Once fieldworkers were ranked, the respondents that they interviewed were called back to ascertain whether they had been interviewed or not, and if they had indeed been interviewed, whether the entire questionnaire had been completed.

In creating the suspicion-based ranking of fieldworkers we considered using nine different methods. The central idea in using each of these possibilities was that fieldworkers who do cheat will do so either to save time, or to earn more money, or both. Fieldworkers

could earn more money as they received a performance-based incentive for each completed individual and household questionnaire , as well as for successfully completing an entire household. This would result in systematic differences in some dimensions of the data that were generated by cheating fieldworkers, when compared either to non-cheating fieldworkers or to externally-motivated benchmarks.

The most successful of these nine methods were the use of Benfords law and anthropometric comparisons, which we discuss in detail in sections 3.2 and 3.3 respectively. Although the other seven methods were not particularly useful for diagnostic purposes, it might nevertheless be worthwhile to document what did not work.

## 3.1 Unsuccessful methods of detection

### Method 1: Number of deaths between waves

One way to speed up the process of completing an entire household would be to falsely classify a household member from wave 1 as deceased. This would allow a fieldworker to 'complete' interviewing the household much faster. Alternatively, fieldworkers could falsely classify someone who had died between waves as being still alive, and then fabricate the data. We compared the mortality rates of respondents by fieldworker, and did not observe any anomalies in the data.

### Method 2: Number of refusals/not available

The method and thought behind using this metric is identical to that for using deaths above. Fieldworkers could save time by not interviewing everyone (for example, by fabricating refusals and non-responses from respondents), or fabricate data for those who had, in fact, refused to be interviewed. We compared the response rates by fieldworker, and did not observe any anomalies in the data.

### Method 3: Fieldworkers who are disproportionately likely to activate substantial skip patterns in the survey

Our thoughts here were that one could save considerable time in some sections by capturing certain responses. As already discussed, this incentive is strongest in the labour market section. We abandoned this method as the levels of unemployment in South Africa are high, levels of labour force participation are low, and these are concentrated in certain neighbourhoods and regions (Leibbrandt et al., 2010). Since fieldwork was co-ordinated geographically, fieldworkers could plausibly have genuinely encountered a

pool of respondents with low levels of employment and labour force attachment in their allotted region. In addition, the unemployment rates and percentage that were not economically active, by fieldworker, yielded several fieldworkers with high values, so this was not a particularly useful tool for discriminating between suspicious and non-suspicious fieldworkers.

## Method 4: Using length of interviews to identify fabrication

If fieldworkers were fabricating data, we expected them to complete the surveys relatively quickly. The software we used had time stamps for both when the interview began and was completed, which theoretically allowed us to calculate time per interview. We also expected that each adult interview would take between 45 minutes and one hour to complete. Unfortunately, the time stamp for completion was activated manually, and several fieldworkers only did so at night prior to uploading the data to the server. This rendered this component of the investigation useless.

## Method 5: Using GPS co-ordinates to verify where the interview took place

Part of the survey captures the GPS co-ordinates of the household. This was required for all households in both wave 1 and wave 2. The co-ordinates were obtained by means of a handheld GPS device which was accurate to a radius of 100m. If interviews were being fabricated, we would expect to find differences between the wave 1 and wave 2 co-ordinates. We encountered two problems with this method. First, there was considerable measurement error in wave 1, so not all differences could be attributed to wave 2 cheating. Second, fieldworkers were given the GPS co-ordinates from wave 1 to assist them in finding the households. Instead of entering the GPS readings from the GPS device in wave 2, a cheating fieldworker could simply re-enter the co-ordinates that they had received.

## Method 6: Comparing wave 1 and wave 2 signatures

Each completed questionnaire in each wave needed to be accompanied by a signed paper-based consent form. We considered comparing wave 1 and wave 2 signatures to identify discrepancies. We abandoned this approach very quickly, as signatures have some variability over time, and the method was far too labour intensive.

**Method 7: Low rates of in-migration or births between waves**

If fieldworkers were fabricating entire households, then they would not be able to know about any new household members that entered between waves. They would then either systematically under-estimate the number of new members, or else have to fabricate these new members as well. We calculated the number of new members per household by fieldworker, but there were no clear patterns or anomalies. If cheating fieldworkers did indeed fabricate new members as well, or if cheating fieldworkers cheated only on some fraction of their households, then it would be much harder for this diagnostic to yield usable information.

## 3.2 Using Benford's Law

In contrast to the methods described above, the use of Benford's law as a ranking mechanism for suspicious fieldworkers proved to be very useful. Following a paper by Schäfer et al. (2004), we used Benford's law as the basis of a test of the distribution of the numerical data reported by each fieldworker. Benford's law is an empirical law that was first described in Benford (1938). It describes the probability distribution of leading digits in tables of numerical data and asserts that the distribution is not uniform, as might be expected *a priori*, but rather follows a certain logarithmic probability distribution given by:

$$Pr(\text{leading digit} = d) = log_{10}\left(1 + \frac{1}{d}\right), \quad d = 1, 2, ..., 9$$

This implies that the probability of a leading digit being 1 is about 30%, the probability of it being 2 is about 17.6%, with the corresponding probabilities of the subsequent leading digits decreasing monotonically until we find that the probability of the leading digit being 9 is approximately 4.6%. The probability distribution of leading digits is shown in Table 1, below.

Table 1: Benford's Distribution of Leading Digits

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|
| 30.1% | 17.6% | 12.5% | 9.7% | 7.9% | 6.7% | 5.8% | 5.1% | 4.6% |

For a long time this phenomenon was viewed as not much more than a numerical curiosity. However, some practical implications began to emerge (Scott and Fasli, 2001;

Durtschi et al., 2004), and Benford's law has since been used to detect fraud in financial statements of companies (Carslaw, 1988; Thomas, 1989). More recently, it has been used in a wide variety of settings in the US (Durtschi et al., 2004).The law has also been found to hold with a large number of other kinds of data, including the population of towns, the length of rivers and the half-life of radioactive atoms. The basic premise is this: If you have a relatively large dataset and you accept that Benford's law holds, then you can identify possible cheating by comparing the realized distribution of leading digits for each fieldworker, to the distribution of leading digits that would be expected if Benford's law holds.
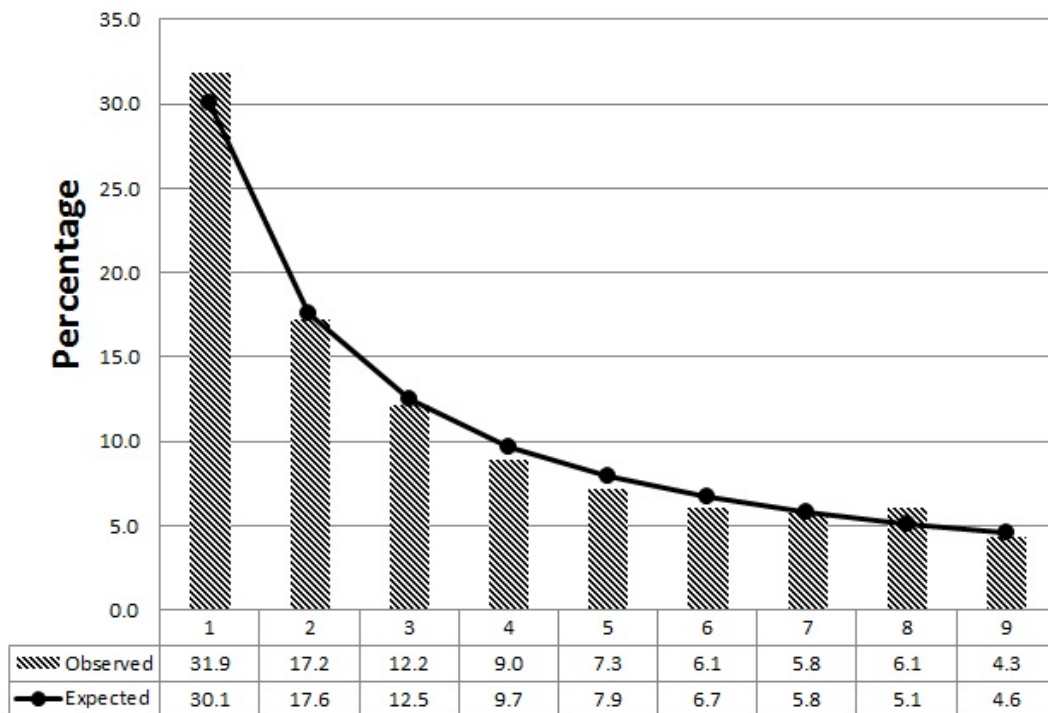
Hill (1995) provided the first theoretical basis for the law and showed that the law applies most accurately to stock market data, some accounting data and census statistics. The intuition underlying the proof is the following: Consider a variable which grows at some constant rate. Regardless of the initial value or the growth rate, the asymptotic distribution of leading digits of this variable (over time) will conform to Benford's law. Thus, a random sample of such variables at a moment in time will also conform to Benford's law.

More recently, Schäfer et al. (2004) and Schräpler (2011) argue that certain survey data also conform to Benford's law, and use this for the express purpose of identifying cheating fieldworkers in the German Socio-economic Panel (SOEP). Schräpler (2011) summarises three requirements that need to be satisfied in order for Benford's law to be a useful diagnostic for detecting fraud in survey data. First, the data should not have a built-in maximum value. Second, there should be no externally assigned values in data. For example, the South African old age pension is a rand amount that is assigned to an individual, and this is an example of data that cannot be used in the diagnostic test. Finally, the distribution of the data should be positively skewed with a median that is lower than the mean. Of all the variables in wave 2 of NIDS, the wage, income and expenditure data satisfied all of these criteria, and we utilized these variables to implement the test.

Within the broad categories of variables that satisfied all of the requirements listed above, there are many candidate variables. These range from labour market income to household expenditure on various goods and services to a self-reported estimate of total household income. Figure 2, below, plots the distribution of leading digits of the variable reflecting the total amount of labour market income received by respondents in the 4 705 households with positive wages in the wave 1 data. This distribution, shown by the bars, is plotted together with the logarithmic distribution that we expect to observe, assuming that Benford's law holds. By comparing the two distributions, it

15

seems that the leading digits of this variable fit the logarithmic distribution very well. The observed proportions of each leading digit are very close to the proportions that we expect to observe *ex ante*, and fall with each successive digit, except for eight which is slightly higher than expected.[12]

Figure 2: Observed and Expected Leading Digits - Wages in Wave 1



| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| Observed | 31.9 | 17.2 | 12.2 | 9.0 | 7.3 | 6.1 | 5.8 | 6.1 | 4.3 |
| Expected | 30.1 | 17.6 | 12.5 | 9.7 | 7.9 | 6.7 | 5.8 | 5.1 | 4.6 |

Source: Own calculations using NIDS Wave 1 2008.

Given that the aggregate data seem to follow this logarithmic distribution for some of the monetary variables, we next sorted the wave 2 observations by fieldworker and considered the conditional leading digits by fieldworker. We ranked how far each fieldworker's distribution of leading digits was from the logarithmic distribution, by computing Chi-squared statistics. Ranking fieldworkers by this method yielded positive results for the detection of cheating. Of the fieldworkers with the five highest Chi-squared values, four were subsequently found to have fabricated entire questionnaires. The top five Chi-squared rankings are shown in Table 2, below. The cheating fieldworkers are highlighted in bold.

The fact that four of the top five most suspicious fieldworkers were identified using

---

[12]See Figure 6 in Appendix B for a comparison of observed versus expected distributions of leading digits for some other wave 1 and wave 2 monetary variables.

Table 2: Most Suspicious Fieldworkers by Chi-squared Ranking

| Ranking | Fieldworker | Chi-squared (no. of interviews) |
|---------|-------------|--------------------------------|
| **1** | **A** | **39.7 (80)** |
| **2** | **B** | **31.2 (49)** |
| 3 | OK | 28.0 (66) |
| **4** | **C** | **27.3 (74)** |
| **5** | **D** | **27.1 (100)** |

Source: Own calculations using pre-public release NIDS Wave 2 data, 2010.

Benford's law suggests that using this method is appropriate for survey data of this nature. Nonetheless, some cheating fieldworkers may have left the monetary variables blank, or set them to missing. In this case we would have the unique problem of missing fake data, which also happens to be fake missing data. We thus exploited the longitudinal dimension of the data and evaluated variables that are difficult to fabricate convincingly in a panel study – namely, height and weight.

### 3.3 Anthropometric Measures

The first two waves of NIDS included modules in which respondents were weighed and measured for height. Weights were obtained using digital scales that were accurate to 0.2kg, while heights were obtained using a portable stadiometer. These data were not pre-populated into the CAPI system, making it almost impossible for fieldworkers to systematically fabricate values that were consistent with wave 1 measures for respondents that they had not seen.

Various diagnostic measures were used to rank fieldworkers according to their likelihood of having cheated. These were:

- **Mean adult body mass index (BMI), by fieldworker.**[13] Our thinking was that fieldworkers who were fabricating data might not be aware of the extent to which the height and weight of respondents are correlated. This would result in 'abnormal' BMI measures. We considered fieldworkers who generated exceptionally high or exceptionally low average BMI values as potentially suspicious.

---

[13]BMI is calculated by dividing a person's mass (in kilograms) by the square of their height (in meters). A BMI above 25 is considered to be overweight by the medical profession.

- **The mean growth in the height of adult respondents between waves, by fieldworker**. We expected that the heights of adults would be stable over a two year span, on average. If the mean growth in height for a particular fieldworker, differed substantially from zero in either direction then we interpreted this as an indication of possible cheating.

- **The mean BMI growth from wave 1 to wave 2, by fieldworker.** If this differed substantially from zero then this was a sign of possible fabrication.

- **Spikeplots of the weight distribution, by fieldworker**. Given that the scales were digital, and that no fieldworker had interviewed hundreds of respondents, we expected to obtain a uniform distribution of heights. Visual inspection of the spikeplots allowed us to relatively quickly identify suspicious patterns such as 'heaping' at natural reference numbers.

The diagnostics above were restricted to adults only, where adults were classified as respondents 20 years old and above. Running the diagnostics on children would have presented a problem as the height and weight variables for children are more volatile, even in a two to three year period.

Table 3, below, shows the list of suspicious fieldworkers generated using the mean adult BMI in the wave 2 cross-section. As before, fieldworkers who were found to have cheated are highlighted in bold. E, who interviewed 97 adults, had the highest mean BMI of 55.3. The highest mean BMI measure that we verified *ex post* was 43.3, but was based on only 20 respondents. At the other end of the distribution, H returned a very low mean BMI of 21.7, far lower than the overall average of 28.6 from 9 821 adults.

Table 3: Suspicious Fieldworkers and Mean Adult BMI

| Fieldworker | N | Mean BMI |
|---|---|---|
| **E** | **97** | **55.3** |
| **F** | **24** | **49.6** |
| **G** | **49** | **48.9** |
| **B** | **104** | **44.7** |
| OK | 20 | 43.4 |
| OK | 33 | 38.8 |
| **H** | **156** | **21.7** |
| OK | 62 | 21.5 |
| Total | 9 821 | 28.6 |

Many of the same fieldworkers also appeared to be suspicious when BMI growth, rather than the mean of BMI, was used for identifying cheating. As shown in Table 4, six of the 12 most suspicious fieldworkers were subsequently found to have fabricated part or all of their interviews. The mean percentage change in BMI for the entire adult sample was 9%. G, who only interviewed 32 adults, returned a BMI growth rate of 173%, followed by E and B with 109% and 99% respectively. At the other end of the distribution, H's 117 respondents showed a decrease in BMI of 19%, on average.

Table 4: Suspicious Fieldworkers and Mean Adult BMI

| Fieldworker | N | Mean % change |
|---|---|---|
| **G** | **32** | **173** |
| **E** | **67** | **109** |
| **B** | **75** | **99** |
| OK | 38 | 33 |
| **I** | **83** | **31** |
| OK | 89 | 25 |
| OK | 35 | 23 |
| OK | 14 | 20 |
| **J** | **80** | **20** |
| OK | 44 | -7 |
| OK | 40 | -15 |
| **H** | **117** | **-19** |
| Total | 5 560 | 9 |

We present the suspicious list obtained by using mean adult height growth between

waves in Table 5, below. Fieldworkers who were subsequently found to have fabricated data are shown in bold once again. Of the 5 710 adults for whom valid data were recorded in both waves, the mean change in height was a rise of 0.11%. H and A recorded the largest mean growth rates of around 5%. The four most suspicious fieldworkers at the other end of the distribution, E, I, B and G recorded very large negative growth in height ranging from -4.95% to -14.70%.

Table 5: Suspicious Fieldworkers and Mean Change in Adult Height

| Fieldworker | N | Mean % change |
|---|---|---|
| **H** | **118** | **5.21** |
| **A** | **120** | **4.86** |
| OK | 35 | 4.81 |
| OK | 41 | 4.72 |
| OK | 45 | 4.62 |
| **E** | **68** | **-4.95** |
| **I** | **83** | **-5.38** |
| **B** | **75** | **-7.14** |
| **G** | **32** | **-14.70** |
| Total | 5 710 | 0.11 |

Source: Own calculations using pre-public release NIDS Wave 2 data, 2010.

The final anthropometric method used to detect suspicious fieldworkers was a visual inspection of spike-plots of the weight distributions. This allowed us to quickly observe heaping at focal points. Moreover, one weakness of the three methods used above was that they would only diagnose cheating if the proportion of interviews that were faked was 'substantial' enough to affect the mean.[14] This method was not dependent on the mean of the weight distribution obtained by the fieldworker.[15] It could thus be informative even in cases where a fieldworker had cheated on only a small fraction of surveys.

We provide two spikeplots as illustrative examples of recorded weights for adults,

---

[14]We also considered fieldworkers who captured anthropometric data that consisted of a 'high' number of outliers. This did not prove to be effective. Almost all fieldworkers had some outliers, which could arise because the respondent really did have an exceptional height or weight, or due to measurement error in either wave 1 or wave 2.

[15]An additional approach that we used was to sort fieldworkers by the $10^{th}$, $25^{th}$, $50^{th}$, $75^{th}$ and $90^{th}$ percentiles of their obtained distributions of height, weight, BMI and the growth of these variables. These did not yield substantially new insights beyond those already obtained from the methods already described and employed.

with values restricted between 48 and 100. Figure 3 shows the spike-plot of the weight distribution of a non-suspicious fieldworker. This fieldworker interviewed 39 adults (with weights between 48 and 100) and recorded two weight values for each of them, hence the uniformity at a frequency of two on the y-axis. Only two adults out of the 39 had the same weight  52kg and 81.1kg. Contrast this to the spike-plot of a suspicious fieldworker, shown in Figure 4. This fieldworker interviewed 125 adults and there is significant spiking at 65kg, 70kg, 75kg, 80kg and 85kg with very few of the other observations falling out of that range. Note that the y-axis goes up to 44, compared with Figure 2 where it only goes up to four. This distribution immediately raised suspicions that the fieldworker had made up the anthropometric data at best, and had fabricated the entire interview at worst.

Figure 3: Spike-plot of Non-suspicious Weight Distribution



Source: Own calculations using pre-public release NIDS Wave 2 data, 2010.

Figure 4: Spike-plot of Suspicious Weight Distribution



Source: Own calculations using pre-public release NIDS Wave 2 data, 2010.

Bringing the three anthropometric measures together allowed us to create a crude index of suspicion. Fieldworkers were assigned scores of zero to three, depending on how many times they were flagged as suspicious in each of the diagnostic methods. Table 6 shows the top 12 most suspicious fieldworkers according to the three anthropometric diagnostics. Every fieldworker who scored three out of three was later found to have fabricated data. Of the fieldworkers who were flagged as suspicious using the anthropometric measures, A and B were also flagged as the two most suspicious fieldworkers using the Benford's law method.

Table 6: Combined Anthropometric Suspicion Index

| Fieldworker | BMI | BMI Growth | Height Growth | Row Total |
|---|---|---|---|---|
| **H** | **1** | **1** | **1** | **3** |
| **E** | **1** | **1** | **1** | **3** |
| **B** | **1** | **1** | **1** | **3** |
| **G** | **1** | **1** | **1** | **3** |
| **I** | | **1** | **1** | **2** |
| OK | 1 | 1 | | 2 |
| **F** | **1** | | | **1** |
| **A** | | | **1** | **1** |
| **J** | | **1** | | **1** |
| OK | | 1 | | 1 |
| OK | | | 1 | 1 |
| OK | 1 | | | 1 |
| Col. Total | 7 | 8 | 7 | 21 |

An important issue to consider is that data in wave 1 could have been fabricated as well. If interviews (or parts of interviews) were faked in the first wave of data, this could feed through and make large changes in anthropometric data look suspicious. This would be problematic as the error is entering in the first period, rather than in the second which is the focus of our concern. On average, however, even if the first wave contains fabricated data, this should be diluted as different fieldworkers interviewed different respondents in both waves. The probability that the majority of a wave 2 fieldworker's valid interviews are combined with majority fake wave 1 data is small but non-trivial, given the spatial logistics under which fieldwork was conducted in both waves.

# 4 The NIDS Operational Response

A meta-list of suspicious fieldworkers was drawn up using a combination of the Benford's law rankings and the anthropometrics rankings. The NIDS operations centre initiated an intensive set of telephonic callbacks in order to verify whether or not the interviews of suspicious fieldworkers had actually been conducted. Priority was given to calling back respondents who were interviewed by the most suspicious fieldworkers and the NIDS team worked down the list systematically, calling every household for which data had been collected by that fieldworker, until there was a high level of confidence about the veracity of the data.

The callbacks comprised a set of questions intended to establish initially whether an interview had taken place or not, and whether the entire interview had been completed. A copy of the list of questions asked in the verification process can be found in Figure 7 in Appendix C.
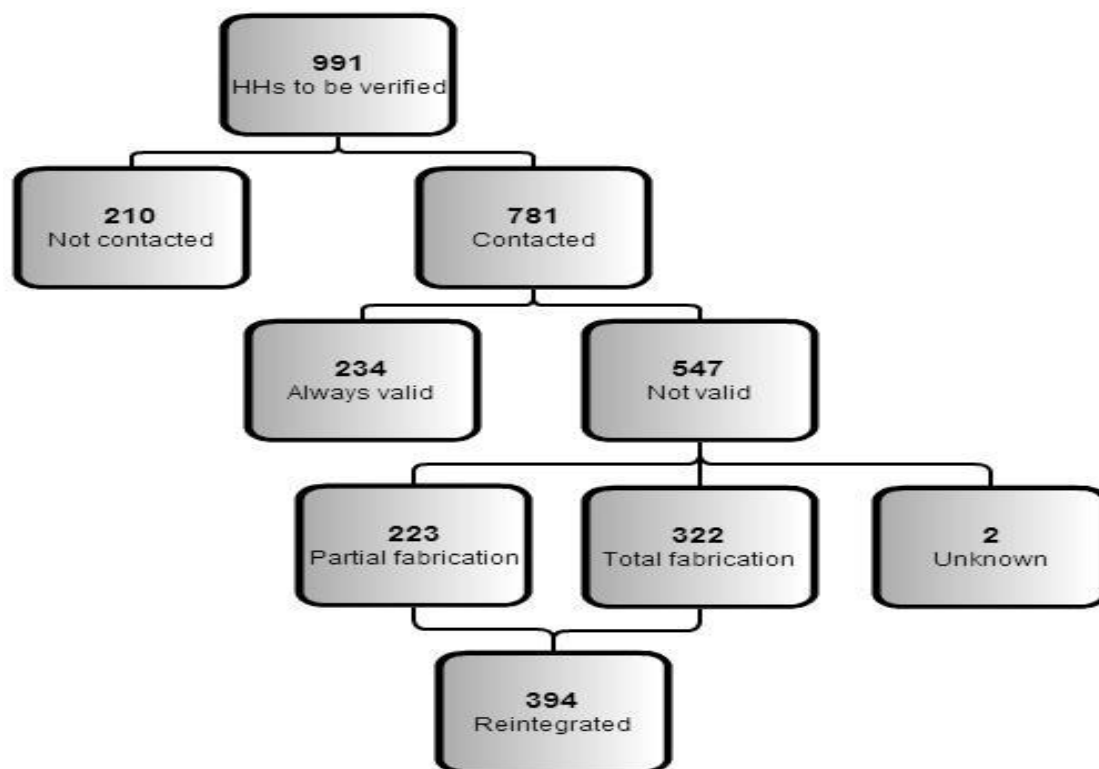
Where fraud was evident, all data from those interviews were rejected at the expense of the company conducting the fieldwork. This was true regardless of whether there was partial or total fabrication of the relevant interview. During the data collection phase, fieldworkers were sent out in teams of three, with a team leader and two additional interviewers. The intensive auditing revealed that data fabrication was generally a team-specific phenomenon. That is to say, if one fieldworker in Team A was found to have cheated, there was a fairly high probability that the other two team members also cheated. In all cases in which the team leader was found to have cheated, the other two fieldworkers also fabricated their data, to a greater or lesser extent. The auditing turned up one interesting case where a suspicious fieldworker was found to have used the scale and measuring device incorrectly for the anthropometrics. In this case, only the anthropometric data was flagged as invalid and the fieldworker was sent for additional training.

Once the NIDS callback team had reached the point where fraud was no longer being uncovered, a new phase of data collection was put in place to re-interview the appropriate respondents. Instead of deleting data that was fabricated, thereby reducing the sample size in the second wave, correctly collected information was re-integrated into the dataset. Figure 5 describes the different outcomes of the verification process. Overall, we identified 991 households that needed verification. Of these, 781 households were successfully contacted, meaning that over 10% of the total sample was called back in the verification process. Of these 781 households, 234 were verified as having been validly interviewed. Furthermore, it was found that 547 households had data that was

not valid, either because of partial or total fabrication. These 547 households were made up of 223 partial fabrications, 322 total fabrications and 2 which could not be classified.[16]

In summation, of all the households on the suspicion list that were successfully contacted, over 70% had large data quality concerns that were driven entirely by fieldworker cheating. This represents 7.3% of the wave 2 households at the time that the verification process started. Of the 547 households with problematic data, a total of 394 households were successfully re-interviewed and re-integrated into the dataset. Thirteen fieldworkers, or approximately 10% of the total number of fieldworkers that were employed to conduct fieldwork for wave 2 of NIDS, had produced some data that was entirely fabricated.[17] The rates of fabrication across fieldworkers who were found to have cheated ranged from 10% to 67% of all households interviewed.

Figure 5: Verification process and outcomes



Source: Own calculations using pre-public release and public release NIDS Wave 2 data, 2010/2011.

---

[16]These two households were successfully contacted and the reference person indicated that there was a problem with the interview process, but refused to provide further information.

[17]Note that not all of their data was fabricated, but that some positive proportion of the data that they generated certainly was fabricated.

# 5 Implications for Analysis

By how much would the presence of the fabricated data have affected our estimates, had the cheating fieldworkers not been discovered? Assuming that some fabrication is probably present in all large South African surveys, and simultaneously, that most field-workers are probably honest, should we be wary of most empirical results? Alternatively, does the measurement error caused by fieldworker cheating have relatively small effects on our estimates, such that, for practical purposes, we may ignore its implications with respect to research findings? In addition to the resources invested in the production of data, considerable time, energy and resources are invested by users of these data, and research findings subsequently feed into important policy making discussions and debates. Measuring the effects of fieldworker cheating on the validity of empirical findings is the objective of this section of the paper.

From an econometric perspective, data fabrication leads to measurement error for potentially all of the variables in some subset of the data. *A priori*, we cannot make a general prediction about the effects of data fabrication on subsequent estimates, as the effects, if any, will depend on multiple factors. These factors include the fraction of the overall dataset that is fabricated, the difference between the fabricated data and the true data that it represents, the type of estimator being implemented, whether the fabrication results in classical[18] or non-classical[19] measurement error in the variable or variables that are being used, and the magnitudes of such measurement error. Moreover, if one is using a multivariate estimator, the empirical effects will depend on the relationship between the variables being used in the fabricated dataset, relative to the true relationship between those variables. Any theoretical predictions thus have to be restricted by quite a specific set of criteria.

Nonetheless, there are some well known and fairly general effects that measurement error in a regressor will have in a regression analysis. First, measurement error in an independent variable will result in a violation of the orthogonality condition. This will induce biased estimates of the $\hat{\beta}$ vector obtained from an OLS regression (Wooldridge,

---

[18]Assuming the true model is $y^* = X^*\beta + \epsilon$ but we measure $X = X^* + \mu$ and $y = y^* + v$, under the conditions of classical measurement error, $u$ and $v$ are i.i.d. and uncorrelated with $X^*, y^*$ and $\epsilon$ and the estimated $\beta$ coefficients are biased in the direction of zero. See Bound et al. (2001) for a comprehensive overview of the literature on bias due to measurement error.

[19]Of course, the misclassification of a categorical variable such as labour market status or a dummy variable such as employed/unemployed cannot be thought of in the same way as classical measurement error, as the error itself cannot be mean zero. In fact, for dummy variables, the measurement error must be negative. The case for measurement error in categorical variables is not as straightforward, but a thorough treatment is beyond the scope of this paper. See Krueger and Summers (1988) for a discussion of the results of measurement error in categorical regressions.

2002). In the case of classical measurement error, this will result in an attenuation bias, that is, a bias of the estimated coefficient towards zero. Second, some common estimators, such as fixed effects estimators and first difference estimators, are more sensitive to a particular endogeneity problem than a standard OLS estimator (Griliches and Hausman, 1986).

In addition, Schnell (1991) and Schräpler and Wagner (2005) find that univariate statistics such as means, medians and variance are generally robust to the presence of fake data, where the prevalence of fake data is less than 5%. However, the negative effects of fake data begin to compound as analysis moves to a multivariate setting, particularly when some of the commonly-used panel data estimators are used (Schnell, 1991; Schräpler and Wagner, 2005).

In order to provide an illustrative example of the effects of cheating in our dataset, we chose to investigate the broad theme of understanding the effects of finding employment on health, as measured by BMI. We chose this area for investigation for two reasons. First, we have spent some time documenting the fabrication that took place in the labour market module as well as the height and weight measurement module. Looking at the effect of finding employment on a measure of well-being such as BMI complements our section on detecting fabrication. Second, the determinants of BMI as well as the effect of employment on BMI are topics that have received a great deal of attention in the recent South African literature (Wittenberg, 2013, 2009; Ardington and Gasealahwe, 2012; Ardington and Case, 2009) and our study provides an important addition to these.[20] We then needed to choose a set of variables and a set of estimators for our analysis. For the variables, we include BMI, age, education and labour market status. Following the theoretical discussion in the preceding paragraphs, we calculate the mean BMI, the labour market transition rates and finally, we fit OLS and first difference regression models of BMI on age, education and employment.[21]

## 5.1   Data

To implement the analysis, we constructed two datasets. The first, which we refer to as the 'Dirty' dataset, is a combination of the 'Always Correct' data combined with the 'Fake' data at the time that our verification process was completed. Essentially,

---

[20]We also modelled the effect of receiving the state old age pension on labour force participation using the dirty and clean datasets. The results of this are available from the authors on request.

[21]We collapse the four labour market states into a binary employed variable for the regressions. We did this as it made more sense theoretically and it made the discussion of the regression results simpler. We also performed the estimations with the labour market states disaggregated and the overall findings do not change substantially (not reported).

it represents what the NIDS wave 2 dataset would have been if the cheating had gone undetected, and the survey was completed at the date that our verification process drew to a close. The second dataset, which we refer to as the 'Clean' dataset, is composed of the same 'Always Correct' data, combined with the subsequently corrected data where such correction was possible.

The variables that we use are all at the individual level. They are:

- **BMI** - This is calculated as a person's mass in kilograms divided by height in metres squared. Since each respondent had either two or three measures of height and weight each, we used the average of all recorded measures. There was no pre-population of this variable in wave 2.

- **Age** - This was measured in integer years. The variable triggered a data confirmation question for fieldworkers if the respondent had aged by less than 1 year or by more than 2 years between wave 1 and wave 2. Fieldworkers had access to the wave 1 roster, hence even fabricated surveys would likely have appropriate data in wave 2.

- **Years of education** - This variable is bounded between 0 and 15. The wave 1 information on education was also given to fieldworkers. Moreover, if the education levels had increased by more than 2 years, or had decreased between wave 1 and wave 2, the software would ask for confirmation from the fieldworkers. Thus, we expect to have only a small amount of measurement error on this variable in the fabricated data.

- **Male** - This is an indicator variable that captures the sex of the respondent. It was pre-populated based on the wave 1 dataset.

The labour market status variables are comprised of four mutually exclusive indicator variables,[22] which represent the labour market state of respondents. These are all derived from the labour market section of the survey. These variables were not pre-populated. They are:

- **Employed** - This is an indicator variable that takes on a value of one if the respondent had any form of employment at the time of the interview.

---

[22]There are some cases where the questions used to drive a person's labour market status were not answered. In these cases, we cannot define their status. Otherwise, the four variables would be mutually exclusive and exhaustive.

- **Unemployed (searching)** - This is an indicator variable that takes on a value of one if the respondent was not employed but was actively looking for work in the month prior to the interview.

- **Unemployed (discouraged)** - This is an indicator variable that takes on a value of one if the respondent was not employed and was not actively looking for work in the month prior to the interview, but stated that he/she would like to have a job. The difference between the searching and non-searching unemployed conforms to the standard ILO definitions for these categories (International Labour Office, 2011).

- **Not economically active** - This is an indicator variable that takes on a value of one if the respondent was not employed, was not actively looking for work in the month prior to the interview and stated that they would not accept a job offer.

There are a few additional data issues that require elaboration. First, for our entire analysis, we restrict our estimation sample to include only the adult African sub-population aged 18 to 65.[23] Second, we restrict the sample to include BMI values of less than 50, as we were concerned that some of the extremely high BMI values were due to the scales being inadvertently set to pounds instead of kilograms. In addition, we exclude any observation with any covariate missing from our samples, as they would not survive into our regression analyses. Third, we do not make use of either the sampling weights or attrition-corrected weights in any of the subsequent analysis. Our objective is not to replicate population level analyses, but merely to compare the differences between estimates obtained from the dirty and the clean dataset. Moreover, we would have had to recalculate all of the post-stratification weights, as the datasets that we use do not represent the full sample due to the time at which we completed our audit. Fourth, in our regressions we re-weight the subsequently corrected data by the inverse of the ratio of the number of corrected fakes to the number of fakes. We do this because we want the weighted fraction of data from the 'Always Correct' data to be the same in the Dirty and Clean datasets. The implicit assumption here is that the group of corrected fakes are representative of the group of fakes that we were unable to subsequently re-interview.

The sample sizes, and how they are affected by our restrictions, are displayed in Table 7 below. We observe that the BMI cutoff of 50 is not too onerous. We lose 106 and 84 observations from the Dirty and Clean datasets respectively. This represents less

---

[23]This is because we have small sample sizes for the other race groups, especially when using the balanced panel members from wave 1 and wave 2. Wittenberg (2013) applies a similar restriction to the NIDS data.

than two percent of either sample, and a substantial fraction of these are observations from the 'Always Correct' subset of the data. Our final sample sizes for the OLS and First Differences analysis are 6 768 and 5 388 observations for the Dirty dataset, and 6 576 and 5 263 for the Clean dataset, respectively. The sample sizes for the First Differences regressions are substantially smaller for two reasons. First, new household members would not have been interviewed in wave 1. Second, any missing data in any covariate in wave 1 would have resulted in that observation being dropped from the First Differences sample as well.

Table 7: Sample sizes

| Number | Dirty | Clean |
|--------|-------|-------|
| **All** | 6 874 | 6 660 |
| **BMI<50** | 6 768 | 6 576 |
| **CS OLS** | 6 768 | 6 576 |
| **FD** | 5 388 | 5 263 |

Source: Own calculations using pre-public release NIDS Wave 2 data, 2010.
Sample restricted to African adults aged 18 to 65 in wave 2.

## 5.2 Summary Statistics

The means of the variables in each of the sub-datasets, as well as the Clean and Dirty datasets, are provided in Table 8 below. Note that the mean of a variable in the Dirty dataset will be a weighted average of the corresponding means in the Fake dataset and the Always Correct dataset, with the weight being determined by the proportion of the data in the Dirty dataset that originates from the Fake and Always Correct datasets, respectively. Similarly, the mean in the Clean dataset will be a function of the means in the Corrected Fake and Always Correct datasets. Any differences in the means between the Clean and Dirty datasets must therefore reflect differences in the means between the Fake and Corrected Fake datasets, combined with the differences in their respective sample sizes.

For the BMI, age, male and years of education variables that we use, the difference in means between the Fake and Corrected Fake datasets is relatively small. Thus, in the aggregated Dirty and Clean datasets, the difference in means for these variables is very small, at less than 0.15 units in each case. For some of the other variables, such as Unemployed (discouraged), the difference in means between the Fakes and Corrected Fakes is somewhat larger, at 2.81 percentage points, but the aggregate difference in

30

Table 8: Means of Variables Used in Analysis

|  | Fakes | Corrected Fakes | Always Correct | Dirty Dataset | Clean Dataset |
|---|---|---|---|---|---|
| BMI | 26.73 | 25.76 | 26.94 | 26.92 | 26.89 |
| Age | 35.45 | 37.08 | 36.16 | 36.11 | 36.20 |
| Education (years) | 6.74 | 7.47 | 8.12 | 8.03 | 8.09 |
| Employed | 19.82% | 28.40% | 33.99% | 33.05% | 33.77% |
| Unemployed (searching) | 5.12% | 15.56% | 10.98% | 10.59% | 11.16% |
| Unemployed (discouraged) | 0.69% | 3.50% | 5.93% | 5.59% | 5.84% |
| Not economically active | 73.27% | 52.53% | 48.60% | 50.24% | 48.75% |
| Male | 40.98% | 39.69% | 40.26% | 40.31% | 40.24% |
| Number | 449 | 257 | 6 319 | 6 768 | 6 576 |

Source: Own calculations using pre-public release and public release NIDS Wave 2 data, 2010.
Samples are restricted to Africans aged 18 to 65 in wave 2, with BMI values less than 50.
The number of fakes do not equal the number of corrected fakes because not all faked respondents were
successfully re-interviewed. The means in the Dirty and Clean dataset do not precisely correspond to
the weighted means obtained from the first three columns due to rounding effects.

means for these variables remains relatively small. This is because the relative weightings contributed by Fakes and Corrected Fakes to the means in the Dirty and Clean datasets are also small.

In contrast, for the Employed, Unemployed (Strict) and NEA variables, the difference in means between the Fakes and Corrected Fakes is substantial. Even these differences, however, get substantially moderated by the weights when we calculate the means of the Dirty and the Clean datasets. For example, let us consider the percentage that are employed in the Fakes and Corrected Fakes datasets. The difference in the means is large, at 8.58 percentage points. Nonetheless, the weight that these contribute to the Dirty and Clean datasets are relatively small, at 6.6 and 3.9 percent. Thus the aggregate difference in mean percentage employed between the Dirty and Clean datasets is only 0.72 percentage points. This is substantially smaller than the corresponding difference in means between the Fakes and Corrected Fakes dataset, and depending on one's interest, may or may not be considered to be 'substantial'. Overall then, we confirm the finding by Schnell (1991) that a univariate statistic such as the mean is generally robust to the presence of a small amount of fake data.

## 5.3 Transition Matrices

We next consider labour market dynamics, by calculating transition matrices across labour market states. Each row in the transition matrix in Table 9 contains the distribution of labour market states observed for respondents in wave 2, conditional on their wave 1 state. For example, of the 2 348 people who were not economically active (NEA) in wave 1 in the Dirty dataset, 73% were NEA in wave 2, while 12.27% had found employment.

The differences between the Dirty and Clean datasets are contained in the third matrix in Table 9, below. A clear pattern emerges, even though the magnitudes are not too large in absolute value. The presence of the faked data would cause us to systematically overstate the likelihood of transiting into the NEA state, and underestimate the probability of transiting or remaining in any of the other states, regardless of the initial wave 1 state. This reconciles well with the cross-sectional differences in means from Table 8, where the difference in the mean percentage that were NEA between the Fakes and Corrected Fakes was more than twenty percentage points, and the corresponding difference between the Clean and Dirty means was about 1.5 percentage points. It also seems consistent with the time-saving hypothesis that we postulated would accompany cheating behaviour.

Nonetheless, the magnitudes of the differences are generally below 1 percentage point, and only one entry is above 2 percentage points. Whether one considers these differences to be material or not will, once again, depend on one's perspective. If one were interested in population-level dynamics, or in long run forecasting of retirement behaviour, then they could well be material. On the other hand, if one were simply interested in the conditional probabilities of transitioning across labour market states, then the Dirty and Clean data would not have yielded a very different understanding of the aggregate levels of churning in the SA labour market.

## 5.4 Regression Results

Our final set of analyses involves estimating the regression coefficient of employment on BMI. We first present the cross-sectional regression results using wave 2 data, and compare the coefficients from the Dirty and Clean datasets. One weakness of this approach, if we think that BMI is a proxy for health, is that we are likely to have a selection problem since healthier people will probably be more likely to find employment. A natural extension would be to estimate a fixed effects model of employment on BMI. We thus also present the First Differences (FD) form of the regression.

Table 9: Transition Matrices Across Labour Market States (%)

| | | NEA | Unemployed (Discouraged) | Unemployed (Strict) | Employed | Total | N |
|---|---|---|---|---|---|---|---|
| | | **Wave 2 State (Dirty Dataset)** | | | | | |
| **W1 State** | **NEA** | 73.00 | 6.05 | 8.69 | 12.27 | 100 | 2 348 |
| | **Unemp. D.** | 53.13 | 10.34 | 12.26 | 24.28 | 100 | 416 |
| | **Unemp. S.** | 43.32 | 6.62 | 22.42 | 27.64 | 100 | 861 |
| | **Employed** | 30.16 | 4.31 | 7.48 | 58.05 | 100 | 2 205 |
| | | **Wave 2 State (Clean Dataset)** | | | | | |
| **W1 State** | **NEA** | 71.47 | 6.32 | 9.58 | 12.63 | 100 | 2 296 |
| | **Unemp. D.** | 52.06 | 10.65 | 12.59 | 24.7 | 100 | 413 |
| | **Unemp. S.** | 41.04 | 6.96 | 23.47 | 28.54 | 100 | 848 |
| | **Employed** | 29.47 | 4.38 | 7.53 | 58.62 | 100 | 2 192 |

**Difference Dirty-Clean (%)**
**Wave 2 State**

| | | NEA | Unemployed (Discouraged) | Unemployed (Strict) | Employed |
|---|---|---|---|---|---|
| **W1 State** | **NEA** | 1.53 | -0.27 | -0.89 | -0.36 |
| | **Unemp. D.** | 1.07 | -0.31 | -0.33 | -0.42 |
| | **Unemp. S.** | 2.28 | -0.34 | -1.05 | -0.90 |
| | **Employed** | 0.69 | -0.07 | -0.05 | -0.57 |

Source: Own calculations using pre-public release and public release NIDS Wave 2 data, 2010.

We do not have strong priors regarding the differences between the Dirty and Clean datasets based on econometric theory. The measurement error in the employment dummy cannot be classical measurement error as the variable is a binary variable. Moreover, in the FD model, the error has a particular distribution that is not symmetric. Nonetheless, we do expect that we have a measurement error problem, and this will cause an endogeneity problem. We also know that the FD estimator is more sensitive to measurement error than the OLS estimator, so we might expect that the presence of fabricated data will have a stronger effect on the FD coefficients than the OLS coefficients.

In Table 10, we present the regression outputs from estimating the OLS model on the Dirty and Clean data. Our dependent variable is BMI and our regressors are age, gender, education and employment status. The overall finding is that the regression results look very similar. The R-squared values differ by 0.01 units, and none of the coefficients are

statistically significantly different at any reasonable significance level. The education variable, where we do have something that resembles classical measurement error within a limited range, is indeed slightly smaller in the regression using the Dirty dataset, but the difference is only 0.03 BMI units, which is quite negligible. The coefficients on the employed dummy are both positive and significant. The coefficient is slightly larger when using the Dirty dataset (0.76 vs 0.70), but again they only differ by 0.06 BMI units.

Table 10: Cross-sectional and First-differenced Regressions

| | Cross sectional | | First differenced | |
| | Dirty | Clean | Dirty | Clean |
| Variables | W2 BMI | W2 BMI | $\Delta$ BMI | $\Delta$ BMI |
|---|---|---|---|---|
| Age | 0.14*** | 0.14*** | 0.00 | 0.00* |
| | (0.01) | (0.01) | (0.00) | (0.00) |
| Male | -4.52*** | -4.60*** | | |
| | (0.15) | (0.15) | | |
| Education | 0.13*** | 0.16*** | 0.08** | 0.02 |
| | (0.02) | (0.02) | (0.04) | (0.04) |
| Employed | 0.76*** | 0.70*** | 0.31** | 0.18 |
| | (0.16) | (0.16) | (0.15) | (0.14) |
| Constant | 22.50*** | 22.00*** | 1.03*** | 1.00*** |
| | (0.37) | (0.38) | (0.08) | (0.07) |
| | | | | |
| Observations | 6,768 | 6,576 | 5,388 | 5,263 |
| R-squared | 0.19 | 0.20 | 0.00 | 0.00 |

Standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

Source: Own calculations using pre-public release and public release NIDS Wave 2 data, 2010.
Samples are restricted to Africans aged 18 to 65 in wave 2, with BMI values less than 50. Columns 2
and 3 present results from cross-section OLS estimation. Columns 4 and 5 present results
first-differenced regressions, and the regressors should be read as differences, rather than levels.

The similarities are not surprising given what was observed in Table 8. In the cross-sectional datasets the differences in the means of the relevant variables were all very small, and most of the data in both the Dirty and Clean datasets are obtained from the Always Correct dataset.

Our final set of results are obtained from the FD regressions and are presented in the last two columns of Table 10. Note that the male dummy gets dropped as it is a time invariant variable. Our findings from this component of our analyses are a bit more

nuanced than those from our earlier analyses.

When we compare the differences in the FD results between the Dirty and Clean datasets, we notice that the Dirty coefficients on education and employment are larger than those obtained from the Clean dataset, and they are statistically significant whereas those obtained from the Clean dataset are not statistically significant. On the other hand, the differences in magnitude are 0.06 and 0.13 BMI units for the education and employed variables, which are not particularly large. Moreover, the differences in the coefficients are not statistically significant. From this perspective, the fabricated data does affect our estimates, but not in a meaningful way.

Alternatively, when we compare the OLS and FD results within each dataset, we observe that the coefficients from both the Dirty and Clean datasets are reduced quite substantially. For example, in the Dirty regressions, the coefficient on education is 0.13 in the OLS regression but decreases to 0.08 in the FD regression. The decrease obtained in the Clean dataset is from 0.16 to 0.02, and also results in a change in the statistical significance of the coefficient.

A similar comparison between the OLS and FD coefficients focusing on the employed dummy yields larger decreases in the coefficients (in absolute value) for both the Dirty and the Clean datasets. If our only dataset had been the Dirty dataset, we would have concluded that using a longitudinal estimator results in a decrease of our estimated coefficient from 0.76 to 0.31, that is a decrease of 0.45 BMI units, although both coefficients are statistically significant at the 5% level. In contrast, if we performed the identical exercise using the Clean dataset, we would have concluded that using a longitudinal estimator results in a decrease of our estimated coefficient from 0.70 to 0.18, that is, a decrease of 0.52 BMI units. Moreover, we would observe that the FD coefficient, unlike the OLS coefficient, is not statistically significantly different from zero.

## 5.5   Discussion

We are aware that we are probably not getting the true 'causal' estimate of labour market status on BMI, but our focus is primarily on measuring the difference between the estimates obtained by using the fabricated data instead of the subsequently corrected data. This is the main contribution of this section of the paper. To our knowledge, all previous research on this topic has amounted to comparing an *ex ante* dataset containing fabricated data to an *ex post* dataset where the fraudulent data has been deleted, but not replaced. In our case, where households with fabricated data were re-integrated into the NIDS dataset, we are in a unique position in that we observe both the fabricated

data as well as the subsequently corrected data.

The overarching question that we set out to answer was whether the fabricated data would have affected our estimates in a meaningful way. Our findings suggest that the answer to this question depends on the estimator being considered and the purposes for which the analysis is being conducted. The general picture that emerges, which is consistent with econometric theory, is that the cross-sectional estimates of means and OLS regressions are not substantially affected by the presence of a relatively small amount of fabricated data. At the same time, the identical amount of fabrication does affect the longitudinal estimators – the transition matrices and first differences regressions. With the transition matrices, the differences are relatively small in absolute value, but may be meaningful depending on the purposes for which they are being calculated. For the FD regressions, the difference in the estimates would have led us to reach quite different overall conclusions.

Regardless of one's perspective, better quality data is always desirable. Thus, other things being equal, we should always get the best quality data that we can. Unfortunately, running a survey is a costly and complex task, and the costs of auditing and monitoring fieldworkers competes with several other tasks for resources in a finite budget. Thus, the 'other things being equal' assumption is not very realistic. Nonetheless, given the overall costs and effort invested in running a large survey, the marginal costs of performing some generic data quality checks for fabricated data seems to be highly warranted, especially in an environment where we now have evidence of fabrication in several studies. In our study, the estimators that were most affected were the longitudinal estimators. At the same time, the marginal cost of detecting fabricated data can be substantially lowered when one has longitudinal data, as one can then look for intertemporal data anomalies in addition to cross-section data anomalies. This makes an even stronger case for the argument that adequate resources be allocated for identifying data fabrication in longitudinal studies.

# 6   Conclusion

In this paper, we argued that the incidence of fieldworker cheating is widespread. We documented cheating and potential cheating in five substantial South African surveys. Of the various methods that we considered to detect fraudulent data, two were more useful in our context than the others.

Looking forward, there may be ways to improve on our process for identifying fabrication. First, survey companies that are using computers with built in GPS devices to fill in questionnaires can use the software to capture the time and place that a survey was conducted. This can be done without the knowledge of the fieldworker, which will aid in detection. Using a GPS software will allow much better monitoring of fieldworkers' whereabouts while they are in the field. In addition, fieldworkers that fabricate data are likely to complete entering the data much faster than an actual survey would take to complete. These two pieces of information alone would greatly improve the data quality auditing process.

Second, wireless networks and cellular technologies are now widespread even in developing countries. It is thus not unrealistic to expect to get data with just a day's lag. Previously, while using paper questionnaires, it would take months to obtain data in an electronic form. After a survey was completed in the field, it would eventually get sent in batches to the head office of the fieldwork company where a cursory data quality check would be performed. The head office would then send the questionnaires to the survey organisation which would do its own data quality controls, after which it would be sent to a double-blind data capturing process. By the time the data was ready to be interrogated for inconsistencies, it would not have been possible to discipline any cheating fieldworkers. With the real-time uploading of the data, one can now check on each fieldworker much earlier in the process, and constantly monitor each fieldworker's performance. This enables survey organisations to fire cheating fieldworkers, as well as compel the fieldwork company to redo the interview.

Other possibilities might be to use built in cameras to take photographs of survey respondents, real-time callbacks to ensure that the interview did in fact take place, and to strategically not pre-populate certain variables in longitudinal studies. In summation, it seems that there are several relatively low cost ways in which survey organizations can use modern technology to minimise both the likelihood of fieldworker cheating, as well as the impact of such cheating on the overall quality of the data.

# 7 Bibliography

Ardington, C. and Case, A. (2009), Health: Analysis of the NIDS wave 1 dataset, NIDS Discussion Paper 2, National Income Dynamics Study, University of Cape Town.

Ardington, C. and Gasealahwe, B. (2012), Health: Analysis of the NIDS wave 1 and 2 datasets, SALDRU Working Papers 80, Southern Africa Labour and Development Research Unit, University of Cape Town.

Benford, F. (1938), 'The law of anomalous numbers', *Proceedings of the American Philosophical Society* **78**(4), 551–572.

Birnbaum, B. (2012), Algorithmic Approaches to Detecting Interviewer Fabrication in Surveys, PhD thesis, University of Washington.

Bound, J., Brown, C. and Mathiowetz, N. (2001), Measurement error in survey data, *in* J. Heckman and E. Leamer, eds, 'Handbook of Econometrics', Vol. 5 of *Handbook of Econometrics*, Elsevier, chapter 59, pp. 3705–3843.

Bredl, S., Winker, P. and Kötschau, K. (2008), A statistical approach to detect cheating interviewers, Technical Report 39, Diskussionsbeiträge: Zentrum für internationale Entwicklungs-und Umweltforschung.

Carslaw, C. A. (1988), 'Anomalies in income numbers: Evidence of goal oriented behavior', *Accounting Review* **63**(2), 321–327.

Cho, M., Eltinge, J. and Swanson, D. (2003), Inferential methods to identify possible interviewer fraud using leading digit preference patterns and design effect matrices, *in* 'Proceedings of the American Statistical Association (Survey Research Methods Section)', pp. 936–941.

Devey, R., Skinner, C. and Valodia, I. (2006), Definitions, data and the informal economy in South Africa: a critical analysis, *in* V. Padayachee, ed., 'The Development Decade?: Economic and Social Change in South Africa, 1994-2004', HSRC Press, chapter 15, pp. 302–323.

Durtschi, C., Hillison, W. and Pacini, C. (2004), 'The effective use of Benford's law to assist in detecting fraud in accounting data', *Journal of Forensic Accounting* **5**(1), 17–34.

Griliches, Z. and Hausman, J. A. (1986), 'Errors in variables in panel data', *Journal of Econometrics* **31**(1), 93–118.

Hill, T. P. (1995), 'A statistical derivation of the significant-digit law', *Statistical Science* **10**(4), 354–363.

International Labour Office (2011), Unemployment, underemployment and inactivity indicators, *in* 'Key indicators of the labour market', ILO, chapter 4.

Krueger, A. B. and Summers, L. H. (1988), 'Efficiency wages and the inter-industry wage structure', *Econometrica* **56**(2), 259–93.

Lam, D., Ardington, C., Branson, N., Maughan-Brown, B., Menendez, A., Seekings, J. and Sparks, M. (2012), The Cape Area Panel Study: Overview and technical documentation waves 1-2-3-4-5(2002-2009), Technical report, The University of Cape Town.

Leibbrandt, M., Woolard, I., Finn, A. and Argent, J. (2010), Trends in South African income distribution and poverty since the fall of apartheid, OECD Social, Employment and Migration Working Papers 101, OECD Publishing.

Li, J., Brick, J., Tran, B. and Singer, P. (2011), 'Using statistical models for sample design of a reinterview program', *Journal of Official Statistics* **27**(3), 433–450.

May, J., Agüero, J., Carter, M. and Timaeus, I. (2007), 'The KwaZulu-Natal Income Dynamics Study (KIDS) 3rd wave: Methods, first findings and an agenda for future research', *Development Southern Africa* **24**(5), 629–648.

Murphy, J., Baxter, R., Eyerman, J., Cunningham, D. and Kennet, J. (2004), A system for detecting interviewer falsification, *in* 'American Association for Public Opinion Research, 59th Annual Conference', pp. 4968–4975.

Porras, J. and English, N. (2004), Data-driven approaches to identifying interviewer data falsification: The case of health surveys, *in* 'Proceedings of the American Statistical Association (Survey Research Methods Section)', pp. 4223–4228.

Schäfer, C., Schräpler, J.-P., Müller, K.-R. and Wagner, G. G. (2004), Automatic identification of faked and fraudulent interviews in surveys by two different methods, DIW Discussion Paper 441, DIW Berlin, German Institute for Economic Research.

Schnell, R. (1991), 'Der Einfluß gefälschter Interviews auf Survey Ergebnisse', *Zetischrift für Soziologie* **20**, 25–35.

Schräpler, J.-P. (2011), 'Benford's law as an instrument for fraud detection in surveys using the data of the socio-economic panel (SOEP)', *Journal of Economics and Statistics* **231**(5-6), 685–718.

Schräpler, J.-P. and Wagner, G. (2005), 'Characteristics and impact of faked interviews in surveys - An analysis of genuine fakes in the raw data of SOEP', *Allgemeines Statistisches Archiv* **89**(1), 7–20.

Schreiner, I., Pennie, K. and Newbrough, J. (1988), Interviewer falsification in census bureau surveys, *in* 'Proceedings of the American Statistical Association (Survey Research Methods Section)', pp. 491–496.

Scott, P. D. and Fasli, M. (2001), 'Benford's law: An empirical investigation and a novel explanation', *Unpublished manuscript* .

Statistics South Africa (2000), 'Time Use Survey: Fieldworker's manual', *Unpublished training manual* .

Swanson, D., Cho, M. and Eltinge, J. (2003), Detecting possibly fraudulent or error-prone survey data using benfords law, *in* 'Proceedings of the American Statistical Association (Survey Research Methods Section)', pp. 4172–4177.

Thomas, J. K. (1989), 'Unusual patterns in reported earnings', *The Accounting Review* **64**(4), 773–787.

Wittenberg, M. (2009), Weighing the value of asset proxies: The case of the body mass index in South Africa, SALDRU Working Papers 39, Southern Africa Labour and Development Research Unit, University of Cape Town.

Wittenberg, M. (2013), 'The weight of success: The body mass index and economic well-being in Southern Africa', *Review of Income and Wealth* **59**, S62–S83.

Wooldridge, J. M. (2002), *Econometric analysis of cross section and panel data*, The MIT press.

# A  Time Use

Table 11: Time Use Survey Selection Grid

| Persons 10 years + | HH1 | HH2 | HH3 | HH4 | HH5 | HH6 | HH7 | HH8 | HH9 | HH10 |
|---|---|---|---|---|---|---|---|---|---|---|
| **1** | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| **2** | 1 2 | 1 2 | 1 2 | 1 2 | 1 2 | 1 2 | 1 2 | 1 2 | 1 2 | 1 2 |
| **3** | 1 2 | 1 3 | 2 3 | 2 3 | 1 3 | 2 3 | 1 2 | 2 3 | 1 3 | 2 3 |
| **4** | 2 4 | 1 3 | 1 3 | 2 4 | 1 3 | 2 4 | 1 3 | 2 4 | 1 3 | 2 4 |
| **5** | 3 5 | 1 4 | 1 3 | 2 4 | 1 5 | 2 4 | 2 4 | 4 5 | 1 2 | 2 4 |
| **6** | 5 6 | 4 6 | 1 2 | 1 2 | 1 5 | 4 6 | 1 5 | 3 5 | 4 6 | 1 3 |
| **7** | 2 6 | 4 6 | 2 5 | 5 7 | 2 4 | 4 7 | 5 7 | 1 4 | 2 6 | 1 4 |
| **8** | 1 5 | 1 3 | 6 8 | 2 5 | 1 4 | 5 6 | 2 3 | 5 7 | 6 8 | 2 8 |
| **9** | 4 9 | 1 3 | 4 9 | 1 5 | 2 7 | 2 9 | 2 3 | 4 5 | 7 8 | 2 6 |
| **10** | 3 9 | 1 6 | 2 3 | 4 9 | 1 3 | 8 10 | 5 6 | 3 7 | 2 5 | 8 9 |

Source: Time Use Survey Fieldworker Manual.

Table 12: Chi-squared Tests for Difference in Distributions

| Number Eligible | Number of HHs | | Person Number | | | | | | Total | Chi-sq. |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | 1 | 2 | 3 | 4 | 5 | 6 | | |
| 3 | 1 045 | Expected % | 50 | 70 | 80 | | | | | $\chi^2(2)$ |
| | | Expected # | 523 | 732 | 836 | | | | 2 090 | |
| | | Actual # | 845 | 844 | 401 | | | | 2 090 | CV: 5.99 |
| | | Difference | 323 | 113 | -435 | | | | 0 | Test stat: 442.7 |
| 4 | 1 104 | Expected % | 50 | 50 | 50 | 50 | | | | $\chi^2(3)$ |
| | | Expected # | 552 | 552 | 552 | 552 | | | 2 208 | |
| | | Actual # | 738 | 731 | 530 | 209 | | | 2 208 | CV: 7.91 |
| | | Difference | 186 | 179 | -22 | -343 | | | 0 | Test stat: 334.7 |
| 5 | 901 | Expected % | 40 | 50 | 20 | 60 | 30 | | | $\chi^2(4)$ |
| | | Expected # | 360 | 451 | 180 | 541 | 270 | | 1 802 | |
| | | Actual # | 529 | 434 | 476 | 265 | 98 | | 1 802 | CV: 9.45 |
| | | Difference | 169 | -17 | 296 | -276 | -172 | | 0 | Test stat: 815.4 |
| 6 | 590 | Expected % | 50 | 20 | 20 | 30 | 40 | 40 | | $\chi^2(5)$ |
| | | Expected # | 295 | 118 | 118 | 177 | 236 | 236 | 1 180 | |
| | | Actual # | 293 | 214 | 263 | 193 | 140 | 77 | 1 180 | CV: 11.1 |
| | | Difference | -2 | 96 | 145 | 16 | -96 | -159 | 0 | Test stat: 403.9 |

Expected numbers rounded to closest integer.

# B   Leading Digits of Other Monetary Variables

Figure 6: Leading Digit Distributions for Other Monetary Variables



Source: Own calculations using NIDS Wave 1 2008 and Wave 2 2010/2011.

# C NIDS Verification Questionnaire

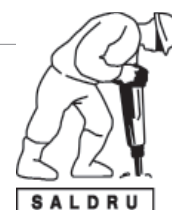Figure 7: NIDS Verification Questionnaire

**INDIVIDUAL QUESTIONNAIRE**

To be used only if MAIN respondent answered **NO**, to either Q11a or b (Was not present or unsure of question details)

| HHID Number: | | PID Number: | |
|---|---|---|---|

| | | Yes/Correct (✐✓) | No/Incorrect (✐✓) |
|---|---|---|---|
| | **INDIVIDUAL QUESTIONS** | | |
| 1. | Were you interviewed in person by an interviewer for the NIDS survey in 2010? | ☐ (Ask Q2a) | ☐ (Ask Q7) |
| 2a. | READ OUT: We would like to confirm some of the questions, to see if we have the correct information on our system. Did the interviewers ask you personally you were born and where you might have lived previously? | ☐ (Ask Q2b) | ☐ (Ask Q2b) |
| 2b. | Where were you born? Interviewer Note: Open UMPC to birth history and confirm where the respondent was born. | ☐ (Ask Q2a) | ☐ (Ask Q2a) |
| 3a. | Did they ask you detailed questions about your mother and father, such as when they were born and what their last work was? | ☐ (Ask Q3b) | ☐ (Ask Q3b) |
| 3b. | When was your mother born? Interviewer Note: Confirm the details on the UMPC Is the answer recorded the same as the answer given by the respondent? | ☐ (Ask Q4a) | ☐ (Ask Q4a) |
| 4a. | Did they ask you questions about your education, such as the grades you completed and any other studies you may have done after school? | ☐ (Ask Q4b) | ☐ (Ask Q4b) |
| 4b. | What is the last school you attended? Interviewer Note: Confirm the details on the UMPC Is the answer recorded the same as the answer given by the respondent? | ☐ (Ask Q5) | ☐ (Ask Q5) |
| 5. | What your employment status at the time you were interviewed? Were you employed, unemployed or self employed? Interviewer Note: Confirm the details on the UMPC Is the answer recorded the same as the answer given by the respondent? | ☐ (Ask Q6) | ☐ (Ask Q6) |
| 6. | Were your measurements taken, i.e. your height, weight, waist and blood pressure? | ☐ | ☐ (Ask Q7) |
| 7. | Interviewer Note: ONLY If No/Incorrect recorded in any of the questions above, READ OUT: We seem to have some missing information. Would it be all right if we completed these sections now of could we come back at another more convenient time to you and other household members? | ☐ (Proceed to individual interview/sections on UMPC dependant on information omitted/incorrect) | ☐ (Ask Q7) |

# southern africa labour and development research unit

The Southern Africa Labour and Development Research Unit (SALDRU) conducts research directed at improving the well-being of South Africa's poor. It was established in 1975. Over the next two decades the unit's research played a central role in documenting the human costs of apartheid. Key projects from this period included the Farm Labour Conference (1976), the Economics of Health Care Conference (1978), and the Second Carnegie Enquiry into Poverty and Development in South Africa (1983-86). At the urging of the African National Congress, from 1992-1994 SALDRU and the World Bank coordinated the Project for Statistics on Living Standards and Development (PSLSD). This project provide baseline data for the implementation of post-apartheid socio-economic policies through South Africa's first non-racial national sample survey.

In the post-apartheid period, SALDRU has continued to gather data and conduct research directed at informing and assessing anti-poverty policy.  In line with its historical contribution, SALDRU's researchers continue to conduct research detailing changing patterns of well-being in South Africa and assessing the impact of government policy on the poor.  Current research work falls into the following research themes: post-apartheid poverty; employment and migration dynamics; family support structures in an era of rapid social change; public works and public infrastructure programmes, financial strategies of the poor; common property resources and the poor.  Key survey projects include the Langeberg Integrated Family Survey (1999), the Khayelitsha/Mitchell's Plain Survey (2000), the ongoing Cape Area Panel Study (2001-) and the Financial Diaries Project.