



Discussion Paper

Estimating survey questionnaire profiles for measurement error risk

The views expressed in this paper are those of the author(s) and do not necessarily reflect the policies of Statistics Netherlands

2016 | 20

**Barry Schouten
Frank Bais
Vera Toepoel**

Content

1. Introduction	4
2. Estimating questionnaire profiles	7
2.1 Item characteristics and questionnaire profiles	7
2.2 Multiple imputation to account for missing coder scores	10
3. A case study: The LISS panel core studies and Dutch Labour Force Survey	12
3.1 The LISS panel data	12
3.2 Results	14
4. Discussion	20
References	21
Appendix A - Questionnaire profiles	23

Summary

Surveys differ in their topics, language, style and design, and, consequently, in their sensitivity to measurement error. Survey literature presents a range of characteristics of survey items that are assumed to be related to the magnitude and frequency of measurement error. In terms of questionnaire design and testing, it would be very useful to have a questionnaire profile that is a summary of the characteristics of the items contained in a questionnaire. This holds especially true in the context of multi-mode surveys where the detection of measurement error is crucial.

The questionnaire profiles may be derived from scores that coders assign to the items in a questionnaire. Given that agreement among coders may be relatively low, as we observe, the number of coders must be large to ensure sufficient precision of the profiles. For multiple surveys, the coding workload may then become infeasible. In this paper, we propose methodology for the estimation of questionnaire profiles when a pool of coders is randomly allocated to a series of surveys. The methodology is based on multiple imputation and applied to eleven general purpose surveys in the Netherlands..

1. Introduction

Measurement error is widely studied in the survey methodology literature, e.g. Alwin and Krosnick (1991), Biemer et al (1991), Fowler (1995) and Tourangeau, Rips and Rasinski (2000), and it is known to be an error that is difficult to measure and predict. Various authors have attempted to develop methodology to predict measurement error from the characteristics of a survey item. The most well-known attempt is the Survey Quality Predictor (SQP) software by Saris and Gallhofer (2007). More recent lists of predictive item characteristics are Schaeffer and Dykema (2010 and 2011), Campanelli et al (2011), Beukenhorst et al (2013) and Bais et al (2015). In this paper, we define and estimate questionnaire profiles, which are summaries of item characteristics over the items in a survey. In order to do so, we employ the scores from a set of coders that independently worked on a range of surveys.

The motivation for the questionnaire profiles comes from the urgent need for a relatively cheap and quick assessment of the overall measurement error risk of a survey, especially in the early (re)design stages of a survey. We are specifically motivated by multi-mode survey (re)designs that are often driven by cost constraints. Budget and time pressure may bring extensive, cognitive questionnaire testing and/or costly experimental studies to assess measurement effects between alternative design into question. In order to decide to what extent such testing and/or experimentation is conducted, an informative but preliminary assessment of measurement error risk is imperative. Questionnaire profiles are, however, not a substitute for in-depth cognitive questionnaire testing. They may form a criterion to decide about the amount and depth of testing in the survey (re)design and may function as a starting point for such tests. Additionally, they may form the incentive to do experimentation and to reserve more time to (re)design a survey. Hence, questionnaire profiles are foremost tools for survey coordinators and management to make decisions about the various (re)design stages, although they may contain valuable information as well for questionnaire designers.

A questionnaire profile summarizes the frequencies of occurrence of a pre-defined set of relevant item characteristics over the items in a survey. What is deemed relevant depends on the context. The motivation for the present study comes from multi-mode survey designs, where mode-specific measurement bias can be unexpectedly high, may slow down redesigns and hamper publication. An example of large biases is given by Schouten, Van den Brakel, Buelens, Van der Laan and Klausch (2013). For this reason, we focus on the most relevant characteristics for mode effects: difficult language in question or in answer categories, question asks for sensitive information, question is sensitive to strong emotions, question is non-central (asks for knowledge that lies outside daily life), and question may be presumed to be a filter question. For a discussion of these characteristics, see for example, Van der Vaart, Van der Zouwen and Zijlstra (1995), Tourangeau and Yan (2007), Kreuter, Presser and Tourangeau (2008), Campanelli et al (2011) and Eckman et al (2014). However, the methodology presented here is not specific to the

selection of item characteristics. Nor is it specific to the purpose for which the item characteristics are used. Recent uses of characteristics are, for example, the explanation of survey response times, see Yan and Tourangeau (2008), Couper and Kreuter (2013) and Olson and Smyth (2015), which are indirectly linked to measurement error. The important message from this paper is that care is needed in estimating and employing the occurrence of the characteristics.

Bais et al (2015) show that intercoder agreement for item characteristics can be low, even for motivated, trained and experienced coders. They conclude that disagreement can only be resolved by restrictive definitions of the characteristics or by very time-consuming item-by-item decisions to reach a consensus. How to summarize scores on characteristics into profiles is, therefore, not as straightforward as it may seem. One may simply estimate the average number of coders that scored an item as having the property and then take the mean of these averages over all items. Doing so, a lot of information about the coding (and hence the characteristics of the questionnaire profiles) is lost. Instead, we propose to construct a distribution over the items in a survey that reflects the variability between coders. We believe this profile to be more useful than simple means, because all information is maintained but structured.

The precision of the resulting profile depends on the number of coders; the more coders, the more precise the estimated frequencies. Since the list of relevant item characteristics may be long and since coding is a time-consuming, and, hence, costly exercise, it is usually infeasible to let a large number of coders work on all surveys. Consequently, it becomes attractive to construct efficient coding and imputation schemes. It is important to stress that the coding of survey items concerns the wording and format of questions and answer categories, but not the actual answers given by the survey respondents.

There is a vast literature on intercoder agreement, e.g. Cohen (1960), Fleiss (1971) and Shoukri (2010), and various measures have been developed to evaluate agreement. The most well-known is Cohen's Kappa and variants of this measure. This literature focusses on reliability of coding, i.e., when applied to our setting, it assumes that the same coder may give different scores for the same item when replicated at different times and in different circumstances. Here, we assume that coders worked conscientiously and reliability is a negligible problem. The focus is on the systematic differences between coders, i.e. the validity of the scores. Given the findings of Bais et al (2015), we believe that the systematic differences between coders dominate the random differences. For this reason, we do not consider the more traditional measures of intercoder agreement.

As a useful by-product of the study, we give the questionnaire profiles of eleven multi-purpose surveys in the Netherlands scored by a group of eight coders. These surveys are conducted in a comparable form in many countries.

This paper reads as follows: In section 2, we define questionnaire profiles and we propose an estimation strategy based on randomly allocated coders. In section 3, we

apply the estimation strategy to the eleven general purpose surveys. In section 4, we end with a discussion..

2. Estimating questionnaire profiles

In this section, we introduce a number of item characteristics, define questionnaire profiles based on these characteristics, and we construct an estimation strategy.

2.1 Item characteristics and questionnaire profiles

In Bais et al (2015) an extensive list of item characteristics is presented. This list is derived from Saris and Gallhofer (2007), Campanelli et al (2011) and Beukenhorst et al (2013). In this paper, we consider a subset of six characteristics that are taken from this list:

- difficult language in question: the question contains one or more difficult words or a complicated sentence structure;
- difficult language in answer: the answer categories contain one or more difficult words, or require a complicated cognitive action (e.g. sliding bars or abstract visual representation);
- non-centrality: the question asks for knowledge or experience that lies outside daily life of the average respondent;
- sensitive to emotions: the question may arouse negative emotions like anger, distress, sorrow or despair;
- sensitive information: the question asks for information that is viewed as sensitive by the average respondent;
- presumed filter question: the average respondent believes that the question is a filter question and some of the answer categories will avoid follow-up questions;

We view the selected item characteristics as very influential on measurement error, and more specifically, as relevant to mode effects. However, in this paper the characteristics serve mostly as examples to demonstrate the derivation of profiles; the methodology set out in this paper can be applied to any set of item characteristics.

Suppose one would like to code all items in a series of S surveys on a given item characteristic by human coders. A group of M coders is randomly assigned to the surveys, and each survey gets assigned A coders. In other words, each coder gets on average AS/M surveys to work on. We, first, assume that $A = M$ and all coders do all surveys.

Let $C_{s,i,m}$ be the 0-1 score of coder m on item i in survey s for a certain item characteristic. The surveys are labeled $s = 1, 2, \dots, S$, the coders are labeled

$m = 1, 2, \dots, M$, and the items within surveys are labeled $i = 1, 2, \dots, I_s$. Let $I = \sum_{s=1}^S I_s$ be the total number of items.

We assume that the scores for a given set of items and coders $\{C_{s,i,m}\}_{1 \leq s \leq S, 1 \leq i \leq I_s, 1 \leq m \leq M}$ are independent and follow Bernoulli distributions with parameters $p_{s,i,m}$. We do not model the selection of items and the clustering of items within surveys, but assume these as given. However, we do view coders as selected from a superpopulation of coders, i.e. there is an underlying $p_{s,i}$ of interest, which may be viewed as the expected coder item probability for item i in survey s . Furthermore, we are interested in the average over the items in a survey p_s , i.e.

$$p_s = \frac{1}{I_s} \sum_{i=1}^{I_s} p_{s,i}.$$

Our basic modelling assumption is that for any $M > 1$, for any pair of survey items (s, i) and (\tilde{s}, j) , and for any vector $(c_1, c_2, \dots, c_M) \in \{0, 1\}^M$, it holds that

$$P[C_{s,i,m} = c_m | C_{s,i,1} = c_1, \dots, C_{s,i,m-1} = c_{m-1}, C_{s,i,m+1} = c_{m+1}, \dots, C_{s,i,M} = c_M] = P[C_{\tilde{s},j,m} = c_m | C_{\tilde{s},j,1} = c_1, \dots, C_{\tilde{s},j,m-1} = c_{m-1}, C_{\tilde{s},j,m+1} = c_{m+1}, \dots, C_{\tilde{s},j,M} = c_M] \quad (1)$$

The model assumption (1) can be formalized as follows: Item probabilities $\{p_{s,i}\}_{1 \leq s \leq S, 1 \leq i \leq I_s}$ are drawn independently from a distribution with scope $[0, 1]$, say G . Conditional on the item probability, the coder item probabilities $\{p_{s,i,m}\}_{1 \leq m \leq M}$ are drawn independently from the same distribution, i.e. $(p_{s,i,m} | p_{s,i} = p) \triangleq F_p$, where the F_p form a class of distribution functions with expectation p . Within a coder, the $p_{s,i,m}$ are allowed to be dependent, i.e. there may be a coder “fixed effect”. Since we have a given set of surveys and items, we do not further parameterize or attempt to explicitly model the coder item probability distributions, but rather resort to empirical distributions. Model assumption (1) directly allows for imputation schemes.

The set of p_s over the multiple item characteristics may be viewed as a profile of a questionnaire. However, the p_s do not express the amount of coder (dis)agreement, i.e. the variability in coder probabilities. Two surveys may have the same average p_s over their items, but may differ strongly in terms of coder consensus. This difference is important in judging if and how survey designers should deal with measurement error risk. For this reason, we include coder variability. We let $f_p(x), x \in [0, 1]$, be the probability density function for distribution F_p . Now, we define the survey average

$$P_s(x) = \frac{1}{I_s} \sum_{i=1}^{I_s} f_{p_{s,i}}(x), \quad (2)$$

as the questionnaire profile for an item characteristic. We have that $\int_0^1 P_s(x) dx = 1$ and $P_s(x)$ may be interpreted as the relative proportion of items in survey s that has the item characteristic according to a fraction x of the coders. If coders would fully agree, then $P_s(x) = 0$ for $0 < x < 1$.

The set of functions P_s over all item characteristics we call the full questionnaire profile or simply the questionnaire profile.

The $p_{s,i}$, p_s and $P_s(x)$ are unknown and they are to be estimated. The obvious estimators are

$$\hat{p}_{s,i} = \frac{1}{M} \sum_{m=1}^M C_{s,i,m}, \quad (3)$$

$$\hat{p}_s = \frac{1}{I_s M} \sum_{m=1}^M \sum_{i=1}^{I_s} C_{s,i,m}, \quad (4)$$

and the empirical density function for (3) is

$$\hat{P}_s(x) = M \frac{1}{I_s} \sum_{i=1}^{I_s} A_{s,i}(x), \quad (5)$$

where $A_{s,i}(x)$ is the observed 0-1 indicator for the event that a fraction x of the coders scored the characteristic, i.e. that $\frac{1}{M} \sum_{m=1}^M C_{s,i,m} = x \in \left\{0, \frac{1}{M}, \frac{2}{M}, \dots, 1\right\}$. For $x \notin \left\{0, \frac{1}{M}, \frac{2}{M}, \dots, 1\right\}$, we simply interpolate. The multiplication by M in (5) results from the bin size $= \frac{1}{M}$.

Although, we are not specifically interested in the item probabilities of individual coders, we do estimate the coder average score, denoted by $\theta_m = \frac{1}{I_s} \sum_{i=1}^{I_s} p_{s,i,m}$, using

$$\hat{\theta}_m = \frac{1}{I_s} \sum_{i=1}^{I_s} C_{s,i,m}. \quad (6)$$

We do this, because in practice the set of available coders will change only gradually over time, with new coders starting and former coders quitting at a low frequency, so that we may need to monitor and maintain the individual average coder probabilities.

The standard errors for estimators (3), $\sigma_{s,i}$, (4), σ_s , (5), $\tau_s(x)$, and (6), σ_m , cannot easily be estimated. We, therefore, use resampling methods to estimate them. Since we observe only one score per coder per item, we cannot account for lack of coder reliability, i.e. due to $p_{s,i,m}$ between 0 or 1. In estimating standard errors for the $\hat{\theta}_m$, we ignore this variability.

In practice, it will usually be too costly to score the items of all surveys by all coders on all item characteristics. In the case study to this paper in section 3, the coders scored on average 35 items per hour on the set of characteristics, and coding all items would have cost roughly 70 hours per coder. However, as was concluded in Bais et al (2015), the coder average scores may vary greatly and, as a consequence, multiple coders are needed to obtain a precise estimate of the item and survey probabilities. This leads to a trade-off between coder costs and coding precision. In the next section, we will show how the various parameters of this section can be estimated using multiple imputation, when part of the coder scores are missing by design.

2.2 Multiple imputation to account for missing coder scores

Instead of coding all surveys, assume the coders work only on a random subset of surveys, i.e. $A < M$. Assume, furthermore, that the coders differ in their maximal workload. Let S_m be the number of surveys that coder m can work on, i.e. $\sum_{m=1}^M S_m = AS$. Let $U_{m,s}$ be the 0-1 indicator for the allocation of survey s to coder m . We have that $\sum_{s=1}^S U_{m,s} = S_m$ and $P[U_{m,s} = 1] = S_m/S$.

We use multiple imputation to fill out the missing scores of coders and to estimate the questionnaire profiles. As an important by-product, we estimate the standard errors following from the missing item scores and the standard errors following from the selection of coders.

The algorithm is:

1. Construct an imputation scheme for the missing surveys;
2. Repeat B times the following steps:
 - a. Perform a (random) imputation given the scheme of step 1;
 - b. Based on the imputed data set, estimate the item probability, $\hat{p}_{s,i}^b$, the survey probability, \hat{p}_s^b , the coder average score, $\hat{\theta}_m^b$, and the questionnaire profile, $\hat{P}_s^b(x)$.
 - c. Based on the imputed data set, estimate the standard errors for the item probabilities, $\hat{\sigma}_{s,i}^b$, using (5), and for the survey probabilities, $\hat{\sigma}_s^b$, and the questionnaire profiles, $\hat{t}_s^b(x)$, using bootstrap;
3. Estimate the mean of the item probabilities, $\hat{p}_{s,i} = \sum_{b=1}^B \hat{p}_{s,i}^b$, survey probabilities, $\hat{p}_s = \sum_{b=1}^B \hat{p}_s^b$, questionnaire profiles, $\hat{P}_s(x) = \sum_{b=1}^B \hat{P}_s^b(x)$, and coder average scores, $\hat{\theta}_m = \sum_{b=1}^B \hat{\theta}_m^b$;
4. Estimate the mean of the standard errors of the item probabilities, $\hat{\sigma}_{s,i}^W = \sum_{b=1}^B \hat{\sigma}_{s,i}^b$, the survey probabilities, $\hat{\sigma}_s^W = \sum_{b=1}^B \hat{\sigma}_s^b$, and the questionnaire profiles, $\hat{t}_s^W(x) = \sum_{b=1}^B \hat{t}_s^b(x)$;
5. Estimate the standard deviation of the item probabilities,

$$\hat{\sigma}_{s,i}^B = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (\hat{p}_{s,i}^b - \hat{p}_{s,i})^2}, \text{ the survey probabilities,}$$

$$\hat{\sigma}_s^B = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (\hat{p}_s^b - \hat{p}_s)^2}, \text{ the questionnaire profiles,}$$

$$\hat{t}_s^B(x) = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (\hat{P}_s^b(x) - \hat{P}_s(x))^2}, \text{ and the coder average scores, } \hat{\sigma}_m^B =$$

$$\sqrt{\frac{1}{B-1} \sum_{b=1}^B (\hat{\theta}_m^b - \hat{\theta}_m)^2};$$

6. Estimate the total standard error using Rubin's rules,

$$\hat{\sigma}_{s,i}^T = \sqrt{(\hat{\sigma}_{s,i}^W)^2 + (1 + \frac{1}{M})(\hat{\sigma}_{s,i}^B)^2}, \hat{\sigma}_s^T = \sqrt{(\hat{\sigma}_s^W)^2 + (1 + \frac{1}{M})(\hat{\sigma}_s^B)^2}, \hat{t}_s^T(x) =$$

$$\sqrt{(\hat{t}_s^W(x))^2 + (1 + \frac{1}{M})(\hat{t}_s^B(x))^2}, \text{ and } \hat{\sigma}_m^T = \hat{\sigma}_m^B;$$

Some side remarks are in place: The variances of the estimators for complete data sets are called within variances, whereas the variances over the imputed data sets are called between variances. For this reason, the superscripts "W" and "B" are used in steps 4 and 5, respectively. The coder average scores have no standard error in a

complete data set. Hence, standard errors arise only from the missing surveys; the within variances are zero by definition. A usual choice for the number of imputed data sets is $B = 10$. For the bootstrap, we use 1000 replications per imputed data set.

The algorithm produces unbiased estimates under four conditions. First, the coders work independently from each other. Second, the coders score the items consistently, i.e. they score each item as isolated from the other items. Third, the surveys need to be allocated randomly to coders. Fourth, in the imputation the matching property holds, i.e. for each missing score combination an observed score combination exists. The first two conditions are about the coders themselves. The third condition implies a missing-completely-at-random mechanism for the missing scores. The fourth condition implies that predictors exist for each missing score. The first three conditions are under control and we assume they hold. However, the fourth condition does not hold necessarily and depends on the imputation scheme. For example, when the scores of nine coders are used to impute the scores of a tenth coder, then it is likely that part of the possible combinations of nine 0-1 scores did not occur in the data set. When such a combination occurs for a missing score on the tenth coder, then there is no observation that can be used to predict that missing score. Survey items are not randomly clustered within surveys which increases the risk that the matching condition does not hold. In practice, therefore, parsimony is needed in the imputation scheme.

The imputation scheme is the most complicated part of the multiple imputation algorithm. The scheme describes the order in which surveys are imputed and the subset of coders that are used to predict the missing score. In all generality, it cannot be stated beforehand what is the most efficient scheme for a given allocation. The reduction in multiple imputation standard error depends on the correlation between the scores of coders, i.e. their mutual agreement, the amount of overlap in surveys between coders, and the amount of non-overlap between coders. When two coders always agree, then they form an ideal couple to impute each other's missing scores. When the agreement is the same between various pairs of coders, then the amount of overlap in items determines the order in which imputations are made; the more overlap the better. Finally, when agreement is the same and overlap is the same, it is the non-overlap that counts. As mentioned before, the matching condition warns against imputation using the scores of all available coders. In our case study, we use a maximum of three coders to impute scores of other coders.

Consider the example in table 1 with eight surveys and five coders. Per survey, three coders are assigned. Coders 1 to 3 can do six surveys, whereas coder 4 and 5 can only do three surveys. In total, 24 coded surveys are produced and 16 coded surveys are missing. Table 1 shows one possible realization of coder allocation. Given that coders 1 to 3 worked on six surveys, they show the largest overlap, and it is natural to impute their missing surveys first. Coders 1 and 2 and coders 1 and 3 worked on five surveys simultaneously. With no knowledge about their agreement, it is an option to impute the missing (coder, survey) cells (1,7), (2,1), (1,8) and (3,6). From there on, cell (2,8) can be imputed. Coders 1, 2 and 3 scored four surveys simultaneously (surveys 2 to 5) and a second option is that pairs of coders are used to first impute

(2,1) and (3,6), and then to proceed to the cells (1,7), (1,8) and (2,8). In either option, the (imputed) scores of the first three coders can then be used to impute surveys for coders 4 and 5. However, depending on the agreement between coders, the optimal scheme could be different. For example, when coder 4 and 5 agreed on all items, then it may be more efficient to impute the missing cells (4,6) and (5,1) with each other's scores.

Table 1: Example of an allocation of five coders to eight surveys. Grey cells are scores by the coder. White cells are missing survey scores.

Coder	Survey							
	1	2	3	4	5	6	7	8
1	Grey	Grey	Grey	Grey	Grey	Grey	White	White
2	White	Grey	Grey	Grey	Grey	Grey	Grey	White
3	Grey	Grey	Grey	Grey	Grey	White	White	Grey
4	Grey	White	White	White	White	White	Grey	Grey
5	White	White	White	White	White	Grey	Grey	Grey

Another feature of the imputation scheme that has not been mentioned is that imputation needs to be performed per item characteristic. In general, the optimal scheme will be different for each characteristic, but different schemes may not be computationally attractive.

In section 3, we discuss an imputation scheme for a case study.

3. A case study: The LISS panel core studies and Dutch Labour Force Survey

The estimation strategy of section 2.2 is applied to core studies of the Dutch LISS-panel and the Dutch Labor Force Survey (LFS). We, first, describe the data, and, then, present results.

3.1 The LISS panel data

The data that we will use, come from ten core study surveys that were all administered in the LISS-panel of CentERdata, Tilburg University, in the years 2008 - 2014. To these ten surveys we added the Labor Force Survey, an on-going monthly survey that is administered by Statistics Netherlands. LISS is short for Longitudinal Internet studies for the Social Sciences and is a government-funded Web panel based on a probability sample drawn by Statistics Netherlands. The panel was established in 2007 and is now running for eight years. It is, generally, considered a high-quality

panel, because of the extensive recruiting and refreshment throughout the panel. The panel has roughly 8000 panel members. Table 2 gives a short list of the surveys and topics contained in the surveys, and the number of items per survey. In total the number of items is 2470. From 2008 to 2014, all core study surveys have been administered annually, except for Assets (AS) and Income (IN). The questionnaires of the last available wave were used for coding.

The coding of the surveys was prepared in four steps: First a preliminary set of item characteristic definitions was made, second this set of definitions was applied by all coders to a small but diffuse set of items in a pilot study, third the definitions were discussed and revised based on the pilot study findings, and fourth they were applied to all items. The coders are two experts from Statistics Netherlands' cognitive questionnaire lab, three experts from CentERdata's questionnaire design department and three survey methodology researchers from Statistics Netherlands and Utrecht University. In total, eight coders were available. To each survey, three coders were randomly allocated. However, their availability in terms of hours was not equal, as is usually the case for coding exercises, so that the number of surveys per coder is very different; one coder did all surveys, one did seven, one did five, two did two and three did one. The coding exercise was time consuming as, apart from the six characteristics on which we focus attention in this paper, ten more characteristics were coded (Bais et al 2015). Hence, a total of 2470 items was coded on 16 characteristics. For this reason, it was not possible to let all coders do all surveys and avoid missing data. Table 3 presents the allocation to the surveys. We refer to Bais et al (2015) for details.

Table 2: The surveys in the case study.

<i>Label</i>	<i>Survey</i>	<i>Topics of the Content</i>	<i>I_s</i>
AS	Assets	Assets, property and investment;	50
FA	Family and Household	Household composition and family relations;	73
HE	Health	Health and well-being;	286
HO	Housing	Housing and household;	409
IN	Income	Employment, labor and retirement, income, social security and welfare;	243
LFS	Labor Force Survey	Education, employment and labor;	200
PE	Personality	Personality traits;	148
PO	Politics and Values	Politics, social attitudes and values, attitudes towards surveys;	71
RE	Religion and Ethnicity	Religion, social stratification and groupings;	396
SO	Social Integration and Leisure	Communication, language and media, leisure, recreation and culture, social behavior, travel and transport;	471
WO	Work and Schooling	Education, employment, labor and retirement;	123

Table 3: The allocation of coders to the case study surveys. The light grey blocks are omitted in the sensitivity analysis.

Coder	AS	FA	HE	HO	IN	LFS	PE	PO	RE	SO	WO
1											
2											
3											
4											
5											
6											
7											
8											

In order to construct an imputation scheme, we looked at the correlations between the scores of coders that worked on the same surveys. These correlations turn out to be relatively low, in general. Table 4 shows the correlations for item characteristic “Sensitive information” for the pairs of coders that worked simultaneously on at least one survey. Given the low correlations and given the practicality of using the same scheme for all item characteristics, we decided to use the amount of overlap as the criterion to build the imputation scheme.

Table 4: Correlations between the scores on characteristic “Sensitive information” for pairs of coders that worked on at least one survey simultaneously.

	2	3	4	5	6	7	8
1	0.19	0.23	0.12	-0.02	0.47	0.17	0.14
2	-	0.31	NA	0.09	NA	NA	NA
3	-	-	0.22	0.23	0.59	0.63	0.02

The imputation scheme we used is

1. Impute the missing survey of coder 3 using the scores of coder 1;
2. Impute the missing surveys of coder 2 using the scores of coders 1 and 3;
3. Impute the missing surveys of all other coders using the scores of coders 1, 2 and 3;

For some of the coders, the majority of surveys is imputed, which has a strong impact on standard errors, as we will see in the next section. This points to the inefficient allocation of coders in table 4 that is the result of the strongly varying maximal workloads.

3.2 Results

We estimated the proportion of coders that would indicate an item to have a characteristic. We did this for all 2470 items and for all six characteristics: difficult language in question (DLQ), difficult language in answer (DLA), non-centrality (CENT), sensitive to emotions (EMO), sensitive information (SENS) and presumed filter

question (FILT). Furthermore, we aggregated the proportions over surveys and over coders, and we estimated questionnaire profiles per characteristic.

Table 5 gives the estimated survey probabilities \hat{p}_s for the six item characteristics. For each estimate, two standard errors are given; the within standard error corresponding to a complete coder data set and the total standard error that also includes the imputation standard error. The within standard errors are large and much larger than the between standard error resulting from incomplete coder data. This points to a large uncertainty that is due to the coders. We will return to this conclusion when we discuss the coder probabilities. The surveys differ substantially in their probabilities. For instance, core study Assets scores highly on many of the characteristics and may be considered a survey that is susceptible to measurement error. Core study Personality (PE) has low scores and may be considered less prone to error. However, as mentioned, standard errors are large and confidence intervals are wide. Nevertheless, differences between surveys frequently test as significant because there is a strong covariance within coders over surveys. As a result, standard errors of differences between surveys are similar in magnitude to standard errors per survey.

Table 5: Estimated probabilities per survey and item characteristic. Within standard errors and total standard errors within brackets.

	AS	FA	HE	HO	IN	LFS	PO	PE	RE	SO	WO
DLQ	32% (5.6) (6.2)	11% (1.4) (1.5)	13% (2.3) (2.4)	18% (2.7) (2.8)	23% (4.0) (4.1)	15% (2.3) (2.5)	18% (3.2) (3.3)	8% (1.3) (1.3)	21% (5.2) (5.3)	11% (1.7) (1.8)	21% (5.1) (5.1)
DLA	2% (1.2) (1.2)	7% (1.9) (1.9)	1% (0.3) (0.3)	4% (1.0) (1.2)	0% (0.3) (0.4)	7% (2.1) (2.1)	5% (1.6) (1.6)	2% (0.6) (0.6)	8% (2.1) (2.3)	2% (0.5) (0.5)	3% (0.7) (0.7)
CENT	33% (9.9) (10.1)	10% (5.7) (5.7)	13% (4.1) (4.2)	21% (5.0) (5.0)	28% (7.9) (8.0)	11% (6.9) (7.1)	21% (5.9) (6.0)	9% (3.1) (3.2)	23% (6.9) (7.1)	18% (5.7) (5.7)	20% (5.8) (5.8)
EMO	15% (6.8) (7.0)	15% (5.4) (5.4)	12% (5.9) (5.9)	11% (5.6) (5.7)	15% (5.8) (5.8)	10% (5.9) (5.9)	21% (7.9) (8.0)	18% (5.8) (5.8)	14% (5.6) (5.7)	10% (5.6) (5.6)	13% (5.5) (5.5)
SENS	58% (10.1) (10.4)	23% (5.0) (5.1)	34% (9.3) (9.4)	34% (7.8) (7.9)	44% (10.8) (10.8)	16% (4.7) (4.8)	35% (8.7) (8.7)	18% (5.4) (5.4)	42% (10.2) (10.4)	29% (7.7) (7.7)	27% (6.5) (6.5)
FILT	35% (4.5) (4.8)	29% (4.0) (4.1)	25% (4.1) (4.3)	25% (4.9) (4.9)	26% (3.4) (3.4)	30% (5.9) (6.0)	16% (3.3) (3.4)	11% (3.6) (3.7)	18% (3.0) (3.4)	30% (4.3) (4.4)	28% (4.5) (4.6)

Table 6 gives the estimated coder probabilities $\hat{\theta}_m$ for the six item characteristics. The standard errors are small. For coder 1, they are zero by definition as all surveys were allocated to this coder. The estimates confirm the large standard errors of table 6; there is a great variability in scores between coders. It must be noted that the estimated probabilities may be biased for coders that did only a small number of

surveys, as noted in section 2.2, despite the random allocation of coders to surveys. This fallacy appears when some combinations of scores over coders are absent in some of the surveys. For a number of coder-survey combinations, this occurred. These combinations are colored light grey in table 5. Remarkably, coders 4 and 5 did not score any of the items as being non-central, and coder 5 did not score any of the items as being sensitive to emotions as well. As a result, for these item characteristics, these coders have an estimated probability equal to zero.

Table 6: Estimated probabilities per coder and item characteristic. Standard errors within brackets. Light grey cells correspond to coder-survey combinations where some of the scores of other coders were absent.

	1	2	3	4	5	6	7	8
DLQ	32% (0.0)	8% (0.6)	13% (0.1)	12% (0.4)	17% (0.7)	13% (1.3)	11% (0.8)	22% (1.5)
DLA	7% (0.0)	1% (0.2)	4% (0.1)	0.4% (0.1)	3% (0.3)	3% (0.3)	3% (0.2)	6% (1.2)
CENT	31% (0.0)	5% (0.3)	25% (0.2)	0% (0.0)	0% (0.0)	27% (0.8)	35% (1.0)	16% (0.8)
EMO	7% (0.0)	0% (0.1)	19% (0.1)	11% (0.6)	0% (0.0)	8% (0.5)	52% (0.9)	13% (0.4)
SENS	25% (0.0)	5% (0.3)	35% (0.2)	9% (0.4)	22% (1.8)	59% (1.6)	57% (0.9)	27% (0.9)
FILT	20% (0.0)	37% (0.7)	25% (0.1)	32% (0.7)	34% (1.6)	30% (1.8)	27% (1.5)	2% (0.3)

Figures 1 and 2 depict the estimated questionnaire profiles $\hat{P}_s(x)$ for all item characteristics for, respectively, the LFS and core study Politics and values (PO). The profiles contain symmetric 95% confidence intervals (in grey) based on a normal approximation. For example, 55% of the LFS items are estimated to be free of complex language in the question according to all coders, and 30% of the LFS items according to all but one coder. For presumed filter question, these two percentages are 20% and 22%, respectively. The ensemble of profiles per item is viewed as the questionnaire profile.

Despite the imprecision due to the coder variability, the profiles give a useful picture of the survey questionnaire. For the LFS, the fraction of items that a substantial amount of coders, say at least 3 out of 8, would score as having the characteristic is only present for difficult language (DLQ) in question and presumed filter question (FILT). For core study Politics and values (PO), four of the characteristics show such fractions: difficult language in question (DLQ), non-centrality (CENT), sensitive to emotions (EMO), and sensitive information (SENS). Hence, the two surveys clearly have different profiles. Appendix A contains the profiles of the other nine surveys. To give some indication of the impact of the estimation strategy, table 7 shows observed and estimated probabilities per characteristic over all surveys. Some of the probabilities were lowered, like difficult language in question (DLQ), others were lifted, like sensitive information (SENS).

Figure 1: Questionnaire profile for the Labor Force Survey (LFS). From left to right and from top to bottom: difficult language question, difficult language answer, centrality, emotional loading, sensitive information and presumed filter question. 95% confidence intervals based on a normal approximation are given in light grey.

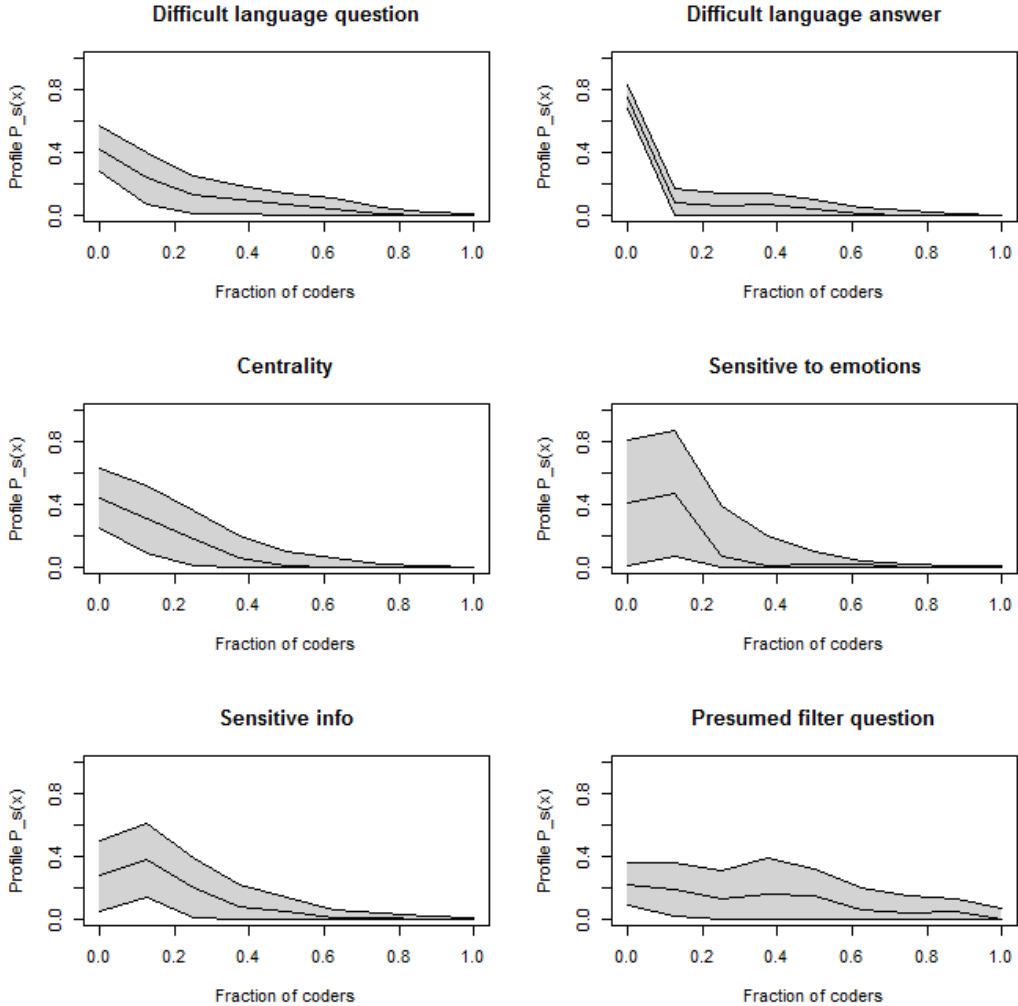


Figure 2: Questionnaire profile for core study Politics and values (PO). From left to right and from top to bottom: difficult language question, difficult language answer, centrality, emotional loading, sensitive information and presumed filter question. 95% confidence intervals based on a normal approximation are given in light grey.

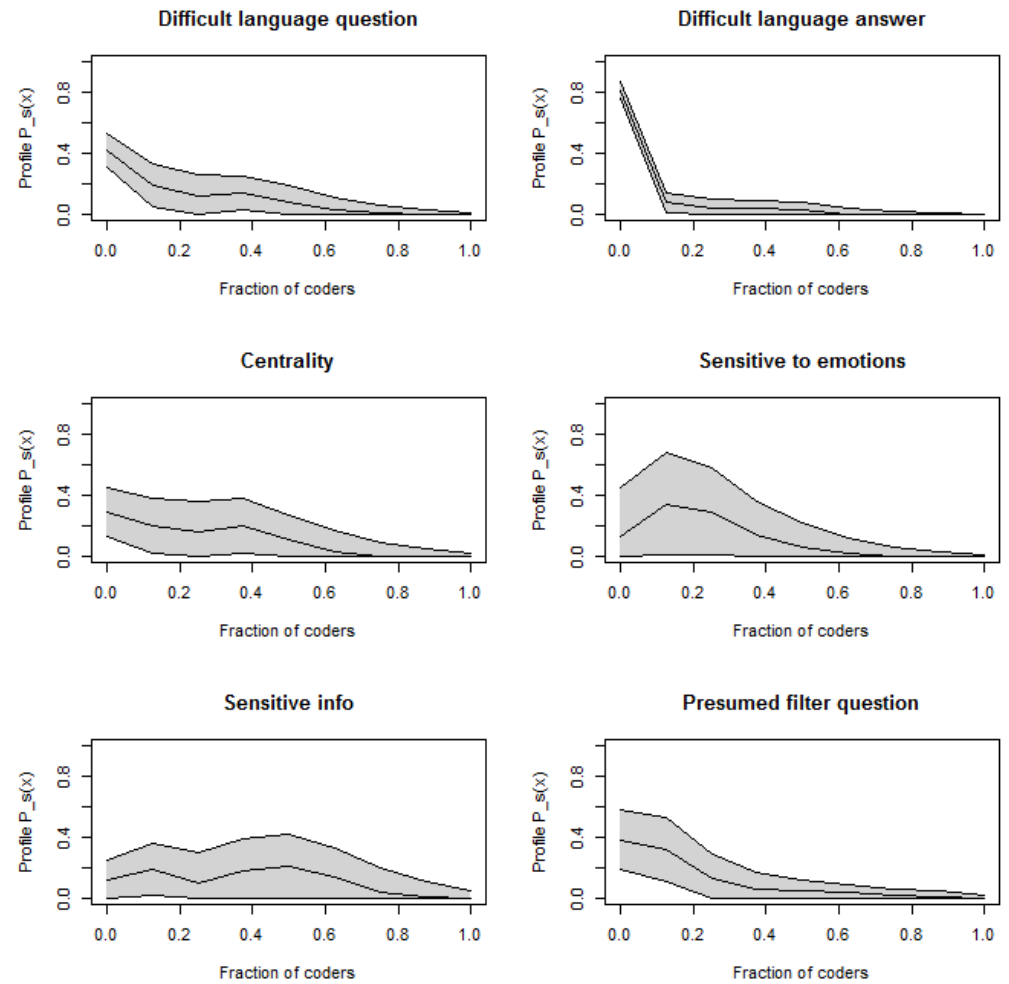


Table 7: Observed and estimated probabilities per item characteristic.

	DLQ	DLA	CENT	EMO	SENS	FILT
Observed	19%	4%	21%	12%	25%	26%
Estimated	16%	3%	17%	14%	30%	26%

In this study, we cannot compare to a complete case analysis. To get some sense of the robustness of the estimation strategy, we omitted five blocks for the coders that worked most of the surveys, the light grey blocks in table 3. Table 8 shows the change in estimated survey probabilities \hat{p}_s relative to all coder data in table 5. For the surveys where a block was omitted, estimates sometimes changed considerably, especially for difficult language in the question (DLQ) and sensitive information (SENS). Four estimates lie outside the original 95% confidence intervals. Hence, estimates are sensitive to the coder allocation scheme and it is advisable that workload is more evenly spread than was done in our study.

Table 8: Differences in estimated survey probabilities when omitting blocks for coders 1 and 3. Underlined differences correspond to values outside 95% confidence intervals.

	AS	FA	HE	HO	IN	LFS	PO	PE	RE	SO	WO
DLQ	<u>-16%</u>	1%	-7%	-10%	0%	-1%	0%	1%	0%	1%	0%
DLA	0%	0%	0%	-3%	0%	0%	0%	0%	0%	0%	0%
CENT	-7%	3%	-6%	-5%	-4%	0%	-4%	0%	-1%	0%	-2%
EMO	0%	0%	-3%	-1%	-1%	1%	-2%	0%	-1%	0%	0%
SENS	-14%	5%	-15%	-6%	-12%	1%	-1%	1%	-6%	0%	0%
FILT	-5%	-2%	-2%	-8%	0%	0%	1%	1%	1%	0%	-1%

For each of the 2470 items, an individual estimate $\hat{p}_{s,i}$ is available per item characteristic. These are not shown (of course), but in future research they will be added as explanatory variables in multi-level models to investigate answering behavior on single items, e.g. social desirable answering, acquiescence, underreporting and don-not-know answers.

4. Discussion

We present methodological tools for a relatively inexpensive and fast preliminary assessment of measurement error risk in surveys. Application of these tools may trigger and inform in-depth cognitive testing and/or experimentation in early (re)design stages. All items of a series of surveys are coded on characteristics that are assumed to be relevant to measurement error. Each survey is assumed to be handled by a limited number of coders. By estimating missing coding data for all other coders for each survey, questionnaire profiles can be constructed.

In our case study, we focused on six item characteristics that are selected for their relevance to mode-specific measurement errors: Difficult language in question, difficult language in answer, risk of non-centrality, sensitive to emotions, sensitive information, and presumed filter question. For each characteristic, the questionnaire profiles for 11 surveys showed for what percentage of all items the characteristic would be present. For instance, the questionnaire profile for the Labor Force Survey showed that the characteristics difficult language in question and presumed filter question appeared to be present for relatively many items according to multiple coders, while the other characteristics did not. This implies that there may be a measurement risk coming from inability of respondents to answer questions or motivated underreporting, a risk that in practice may be mediated by the assistance of interviewers. When the questionnaire profile of Politics and Values survey is inspected, it follows that the characteristics sensitive to emotions, and sensitive information are relatively present. These characteristics point at measurement risk due to socially desirable answering, which may be stronger in the presence of interviewers. In sum, the questionnaire profiles can be used as a starting point for further evaluation.

When using questionnaire profiles specifically as a basis to further investigate measurement error, however, caution is urged for two reasons. First, we observed a large variability in the assigned codes between coders. Some coders were, generally, conservative in coding a characteristic as present, while other coders were generally liberal in doing so. This variability results in large within standard errors accompanying the estimated probabilities and points to a large uncertainty in the judgment of the presence of the characteristics. Second, the estimated probabilities and profiles may be biased due to a selective clustering of items within surveys. Coders may only work on the items of a restricted amount of surveys, and, in the case of selective clustering, may have coded only a small number of items with certain characteristics. Our case study clearly had a suboptimal design in that two coders only coded the items of two surveys and three coders coded the items of only one survey. In general, however, our imputation method is a useful extension of the method with only actual coding data, giving a more informative estimation of a questionnaire profile.

Apart from the methodology, the application in this paper may be relevant to questionnaire designers; the surveys included in the study are general purpose surveys with topics that are used in many countries and in many settings (cross-sectional or panel). Other questionnaire designers may perform similar exercises and compare their profiles to ours. We plan to extend the list of surveys, but to maintain the necessary overlap in coders.

It has yet to be shown that questionnaire profiles indeed help identifying measurement error risk. For this purpose, prospectively or retrospectively questionnaire profiles should be linked to surveys where measurement error is assessed through validation data, paradata and/or re-interview. Obviously, the coders should prepare the questionnaire profiles independently of researchers analysing the measurement error properties. Strongest evidence for the utility of the profiles may come from coding surveys where strong measurement effects have been found. Future study is directed at linking profiles to in-depth measurement error studies.

Future research may focus also on optimal coder allocation schemes and ways to make trade-offs between coding hours and coding accuracy. The relatively low agreement between coders, and the resulting uncertainty in item characteristic probabilities is a reason for concern and further research. Given a specified accuracy of questionnaire profiles, the number of coders must be larger when agreement is smaller. Future research may attempt to improve agreement while maintaining relevance in the definition of the item characteristics. In the case study of this paper, standard errors were relatively large, so that normal approximation is invalid and confidence intervals had to be cut off at zero. More efficient coding schemes may remove the need to have better approximations. Nonetheless, future research may address this methodological issue.

Acknowledgements

The authors like to thank Corrie Vis, Edith de Leeuw, Joop Hox, Judit Arends, Mattijn Morren, Natalia Kieruj, Peter Lugtig, Rachel Vis and Salima Douhou for their help in preparing a coding scheme and in coding the LISS panel core studies.

References

Alwin, D.F., Krosnick, J.A. (1991), "The Reliability of Survey Attitude Measurement: The Influence of Question and Respondent Attributes", *Sociological Methods and Research*, 20, 138 – 181.

Bais, F., Schouten, B., Lugtig, P., Toepoel, V., Arends-Toth, J., Douhou, S., Kieruj, N., Morren, M., Vis, C. (2015), "Can Survey Item Characteristics Relevant to Mode-Specific Measurement Error be coded Reliably?", Discussion paper 2015xx, Statistics Netherlands, available at www.cbs.nl .

Beukenhorst, D., Buelens, B., Engelen, F., Van der Laan, J., Meertens, V. and Schouten, B. (2013), "The Impact of Survey Item Characteristics on Mode-Specific Measurement Bias in the Crime Victimization Survey", Discussion paper 201416. Statistics Netherlands, The Hague, available at www.cbs.nl.

Biemer, P.P., Groves, R.M., Lyberg, L.E., Mathiowetz, N.A., Sudman, S. (1991), *Measurement Error in Surveys*, New Jersey: John Wiley & Sons.

Campanelli, P., Nicolaas, G., Jäckle, A., Lynn, P., Hope, S., Blake, M., Gray, M. (2011), "A Classification of Question Characteristics relevant to Measurement (Error) and Consequently Important for Mixed Mode Questionnaire Design", Paper presented at the Royal Statistical Society, October 11, London, UK.

Cohen, J. (1960), "A Coefficient for Agreement for Nominal Scales", *Education and Psychological Measurement*, 20, 37 – 46.

Couper, M.P., Kreuter, F. (2013), "Using Paradata to Explore Item-Level Response Times in Surveys" *Journal of the Royal Statistical Society Series A*, 176, 271 – 286.

Eckman, S., Kreuter, F., Kirchner, A., Jäckle, A., Tourangeau, R., Presser, S. (2014), "Assessing the Mechanisms of Misreporting to Filter Questions in Surveys", *Public Opinion Quarterly*, 78 (3), 721 – 733.

Fleiss, J.L. (1971), "Measuring Nominal Scale Agreement among Many Raters", *Psychological Bulletin*, 76 (5), 378 – 382.

Fowler, F. J., Jr. (1995), *Improving Survey Questions: Design and Evaluation*. Applied Social Research Methods Series, 38. Thousand Oaks, CA: Sage Publications.
Gallhofer, I. N., Scherpenzeel, A., Saris, W.E. (2007), "The Code-Book for the SQP Program", available at <http://sqp.upf.edu>.

Kreuter, F., Presser, S., Tourangeau, R. (2008), "Social Desirability Bias in CATI, IVR, and Web Surveys: The Effects of Mode and Question Sensitivity", *Public Opinion Quarterly*, 72(5), 847 – 865.

Olson, K., Smyth, J.D. (2015), "The Effect of CATI Questions, Respondents and Interviewers on Response Time", *Journal of Survey Statistics and Methodology*, 3, 361 – 396.

Saris, W. E., Gallhofer, I. (2007), "Estimation of the Effects of Measurement Characteristics on the Quality of Survey Questions", *Survey Research Methods*, 1(1), 29 – 43.

Schaeffer, N.C., Dykema, J. (2010), *Characteristics of survey questions. A review*, American Association for Public Opinion Research, May, Chicago, USA.

Schaeffer, N.C., Dykema, J. (2011), Questions for surveys. Current trends and future directions, *Public Opinion Quarterly*, 75 (5), 909 - 961.

Schouten, B., Brakel, J. van den, Buelens, B., Laan, J. van der, Klausch, L.T. (2013), "Disentangling Mode-Specific Selection and Measurement Bias in Social Surveys", *Social Science Research*, 42, 1555 – 1570.

Shoukri, M.M. (2010), Measures of Interobserver Agreement and Reliability, Chapman and Hall, CRC Biostatistics Series, 2nd edition.

Tourangeau, R., Rips, L.R., Rasinski, K. (2000), *The Psychology of Survey Response*, Cambridge University Press, UK.

Tourangeau, R. and Yan, T. (2007), "Sensitive Questions in Surveys", *Psychological Bulletin*, 133(5), 859 – 883.

Van der Vaart, W., Van der Zouwen, J., Dijkstra, W. (1995), "Retrospective Questions: Data Quality, Task Difficulty, and the Use of a Checklist", *Quality and Quantity*, 29 (3), 299 – 315.

Yan, T., Tourangeau, R. (2008), "Fast Times and Easy Questions: The Effects of Age, Experience, Question Complexity on Web Survey Response Times", *Applied Cognitive Psychology*, 22, 51 – 68.

Appendix A - Questionnaire profiles

We include the questionnaire profile estimates for the surveys not shown in section 3.2. In all figures, the item characteristics are organized from left to right and from top to bottom: difficult language question, difficult language answer, centrality, emotional loading, sensitive information and presumed filter question. 95% confidence intervals (light grey) are based on a normal approximation.

Figure A.1: Questionnaire profile for core study Assets (AS).

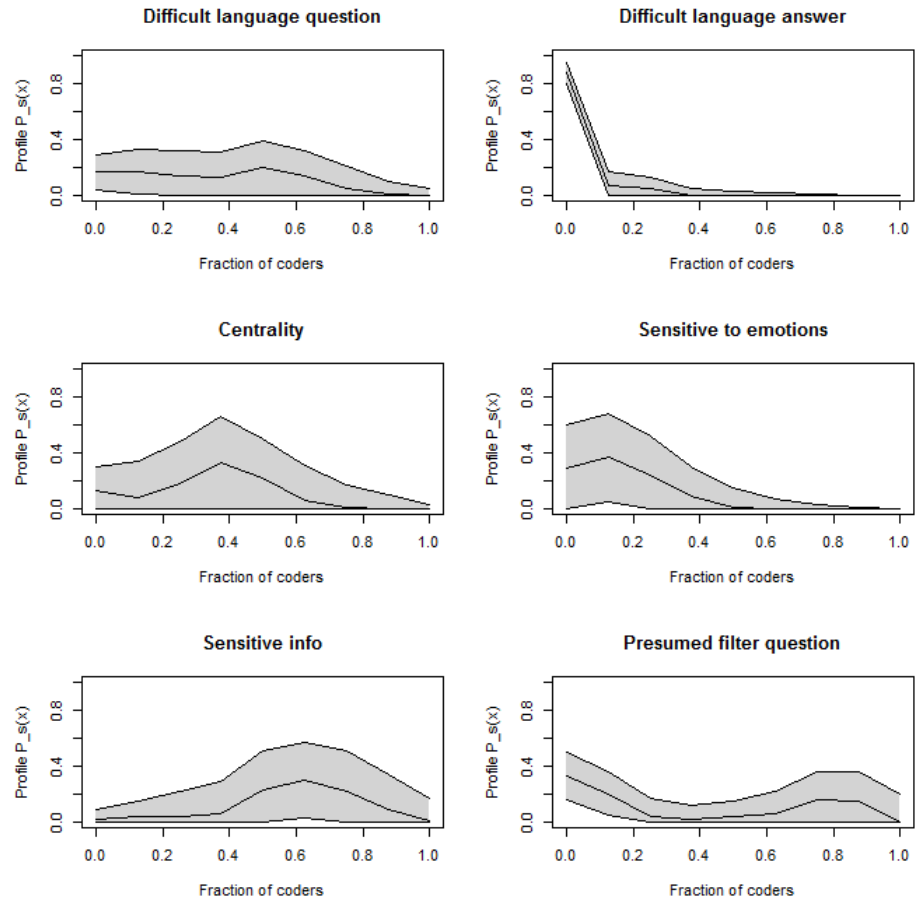


Figure A.2: Questionnaire profile for core study Family and household (FA).

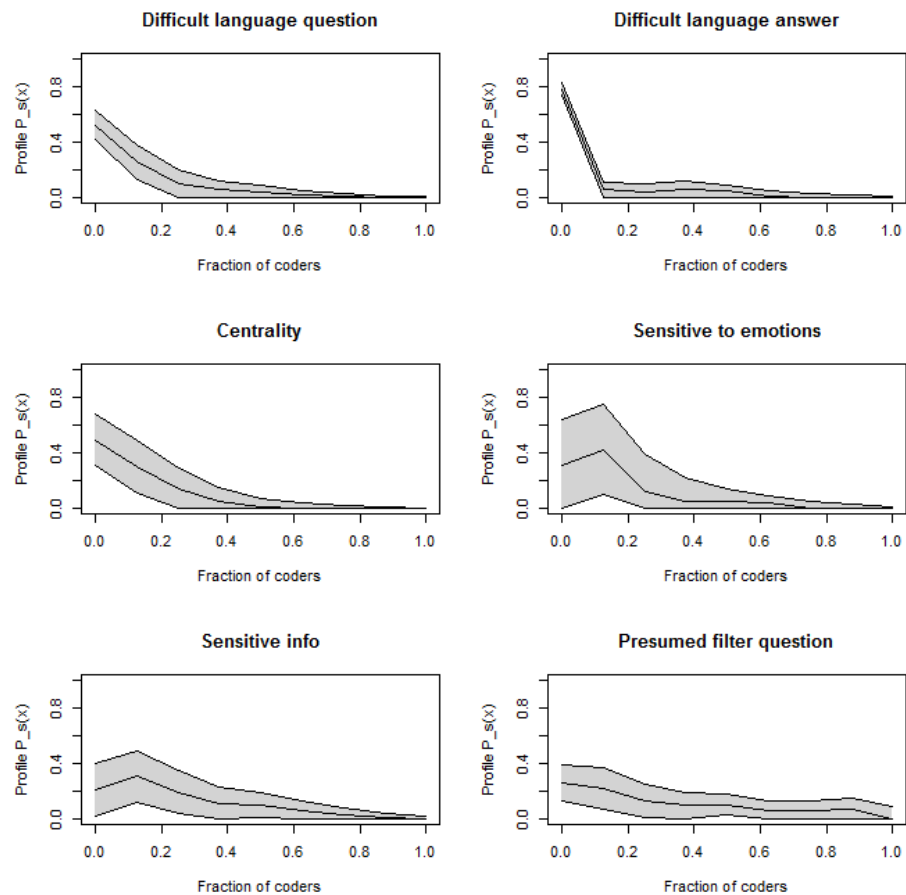


Figure A.3: Questionnaire profile for core study Health (HE).

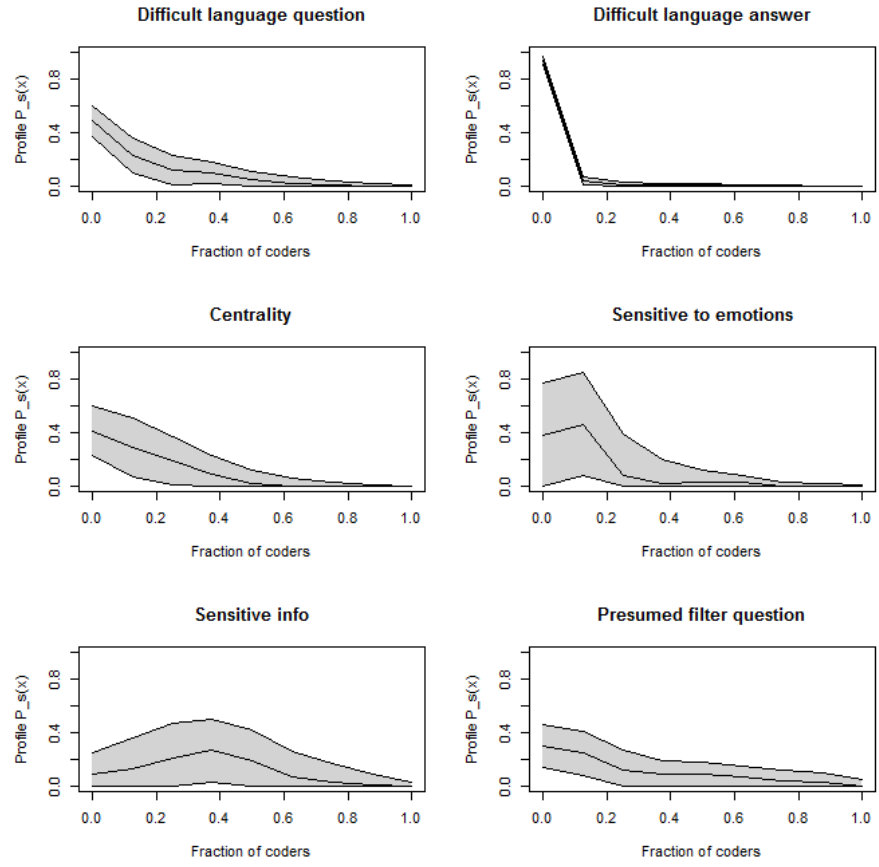


Figure A.4: Questionnaire profile for core study Housing (HO).

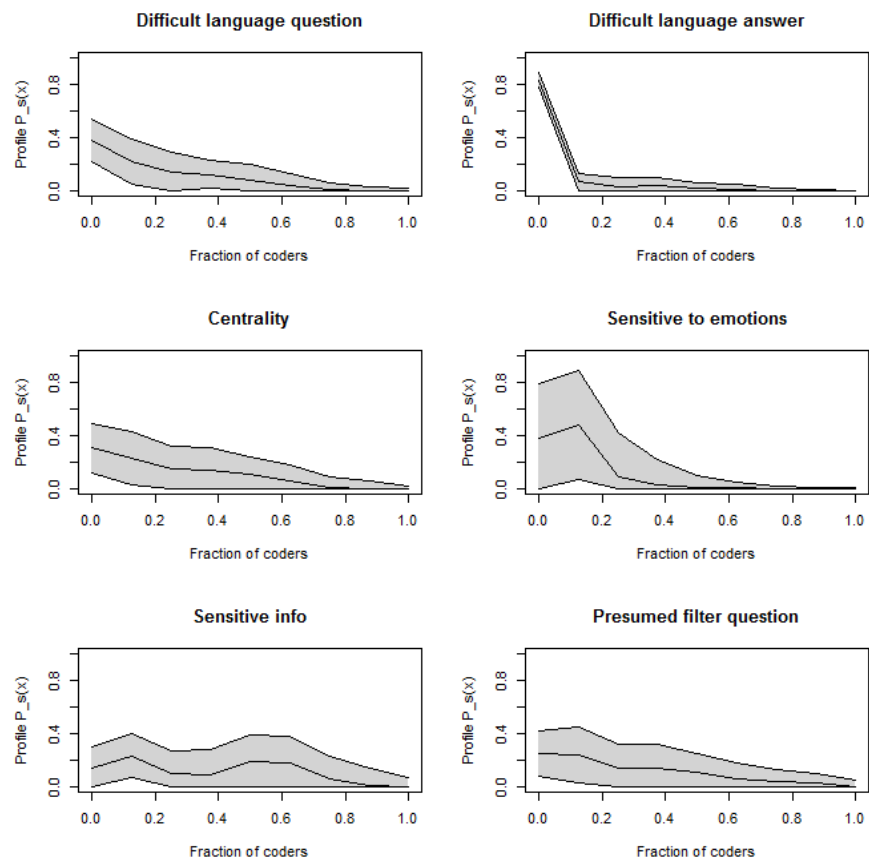


Figure A.5: Questionnaire profile for core study Personality (PE).

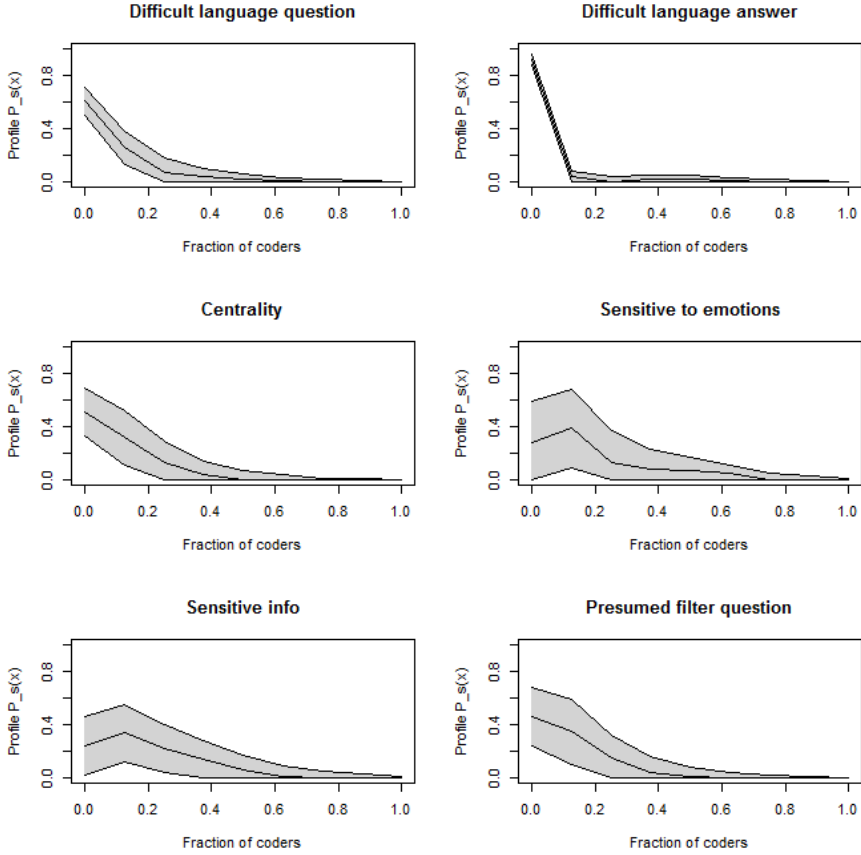


Figure A.6: Questionnaire profile for core study Religion and ethnicity (RE).

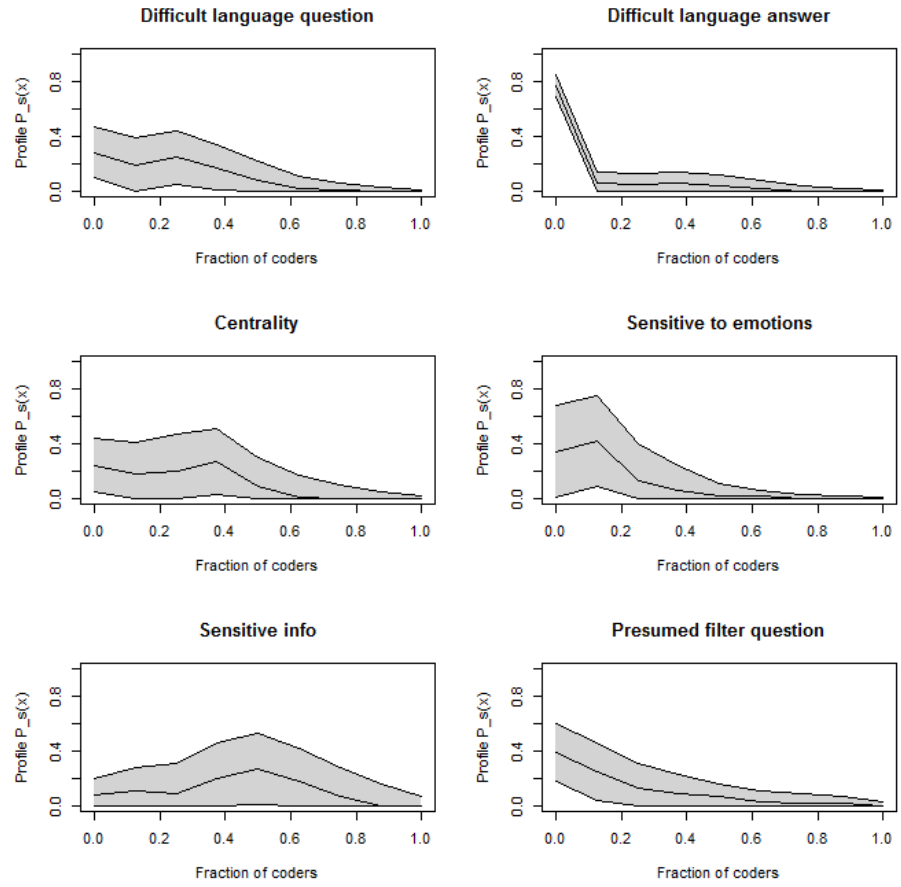


Figure A.7: Questionnaire profile for core study Social integration and leisure (SO).

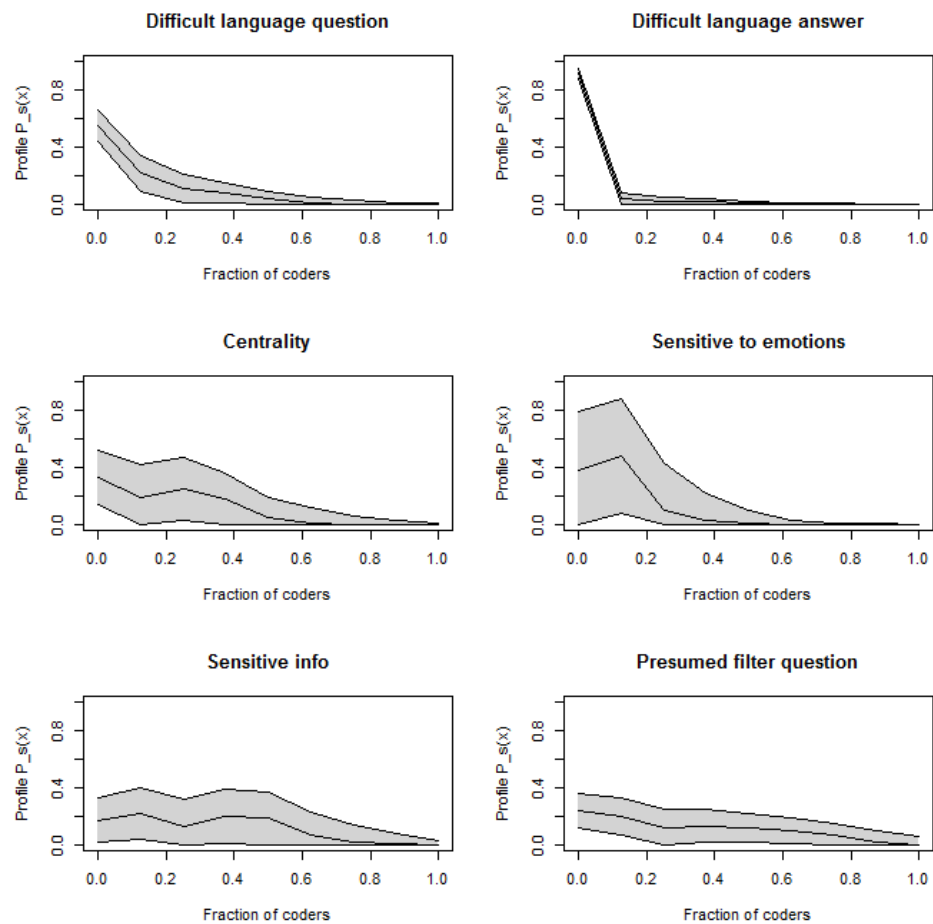


Figure A.8: Questionnaire profile for core study Work and schooling (WO).

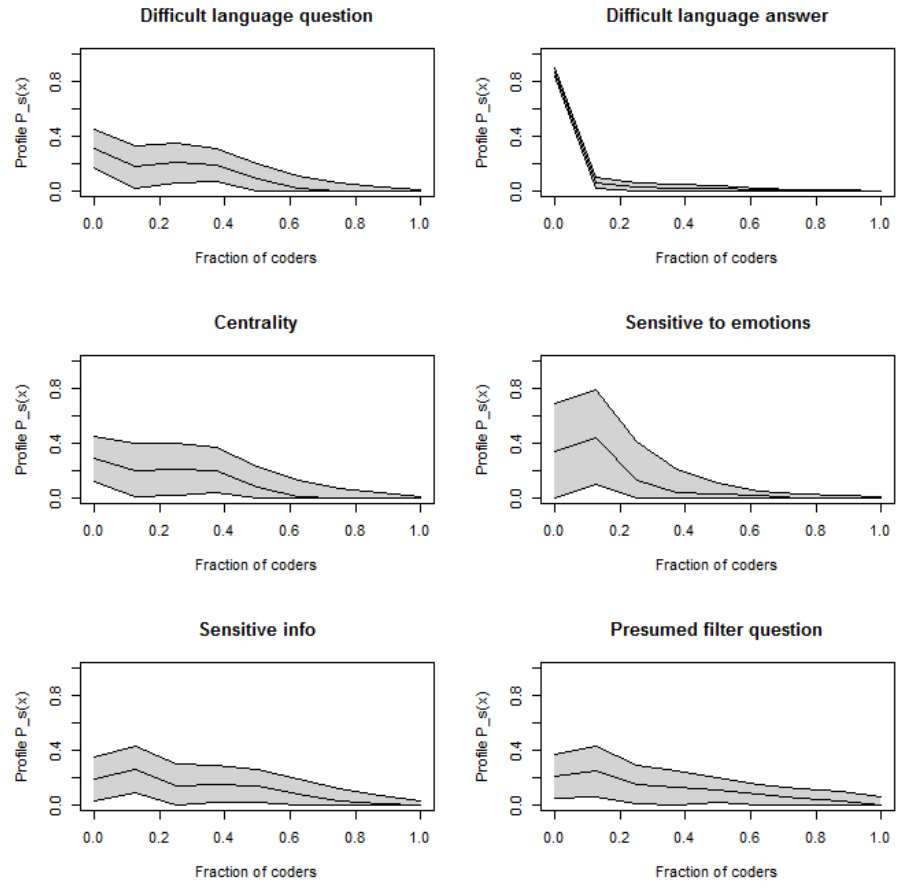
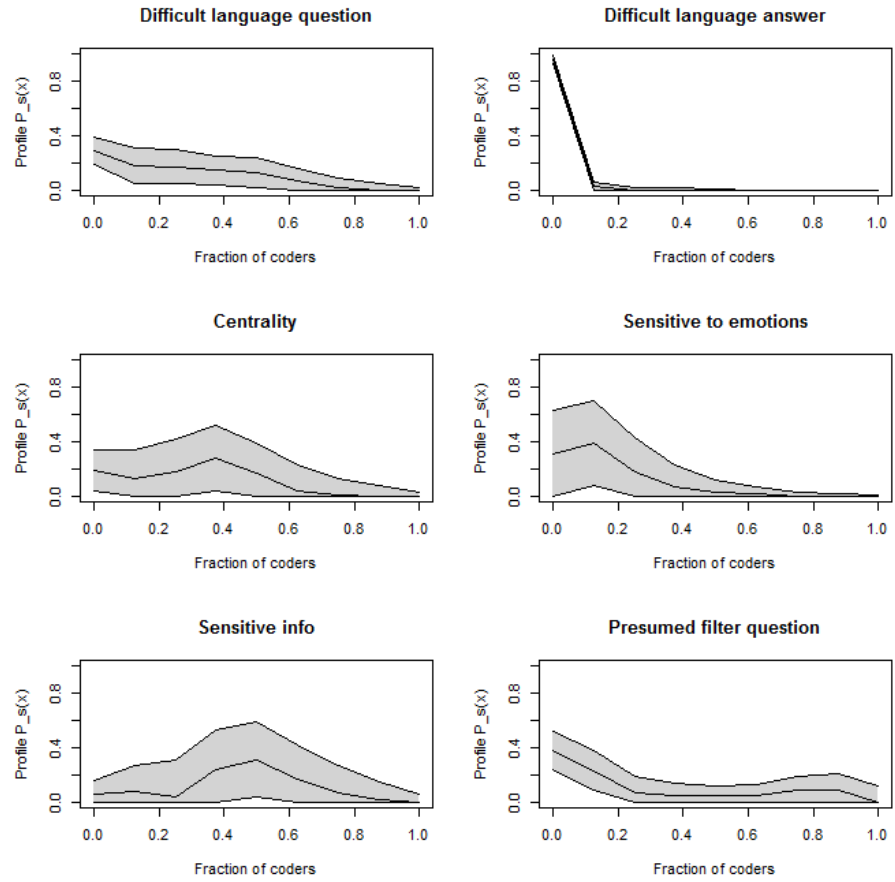


Figure A.9: Questionnaire profile for core study Income (IN).



Explanation of symbols

Empty cell	Figure not applicable
.	Figure is unknown, insufficiently reliable or confidential
*	Provisional figure
**	Revised provisional figure
2014–2015	2014 to 2015 inclusive
2014/2015	Average for 2014 to 2015 inclusive
2014/'15	Crop year, financial year, school year, etc., beginning in 2014 and ending in 2015
2012/'13–2014/'15	Crop year, financial year, etc., 2012/'13 to 2014/'15 inclusive

Due to rounding, some totals may not correspond to the sum of the separate figures.

Colofon

Publisher

Statistics Netherlands
Henri Faasdreef 312, 2492 JP The Hague
www.cbs.nl

Prepress

Statistics Netherlands, Studio BCO

Design

Edenspiekermann

Information

Telephone +31 88 570 70 70, fax +31 70 337 59 94
Via contactform: www.cbsl.nl/information

© Statistics Netherlands, The Hague/Heerlen/Bonaire 2015.

Reproduction is permitted, provided Statistics Netherlands is quoted as the source.