

# Missing Data in the Health Interview Survey 1997 in Belgium

by

Burzykowski T., Molenberghs G., Tafforeau J.,  
Van Oyen H., Demarest S., Bellamammer L.

---

## Abstract

*Health surveys, as other types of population surveys, inevitably face the problem of incomplete data, which may influence the final results and therefore should receive careful attention, starting from the design phase. Also, an investigation of the influence on the results should be undertaken. The Belgian Health Interview Survey (HIS) was conducted in 1997. In this paper the methods used in design and conduct of the HIS to reduce incompleteness are described and an analysis of the influence of the missing data on results of the survey is presented. Some aspects of such a survey are difficult (e.g. household participation is difficult to study), and therefore we carefully discuss the advantages and possible limitations of such efforts.*

## Key-Words

Beta-Binomial Model; Missing Data; Multiple Imputation; Health Survey.

---

Biostatistics, Limburgs Universitair Centrum, Universitaire Campus, B-3590 Diepenbeek, Belgium. e-mail: tomasz.burzykowski@luc.ac.be and geert.molenberghs@luc.ac.be  
Scientific Institute of Public Health, Rue J. Wytsman 14, B-1050 Brussels, Belgium.  
e-mail: Jean.Tafforeau@ihe.be, Herman.VanOyen@ihe.be, and Demarest@epinov.ihe.be  
National Institute of Statistics, Rue de Louvain 44, B-1000 Brussels, Belgium.

## **1. Introduction**

Population surveys, such as health surveys, are always confronted with the problem of incomplete data. Since the occurrence of missing data may influence the final results, they should receive careful attention during the design of a health survey. Also, an investigation of the influence on the results should be undertaken. In health surveys organized by different countries in the past, missing data received generally rather limited attention and treatment, although there are important exceptions (e.g., the Third National Health and Nutrition Examination Survey, NHANES III) (1).

The Belgian Health Interview Survey (HIS) was conducted in 1997. It was the first such survey in Belgium. The aim of this paper is to describe the methods used in design and analysis of the HIS to reduce the risk of occurrence of missing data and to study their influence on results of the survey.

Section 2 is devoted to the general design of the HIS, while Section 3 focuses on the design measures aimed at reducing incompleteness and enabling the study thereof. A taxonomy of missing data types, useful to facilitate further study, is reviewed in Section 4. A quantitative description of the extent of missingness at the household level is given in Section 5, while its impact on the scientific conclusions is aimed in Section 6. Finally, Section 7 is devoted to item-level missingness.

## **2. Design of the Health Interview Survey**

In the HIS a total sample of 10.000 interviews (0.1% of the Belgian population) was planned, to be spread over the year 1997, evenly among the four quarters. Here and in the remaining part of the paper the term "quarters" will be used in reference to the calendar quarters of the year 1997. For the three regions of Belgium (Flemish region, Walloon region and Brussels region) the number of individuals to be successfully interviewed was preset at 3500, 3500 and 3000, respectively. An over-sampling was planned for the German Community of Belgium (in the district Eupen-Malmédy), with 300 successful interviews. A detailed description of the sampling scheme used in the HIS was published elsewhere (2, 3). We now summarize its most important features. Sampling was based on a combination of stratification, multistage sampling, and clustering (4).

1. There were two stratification levels:

- (a) First, stratification was done at the regional level, to ensure that the preset precision at the regional level could be reached.
- (b) Second, stratification was conducted at the level of provinces, proportional to their size.

2. As a result, the population of Belgium was divided into 12 strata defined by five provinces in the Flemish region, five provinces in the Walloon region (excluding the German Community from the province of Liège), the Brussels region and the German Community. Selection of individuals within each of the 12 strata was done in three stages:

- (a) Within each stratum one first selected municipalities, with probability proportional to their size.
- (b) Within each municipality, households (HH) were selected, with equal probability. As drop-out of HHs was expected, HHs were sampled in clusters of four, so that for each HH there were three replacement HHs (see also Section 3).
- (c) Finally, within each HH, individuals were selected based on a few rules. In each HH at most four individuals were interviewed. This always included the reference person and partner (if applicable). (The reference person is the administrative reference of the HH within the National Register; although it may be in general an adult, it is not necessarily so.) The other respondents were selected using a birthday rule.

Figure 1 schematically summarizes the sampling scheme of the HIS.

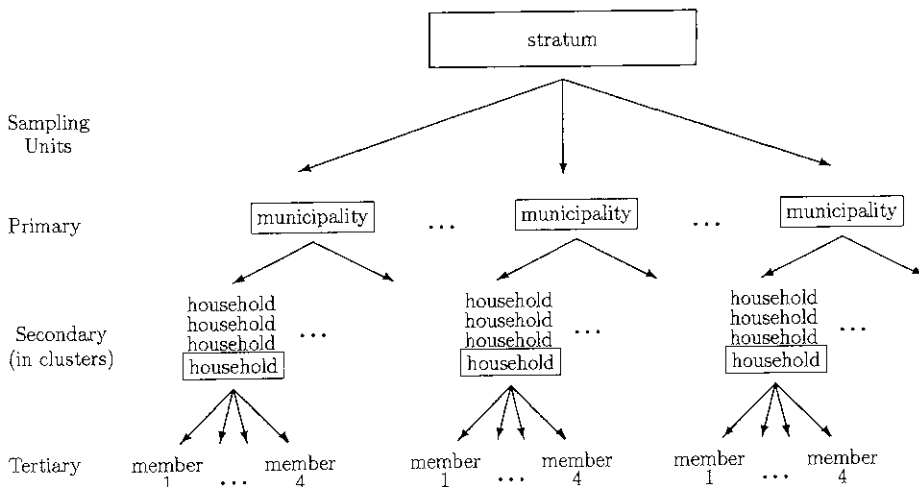


Fig. 1: Sampling scheme of the Belgian Health Interview Survey

### **3. Design Measures Towards Missing Data**

It was expected that due to various reasons (e.g., impossible to locate, refusal) not all sampled HHs would result in an interview. To compensate for unsuccessful attempts, it was decided to increase the number of sampled HHs until the target was reached. Additionally, to tackle, at least partially, systematic trends in drop-out of HHs, it was decided not to replace non-interviews in a simple random fashion, but to seek matches based on:

1. the statistical sector within the primary sampling unit (the Belgian municipalities are divided into statistical sectors, which are thought to be homogenous areas with respect to demographic, economic and social characteristics, as well as related to housing and transport);
2. the size of the HH;
3. the age of the reference person.

To achieve the aims mentioned above, clustered systematic sampling of HHs was performed within each primary sampling unit, using the National Register as the sampling frame. Clusters of four matching HHs were selected. To anticipate fluctuations due to variable number of HH members and drop-out of HHs, the number of sampled clusters was taken twice its expectation. As a consequence, the number of sampled HHs was eight times higher than the expected number of HHs.

The order of contacting the selected clusters of HHs was randomized. In case of failure to obtain an interview for a HH in a cluster, the next HH in the cluster was contacted. The order of contacting the HHs within a cluster was randomized. If none of the HHs in a cluster resulted in an interview, the next available cluster within the primary sampling unit was initiated. Clearly, in this case, no matching was done.

Individuals participating in the survey had to fill out a self-administered questionnaire and had to answer a set of questions in a face-to-face interview. To reduce effects of possible non-response of individuals within participating HHs, the face-to-face interview allowed proxies, following a set of rules:

- Mandatory:
  - for persons younger than 15 years;
  - for persons older than 60 years with a negative score on a set of introductory questions (to assess cognitive ability);
  - for persons too sick or with mental disabilities.
  
- For persons that could not be reached for an extended period (at least 1 month).
  
- For persons refusing an interview but not refusing proxy use.

Additional measures were taken to simplify the fieldwork in general and to reduce missingness in particular. For example, a letter describing the aims of the survey and inviting a HH to participate was sent two weeks prior to any attempt to contact the HH by an interviewer. Also, multiple attempts to contact HHs were to be undertaken, with minimum numbers of phone and personal contact attempts set to 5 and 3, respectively. Whenever it was possible, reasons for failure to obtain an interview were recorded on special forms.

#### 4. Types of Missingness

In order to better describe missingness and to study the impact thereof, we will now give the taxonomy to be used throughout the remainder of the paper:

**Household-Level Non-Participation (NP):** No interview, irrespective of reasons, was obtained for a HH.

**Household-Level Non-Availability (NA):** No interview was obtained due to problems with contacting the HH (e.g., difficulty in contacting the reference person; administrative decision to stop the attempts to contact the HH because of change of the interviewer or because the sample size for that particular quarter was reached).

**Household-Level Non-Response (NR):** The HH was contacted but it explicitly refused to participate in the survey.

“Missingness” will be used as a generic name to indicate all or any one of these three categories.

To obtain the required sample size of 10.000 individuals, 35.023 HHs were sampled based on information contained in the National Registry. 11.568 of these HHs were attempted to be contacted in order to obtain an interview. In 4664 (40.3%) cases the HH was successfully interviewed. In 6904 (59.7%) cases the attempt failed (NP): in 3358 (29.0%) cases the HH refused (NR), while in 3546 (30.7%) cases difficulties with contacting the HH (NA) occurred.

These figures indicate that HH-level missingness is an important issue: it needs to be taken into account carefully in order to reach the required sample size and it may influence the conclusions. For future surveys it might be of interest to study which HHs have a higher risk to be missing and to adjust the design accordingly. Further, it may be argued that the results are based on a self-selected sample of HHs, namely those that agreed to participate in the HIS. These HHs may differ systematically from those that did not agree to take part in the survey and the difference may bias the results. Therefore, HH-level missingness deserves a more detailed investigation.

At the individual level, we may distinguish between NP, NA, and NR in a similar fashion. In total, 10.339 members of the HHs that participated in the HIS were selected for the interview. 210 of them refused to the face-to-face part (individual-level NR). For 96 of the refusing individuals a proxy was used; for the remaining 114 there was no proxy and these interviews are missing. The results of the face-to-face part are also missing (due to unknown reasons) for four additional persons. Hence, for 118 individuals no data for the face-to-face interview were obtained. The remaining 10.221 individuals were included into the final analysis of the HIS. Among those, 7889 responded personally while 2332 used a proxy. The main reason for use of the proxy was age less than 15 years (1661 cases). In 408 cases the use of a proxy was due to difficulties with contacting the person which was to be interviewed (individual-level NA). In total, individual-level NP was noted for 785 persons.

The above data indicate that, if a HH agreed to participate in the survey, interviews for the HH members were obtained in the majority of cases. It can therefore be concluded that characteristics of individual-level NP, NA and NR are of secondary importance as compared to their HH-level counterparts. Hence, they will not be studied further in this article.

Finally, individuals who took part in the survey occasionally refused to answer one or more questions, thereby generating missing values

for a subset of the questions. This is referred to as item-level NR. Frequencies of item-level NR depended on the question under consideration, but never exceeded 11%.

In the main analysis of the HIS, an available-case strategy was used (5-7), meaning that all cases for which the subset of variables under study were recorded contributed to the analysis (8). It is of interest to investigate more formally the effect of the choice of this strategy on the results of the survey.

## **5. Analysis of Missingness at the Household Level**

In this section methods and results of the statistical analysis of household-level non-participation (NP), non-availability (NA) and non-response (NR) are presented and discussed in turn.

### *5.1 Household-Level Non-Participation*

#### **5.1.1 Methods**

To investigate the data on HH-level NP in more detail, a modelling approach was considered. Ideally, such an analysis should take into account the possible clustering effect, which can be thought of as arising from matching HHs. Indeed, it is a priori conceivable that HHs within a cluster would tend to exhibit correlation. Also, it is worth noting that for HHs which were not attempted to contact in a cluster (because an interview for this cluster had already been obtained or because the overall sample size for that particular quarter had been reached) there is no information whether they would participate in the survey or not. Thus formally, we have the somewhat extraordinary feature that, from a NP view-point, those HHs should be considered missing. Ideally, they should also be incorporated into the analysis.

A well-known model for analysis of clustered binary data is the beta-binomial model (9, 10). An extension of the model, as described in the Appendix, was developed to take HHs with missing NP-status into account. The probability of HH-level NP was allowed to depend on cluster-level covariates (region, HH size based on the mean theoretical size of HHs in the cluster, age group based on the mean age of reference persons of HHs included in the cluster) that corresponded to the variables used for matching HHs. A common value for the intra-cluster correlation was assumed.

In the model described above the value of the intra-cluster correlation coefficient was estimated to equal 0.07, with 95% c.i. equal to (0.04, 0.09), suggesting that there was little dependence among HHs with respect to their NP. Hence, a simpler logistic regression model describing the probability of HH-level NP can be considered, without altering the results in a meaningful way. The HH-level covariates used are those defining the rows in Table 1. HHs with missing NP-status were excluded from the analysis. The deviance for the model was equal to 36.54, which was almost exactly the number of its degrees of freedom, 36. Thus, there was no evidence for overdispersion, which agreed with the conclusion drawn from the beta-binomial model.

As mentioned before, the logistic model excluded the HHs that were not attempted to contact in clusters in which an interview was obtained for one of the previously contacted HHs, or in which the process of contacting HHs was stopped due to an administrative decision (e.g., associated with interviewer replacement). For such HHs no information is available whether they would have participated in the survey or not. However, since for these HHs missingness occurs due to external reasons which are unrelated to HH's characteristics, the missing data can be ignored safely when fitting the model (11).

### **5.1.2 Results**

For 6904 out of 11.568 HHs that were attempted to contact no interview was obtained, due to different reasons. Table 1 presents the observed percentages of HH-level NP by quarter (when the cluster of HHs was initiated), region, HH size, as well as age, sex and nationality of a HH reference person. The values of covariates (except for quarter) are based on those recorded in the National Register.

Table 1 suggests that the probability of HH-level NP was higher for the Brussels region than for Flanders and Wallonia, for female reference persons than for male, and for non-Belgian reference persons than for Belgians. The NP probability decreased with HH size. Neither quarter nor reference person's age effects seem to be present.

The results of the best fitting logistic regression model, with HH-level NP as the response, are presented in Table 2. They confirm the remarks drawn from inspecting Table 1. According to the model, the HH-level NP probability depended on region, HH size, sex and nationality of a HH reference person. In addition, for non-Belgian reference persons, the



TABLE 1  
Observed percentages (NP % out of N) of household-level non-participation

Covariate	Quarter								Overall	
	1st		2nd		3rd		4th			
	N	NP %	N	NP %	N	NP %	N	NP %	N	NP %
BELGIUM	2647	57.5	2923	59.8	2844	60.8	3154	60.3	11568	59.7
AGE										
0-14	20	75.0	23	95.7	2	50.0	3	66.7	48	83.3
15-24	185	68.1	252	70.6	133	66.9	123	76.4	693	70.3
25-34	482	53.3	527	57.3	481	58.6	548	60.0	2038	57.4
35-44	477	50.3	509	56.4	474	57.0	551	55.5	2009	54.9
45-54	426	59.2	496	58.7	467	59.3	505	59.4	1894	59.1
55-64	370	56.5	375	58.4	429	57.6	456	59.7	1630	58.1
65-74	407	59.7	418	61.0	448	60.0	531	59.0	1804	59.9
75+	280	64.4	323	71.8	412	71.8	437	65.5	1452	66.0
GENDER										
Female	877	61.7	936	66.6	934	66.0	1043	65.3	3790	64.9
Male	1770	55.5	1987	56.7	1910	58.3	2111	57.8	7778	57.1
NATIONALITY										
Belgian	2189	55.8	2426	58.0	2378	58.2	2619	59.3	9612	57.9
Non-Belgian	458	65.9	497	68.8	466	74.3	535	65.2	1956	68.5
REGION										
Brussels	1024	65.3	1212	68.1	1081	67.8	1368	66.8	4685	67.0
Flanders	776	51.0	757	48.8	858	55.6	819	56.2	3210	53.0
Wallonia	847	54.1	954	58.2	905	57.5	967	54.6	3673	56.1
HH SIZE										
1	1151	63.9	1250	67.4	1249	68.3	1320	67.6	4970	66.9
2	704	56.8	827	58.2	789	57.5	896	58.4	3216	57.8
3	363	55.9	416	54.1	391	56.5	424	54.0	1594	55.1
4+	429	42.9	430	46.5	415	48.7	514	50.2	1788	52.8

odds of HH-level NP were approximately 1.7 higher than the odds for a Belgian reference person. Moreover, the nationality effect was more pronounced in Wallonia and Flanders. In Wallonia the odds of HH-level NP for non-Belgian reference persons were additionally increased 1.3 times, and were 2.2 times higher than the odds of Belgian reference persons; for Flanders the corresponding numbers are 2.0 and 3.4, respectively. Another form of summarizing these differences is shown in Table 3, which contains expected probabilities of HH-NP for different combinations of covariates, calculated for the model presented in Table 2.

TABLE 2

Odds ratios, with 95% confidence intervals and p-values, for a multiple logistic regression model for the household-level non-participation probability.  
(NOTE: Household (HH) size included into the model as a continuous covariate with unit-spaced levels.)

Covariate	OR (95% C.I.)	p-value
REGION		
Brussels	1	–
Flanders	0.76 (0.62, 0.94)	0.01
Wallonia	0.59 (0.49, 0.72)	< 0.001
GENDER		
Female	1	–
Male	0.90 (0.83, 0.99)	0.02
NATIONALITY		
Belgian	1	–
Non-Belgian	1.68 (1.32, 2.12)	< 0.001
HH SIZE	0.83 (0.77, 0.89)	< 0.001
REGION*NATIONALITY		
Brussels*Non-Belgian	1	–
Flanders*Non-Belgian	2.00 (1.39, 2.88)	< 0.001
Wallonia*Non-Belgian	1.29 (1.01, 1.66)	0.04
REGION*HH SIZE		
Brussels*HH size	1	–
Flanders*HH size	0.90 (0.82, 0.99)	0.03
Wallonia*HH size	1.06 (0.98, 1.16)	0.16
NATIONALITY*HH SIZE		
Belgian*HH size	1	–
Non-Belgian*HH size	0.87 (0.79, 0.95)	0.003

## 5.2 Household-Level Non-Availability

### 5.2.1 Methods

In case of HH-level NA, a stronger clustering effect than the one for HH-NP was expected. Indeed, if a particular region of a municipality was difficult to reach, this problem might influence contacts with all HHs from a cluster located in the area. Also, we need to realize that from the technical point of view of HH-level NA, HHs which were not attempted to contact in a cluster are set equal to missing.

Given this potential clustering, we resort again to the extended beta-binomial model. The probability of HH-level NA was allowed to depend on cluster-level covariates that corresponded to the variables used for matching HHs. Specifically, these are: region, HH size based on the

TABLE 3  
*Predicted probabilities of household-level non-participation based  
 on the model presented in Table 2.*

Region	HH size	Probability (%)			
		Female		Male	
		Belgian	Non-Belgian	Belgian	Non-Belgian
Brussels	1	70.5	77.7	69.3	75.8
	2	66.4	71.5	64.1	69.3
	3	62.1	64.3	59.6	61.9
	4+	57.6	57.5	55.1	53.9
Flanders	1	62.2	82.7	59.7	81.8
	2	55.1	75.7	52.5	73.7
	3	47.8	66.9	45.2	64.5
	4+	40.6	56.7	28.2	54.1
Wallonia	1	60.1	73.9	57.6	71.9
	2	57.0	68.4	54.5	66.2
	3	53.9	62.4	51.3	59.9
	4+	50.7	56.0	48.2	53.4

mean theoretical size of HHs in the cluster, age group based on the average age of reference persons in the cluster. The model also includes a cluster-correlation parameter.

This model properly accounts for incompleteness, provided that missingness does not depend on the unobserved data. If deemed necessary, more complicated models can be considered as well (12,13).

### 5.2.2 Results

Among 6904 HHs for which attempts to obtain an interview failed, in 3546 cases the reason was difficulty with contacting the reference person. Table 4 presents the observed percentages of HH-level NA by quarter of the year 1997, region, HH size, as well as age, sex and nationality of a HH reference person. It indicates that the proportion of HHs that were difficult to contact generally decreased with HH size and age of the reference person. It was higher for non-Belgian reference persons and for females, and differed by region.

Table 5 presents results obtained for the beta-binomial model. The estimated value of the intra-cluster correlation coefficient equals  $\rho = 0.12$ , with the 95% c.i. equal to (0.10, 0.15). As expected, the clustering effect is somewhat higher than that observed in the HH-NP analysis (Section 5.1.1). It is still relatively small, though. The results shown in Table 5 con-

TABLE 4  
Observed percentages (NA % out of N) of household-level non-availability

Covariate	Quarter								Overall	
	1st		2nd		3rd		4th		N	NA %
	N	NA %	N	NA %	N	NA %	N	NA %		
BELGIUM	2647	26.5	2923	30.3	2844	33.1	3154	32.23	11568	30.7
AGE										
0-14	20	65.0	23	65.2	2	0.0	3	66.7	48	62.5
15-24	185	51.4	252	49.2	133	53.4	123	56.1	693	51.3
25-34	482	34.0	527	38.3	481	38.1	548	37.8	2038	37.1
35-44	477	23.3	481	27.5	474	35.4	551	32.1	2009	29.6
45-54	426	19.5	496	25.8	467	35.3	505	30.9	1894	28.1
55-64	370	22.2	375	24.5	429	26.8	456	30.3	1630	26.2
65-74	407	17.9	418	22.7	448	25.2	531	26.6	1804	23.4
75+	280	28.9	323	27.6	412	31.1	437	29.1	1452	29.3
GENDER										
Female	877	27.6	936	22.7	934	36.0	1043	36.8	3790	34.5
Male	1770	26.0	1987	27.2	1910	31.7	2111	30.0	7778	28.8
NATIONALITY										
Belgian	2189	23.4	2426	28.1	2378	29.6	2619	29.2	9612	27.7
Non-Belgian	458	41.3	497	41.1	466	51.3	535	47.3	1956	45.3
REGION										
Brussels	1024	18.4	1212	17.3	1081	23.9	1368	24.2	4685	21.1
Flanders	776	36.6	757	41.9	858	42.7	819	40.1	3210	40.4
Wallonia	847	21.7	954	25.8	905	30.4	967	28.0	3673	26.6
HH SIZE										
1	1151	36.4	1250	41.8	1249	42.9	1320	43.2	4970	41.2
2	704	21.9	827	25.4	789	25.9	896	25.6	3216	24.5
3	363	20.1	416	20.2	391	23.0	424	24.3	1594	22.0
4+	429	13.1	430	16.1	415	27.0	514	22.4	1788	19.7

firm that the probability of HH-level NA decreased with HH-size and age of the reference person, and it depended on the region. They also indicate that the probability was higher for quarters 2-4 than for the first quarter.

### 5.3 Household-Level Non-Response

#### 5.3.1 Methods

Among 6904 HHs for which no interview was obtained, 3358 explicitly refused to take part in the survey. The remaining 3546 HHs could not be contacted and therefore it is not known whether they would have agreed to participate in the survey. Thus formally, they are missing data from the point of view of an HH-level NR analysis. Similarly, the HHs that were not

TABLE 5

Odds ratios, with 95% confidence intervals and p-values, for the extended beta-binomial model for the probability of household-level non-availability. Estimated intracluster-correlation coefficient:  $\rho = 0.12$ ; 95% c.i. (0.10, 0.15). (NOTE: Household (HH) size and age-group included into the model as continuous covariates with unit-spaced levels.)

Covariate	OR (95% C.I.)	p-value
REGION		
Brussels	1	–
Flanders	0.44 (0.40, 0.48)	< 0.001
Wallonia	0.60 (0.54, 0.65)	< 0.001
HH SIZE	0.67 (0.65, 0.70)	< 0.001
AGE-GROUP	0.84 (0.82, 0.86)	< 0.001
QUARTER		
1st	1	–
2nd	1.14 (1.02, 1.28)	< 0.001
3rd	1.47 (1.30, 1.65)	< 0.001
4th	1.40 (1.25, 1.57)	< 0.001

attempted to be contacted in clusters in which interview was obtained for one of the previously contacted HHs, or in which the process of contacting HHs was stopped due to an administrative decision (e.g., associated with interviewers replacement), are also formally considered as missing data.

In order to analyse HH-level NR in more detail, we considered a logistic regression model. It can be shown (details omitted), that the extended beta-binomial model does not provide additional insight.

The study of this problem is complicated by the following issue. If a HH is invited to participate but is not successfully contacted, the reason for this failure may be influenced by the data *that would have been obtained* had the HH been contacted. This implies that the results in Section 5.3.2, which are based on the assumption that such an influence is not present, given covariate data, should be treated with caution and further exploration might be advisable (13).

### 5.3.2 Results

Among 6904 HHs for which no interview was obtained, 3358 explicitly refused to take part in the survey. Table 6 presents observed percentages of HHs which refused to participate in the survey.

TABLE 6  
Observed percentages (NR% out of N) of household-level non-response

Covariate	Quarter								Overall	
	1st		2nd		3rd		4th			
	N	NR %	N	NR %	N	NR %	N	NR %	N	NR %
BELGIUM	2647	31.0	2923	29.6	2844	27.7	3154	28.1	11568	29.0
AGE										
0-14	20	10.0	23	30.4	2	50.0	3	0.0	48	20.8
15-24	185	16.8	252	21.4	133	13.5	123	20.3	693	18.5
25-34	482	19.3	527	19.0	481	20.6	548	22.3	2038	20.3
35-44	477	27.0	481	28.9	474	21.6	551	23.4	2009	25.2
45-54	426	39.7	496	32.9	467	24.0	505	28.5	1894	31.1
55-64	370	34.3	375	33.9	429	30.0	456	29.4	1630	31.9
65-74	407	41.8	418	38.3	448	34.8	531	32.4	1804	36.5
75+	280	35.7	323	32.8	412	40.8	437	36.4	1452	36.7
GENDER										
Female	877	34.1	936	29.8	934	30.0	1043	28.5	3790	30.5
Male	1770	29.5	1987	29.4	1910	26.6	2111	27.9	7778	28.3
NATIONALITY										
Belgian	2189	32.3	2426	29.9	2378	28.6	2619	30.1	9612	30.2
Non-Belgian	458	24.7	497	27.8	466	23.0	535	17.9	1956	23.2
REGION										
Brussels	1024	32.6	1212	31.4	1081	31.7	1368	32.0	4685	31.9
Flanders	776	28.7	757	26.2	858	25.1	819	26.8	3210	26.6
Wallonia	847	32.4	954	32.4	905	27.1	967	26.6	3673	29.5
HH SIZE										
1	1151	27.5	1250	25.7	1249	25.4	1320	24.4	4970	25.7
2	704	34.9	827	32.8	789	31.7	896	32.8	3216	33.0
3	363	35.8	416	33.9	391	33.5	424	29.7	1594	33.1
4+	429	29.8	430	30.5	415	21.7	514	27.8	1788	27.2

The results of the logistic regression model shown in Table 7 suggest that the odds of NR of a contacted HH generally decreased with HH size. The decrease was more pronounced in Flanders, where the odds decreased additionally by approximately 50% as compared to the effect of HH size in Brussels and Wallonia. The odds for Brussels were generally higher than those for Flanders and lower than those for Wallonia. The effect of age depended on the region. It generally increased the odds in Flanders and Brussels, but slightly decreased the odds in Wallonia.

## 6. Influence of Household-Level Missingness on Study Conclusions

It is crucial to investigate to what extent HH-level missingness influences conclusions about substantive questions of the HIS. From the analysis presented in Section 5.1 it follows that the probability of NP

TABLE 7  
*Odds ratios, with 95% confidence intervals and p-values, for a multiple logistic regression model for the household-level non-response probability.*  
 (NOTE: Age-group included into the model as a continuous covariate with unit-spaced levels.)

Covariate	OR (95% C.I.)	p-value
REGION		
Brussels	1	—
Flanders	0.74 (0.41, 1.31)	0.30
Wallonia	1.17 (0.69, 1.97)	0.56
HH SIZE		
1	1	—
2	0.80 (0.60, 1.06)	0.13
3	0.92 (0.61, 1.35)	0.67
4+	0.49 (0.32, 0.74)	< 0.001
REGION*HH SIZE		
Flanders:		
1	1	—
2	0.48 (0.29, 0.78)	0.003
3	0.45 (0.23, 0.84)	0.01
4+	0.42 (0.20, 0.86)	0.02
Wallonia:		
1	1	—
2	1.05 (0.67, 1.65)	0.84
3	0.80 (0.45, 1.44)	0.46
4+	1.00 (0.54, 1.86)	0.99
AGE-GROUP	1.07 (1.00, 1.14)	0.06
AGE-GROUP*REGION		
Flanders	1.03 (0.92, 1.16)	0.58
Wallonia	0.88 (0.79, 0.97)	0.01

increased with smaller HH size. Further, it is plausible that the smaller HHs are different. For example, they may have a higher income and be more difficult to contact because their working members travel a lot. As a result, the HIS sample might include more HHs with lower income, which in turn might have influenced the health-related results of the survey.

A formal incomplete data approach to investigate such effects is to compose datasets that would have been observed had the original HHs responded. Multiple imputation is devised to achieve this goal (14-16). Briefly, multiple imputation is a technique in which each missing value is replaced by several simulated values. As a result, several sets of complete data are obtained. They may be analysed using any complete-data method. Results of these analyses are then combined to yield proper point estimates and standard errors (14-16).

However, any attempt to use multiple imputation would face the problem of lack of information regarding non-participating or non-responding HHs and their members, while such information is absolutely vital to “fill in” sensible values. Unfortunately, data available in the National Register offer only very limited information, which additionally need not be accurate. For example, it does not include information about recent changes in the size and structure of the HH, etc. It hardly provides any information that could be used to impute, for example, data on the health status of members of a HH. Clearly, such an analysis would be very imprecise and therefore of limited value.

In its commonly used form, multiple imputation is valid if the missing data are missing at random (15-16), given the observed data. While this assumption may seem strong and multiple imputation therefore questionable, it is possible to extend the technique to the missing not at random situation (16). In any case, a formal analysis of influence of HH-level missingness on the HIS results is very difficult due to the problems with scarcity of information on non-participating HHs stated above. Therefore, it was not made the primary focus of this contribution and emphasis has been placed on item-level NR, as discussed in the next section.

## **7. Influence of Item-Level Non-Response on Inference**

In this section methods and results of an investigation of influence of item-level NR on results of several selected items studied in the HIS are presented.

### *7.1 Methods*

In the HIS in total, 10.221 individuals were interviewed. As it might be expected, for certain items missing data occurred. The percentage of missing data for any of the items did never exceed 11%. However, when several items are considered together, e.g., for modelling purposes, the missing data percentages accumulate. In the available-case approach, which was used in the main analysis of the HIS, cases with incomplete data for variables under study are excluded from consideration. This results in decrease of precision, since conclusions are based on a smaller sample. Moreover, to yield valid (unbiased) results, available-case analysis requires that missing data are missing completely at random (11). That is, the missingness mechanism should depend neither on the



unobserved nor on the observed covariates. This assumption is very strong and rather unrealistic in practice. It is therefore of interest to investigate potential effects of ignoring missingness on results of the HIS analysis. In the HIS, a large amount of data (more than 1000 variables) was collected. Consequently, the scope of investigation of influence of missing data on results of the survey by means of multiple imputation had to be limited. Three variables were chosen for the analysis: body-mass index (BMI), VOEG score, and subjective opinion on health status (SOHS).

1. BMI is a genuinely continuous variable. It is an index that relates weight (in kilograms) to the square of height (in meters) of a person. Higher values of BMI (greater than 25) may be interpreted as an indicator of obesity.
2. VOEG score („Vragenlijst voor Onderzoek naar de Ervaren Gezondheid”; Questionnaire for Research on Self-Perceived Health Status) can be taken as an overall indicator of an individual’s health status. Higher values indicate increased health problems. The score is constructed based on a questionnaire containing 23 binary items concerning self-perception of the health status of an individual. Strictly speaking, VOEG is an ordered categorical variable. However, the number of categories is large (24). Therefore, VOEG score was also treated as a continuous covariate.
3. SOHS is a binary variable, based on the question: “What is your general state of health?”. Persons who evaluated their health status as “very good” or “good” (favorable response) are contrasted with those who ranked it as “fair”, “bad” or “very bad” (unfavorable response).

Several covariates were considered to explain differences in observed responses (means of BMI or VOEG or percentages of unfavorable responses for SOHS): person’s sex, age (children younger than 15 years of age were excluded), region of residence, educational level, income level, smoking status.

Since observed distributions of BMI and VOEG were skewed to the right, prior to modelling and imputation they were transformed logarithmically. They were analyzed using a weighted multiple normal regression model. SOHS was analyzed using a weighted multiple logistic regression model for probability of unfavorable response. In all cases the weights reflected the sampling probabilities for individuals as implied by the chosen sampling scheme.

Multiple imputation for BMI and VOEG was performed using the S+ software for imputing mixed continuous-categorical data written by

J. Schafer (16). Imputation was based on the general location model (16), which included  $\ln(\text{BMI})$  and  $\ln(\text{VOEG}+1)$  as continuous variables, and the covariates mentioned earlier as categorical variables. A restricted main-effects model (for the covariates) was applied. The restriction was adopted because the data were sparse, with large number (756) of cells for the log-linear part of the model. Consequently, fitting of the unrestricted model was impossible, and we restrict ourselves to main effects only.

To investigate the effects of item-level NR on inference for BMI and VOEG, coefficients of a weighted multiple normal regression model were obtained by appropriate combination of coefficients of the model fitted to 5 data sets with imputed missing data (16). The same 5 data sets were used for BMI and VOEG.

For SOHS, multiple imputation was performed using the *S+* software for imputing categorical data written by J. Schafer (16). In the imputation a saturated log-linear model for SOHS and covariates was used (16). To investigate effects of item-level NR on inference for SOHS, coefficients of a weighted logistic regression model for unfavorable response for SOHS were obtained by appropriate combination of coefficients of the model fitted to 5 data sets with imputed missing data.

## *7.2 Results*

BMI, VOEG and SOHS were to be measured for 8564 individuals older than 15 years of age that were included in the survey. They were obtained for 8384 (97.9%), 8250 (96.3%), and 7953 (92.9%) persons, respectively. When the covariates education, income, and smoking status are also taken into account (there were no missing observations for region, sex and age), the numbers of individuals available for a available-case analysis are reduced to 7389 (86.3%), 7272 (84.9%), and 7109 (83.0%), respectively.

Results of the analysis of BMI and VOEG score are presented in Table 8. The table shows coefficients of a weighted normal multiple regression model for  $\ln(\text{BMI})$  and  $\ln(\text{VOEG}+1)$  fitted to the available-case (AC) data sets of 7389 and 7272 observations, respectively, and obtained using multiple imputation (MI).

Table 8 it can be observed that, generally, coefficients estimated using the imputed data are slightly smaller than those obtained excluding missing data. There appears to be little difference in estimated standard errors of the coefficients. Thus, compared to the results of the available-case analysis, qualitative conclusions remain essentially unchanged.

TABLE 8

*Coefficients (with standard errors) for a weighted multiple normal regression models for ln(BMI) and ln(VOEG+1) obtained when the missing observations are excluded (available-cases analysis; AC) or multiply imputed (MI).*

*(NOTE: Age-group included into the model as a continuous covariate with unit-spaced levels.)*

Covariate	ln (BMI)		ln (VOEG+1)	
	AC (7272 obs.)	MI (8564 obs.)	AC (7389 obs.)	MI (8564 obs.)
REGION				
Brussels	0	0	0	0
Flanders	0.007 (0.006)	0.009 (0.006)	-0.264 (0.032)	-0.268 (0.031)
Wallonia	0.023 (0.007)	0.027 (0.006)	0.015 (0.033)	0.002 (0.033)
GENDER				
Male	0	0	0	0
Female	-0.050 (0.004)	-0.054 (0.003)	0.296 (0.019)	0.284 (0.018)
EDUCATION				
Primary	0	0	0	0
Secondary	-0.011 (0.005)	-0.013 (0.004)	-0.072 (0.023)	-0.069 (0.023)
Higher	-0.046 (0.005)	-0.045 (0.005)	-0.099 (0.025)	-0.088 (0.025)
INCOME LEVEL				
< 40,000	0	0	0	0
40,000 – 60,000	0.008 (0.004)	0.006 (0.004)	-0.049 (0.021)	-0.039 (0.021)
> 60,000	0.003 (0.006)	-0.001 (0.006)	-0.107 (0.030)	-0.094 (0.034)
SMOKING				
Non-smoker	0	0	0	0
Smoker	0.003 (0.004)	0.004 (0.004)	0.238 (0.019)	0.220 (0.019)
AGE				
Age-group	0.030 (0.001)	0.001 (0.001)	0.051 (0.006)	0.050 (0.005)

Table 9 presents estimated odds ratios for a weighted multiple logistic regression model for probability of unfavorable response for SOHS fitted to the available-case data set of 7109 observations, as well as obtained by combining results from fitting the model to 5 different data sets with missing data imputed. Similarly to the results obtained for BMI and VOEG score, it can be observed that, generally, effects estimated using the imputed data are slightly smaller than those obtained excluding missing data. Again, qualitative conclusions remain essentially unchanged.

## 8. General Conclusions

The design of the HIS implied specific issues regarding the study of the impact of missingness on the study results. Due to the multi-stage sampling, there were several reasons why attempts to interview an individual might have failed.

TABLE 9

Odds ratios (with 95% confidence intervals) for a weighted multiple logistic regression model for probability of unfavorable opinion on individual health status (SOHS) obtained when the missing observations are excluded (available-case analysis; AC) or multiply imputed (MI). (NOTE: Age-group included into the model as a continuous covariate with unit-spaced levels.)

Covariate	AC (7109 obs.)	MI (8564 obs.)
REGION		
Brussels	1	1
Flanders	0.68 (0.55, 0.83)	0.63 (0.52, 0.77)
Wallonia	1.15 (0.94, 1.40)	1.05 (0.87, 1.27)
GENDER		
Male	1	1
Female	1.52 (1.27, 1.82)	1.55 (1.30, 1.84)
EDUCATION		
Primary	1.98 (1.55, 2.53)	1.98 (1.57, 2.49)
Secondary	1.22 (0.95, 1.57)	1.25 (0.98, 1.59)
Higher	1	1
INCOME LEVEL		
< 40,000	2.57 (1.79, 3.68)	2.34 (1.66, 3.31)
40,000 – 60,000	1.78 (1.24, 2.56)	1.68 (1.17, 2.43)
> 60,000	1	1
SMOKING		
Non-smoker	1	1
Smoker	1.40 (1.16, 1.71)	1.39 (1.16, 1.66)
AGE		
Age-group	1.49 (1.41, 1.56)	1.46 (1.39, 1.55)

Analysis of HH-level missingness indicates that HHs that did not participate or refused to participate in the survey differed from those that eventually were interviewed. If similar studies are to be undertaken in the future, these differences could be used to improve organization of the studies. For example, in certain regions where missingness may be expected to be high, the number of interviewers may be increased to allow more time/flexibility for attempts to contact HHs. To this end, Tables 1-7 can be useful.

From an interpretational point of view, HH-level missingness might imply, for example, problems with generalization of HIS results to the Belgian population (unless the missingness mechanism was missing at random with respect to the National Register data). While investigation of this issue is of interest, the lack of (reliable) data on HHs that did not participate renders such a study at present virtually impossible. If it

is deemed of interest, specific action should be undertaken to collect reliable data on the missingness.

Analysis of item-level NR was performed based on a few examples. No important differences were found between available-case and multiple-imputation based results. It may be concluded that influence of item-level missing data on the HIS results is negligible, at least in the examples analysed. It is worth reminding here that multiple imputation methods assume that missing data are missing at random, i.e. that the missingness mechanism may depend on observed data (but not on the unobserved ones). This assumption is much more flexible than the stringent missing completely at random condition, needed for the available-case analysis, which requires that the missingness mechanism has to be independent of both observed and unobserved data. Consequently, our findings give additional credibility to the results of the latter analysis. More extensive investigation, however, might be needed.

Also, it should be mentioned that while missing at random mechanism is considered very plausible by many authors (17), a sensitivity analysis to assess the impact of non-random missingness may be worthwhile. This might be achieved by analysing the data using a series of models assuming different mechanisms of missingness and comparing results obtained under missing at random assumption with those obtained for the models assuming non-random mechanisms (13). In the latter models assumptions regarding the form of the dependence of the missingness mechanism on the unobserved values for the dependent variable would have to be made. Necessarily, these models are complex and in this investigation we decided to limit ourselves to use of the models described in Sections 5 and 7.

It has to be stressed that in general the influence of missing data on survey results cannot be a priori judged negligible. In each separate case it needs a careful investigation. Only after such an investigation has been conducted can it be decided whether the influence should be formally taken into account or could safely be ignored in the analysis of a survey. The results of an item-level NR analysis suggest the latter conclusion is plausible at least for the several analyzed items of the HIS.

While it is not possible to explore the entire impact of missingness (e.g., the effect of HH-level missingness is difficult to study), it is evident that a careful exploration of missingness is important (1) to assess the impact on the current study and (2) to draw lessons that can be of use in future surveys. It is our hope that the current study can contribute to these goals.

## Appendix: Extension of the Beta-Binomial Model to Missing Data

We present the extension of the beta-binomial model for the HH-level NP analysis. Suppose that for each of the four HHs in a cluster there exists a binary random variable,  $Z_i^*$  say ( $i = 1, 2, 3, 4$ ), that equals 1 if the HH will not participate in the survey, and 0 otherwise. Assume that  $Z^* = \sum_{i=1}^4 Z_i^*$  is distributed according to the binomial distribution with a random success probability  $p$ , which comes from a beta distribution with mean  $\pi$  and variance  $V$ . Then  $Z^*$  has the beta-binomial distribution:

$$P(Z^* = z^*) = \binom{4}{z^*} \frac{B(\pi(\rho^{-1} - 1) + z^*, (1 - \pi)(\rho^{-1} - 1) + 4 - z^*)}{B(\pi(\rho^{-1} - 1), (1 - \pi)(\rho^{-1} - 1))}, \quad (1)$$

where  $\rho = V(\pi(1 - \pi))^{-1}$  and  $B(a, b)$  is the beta function.  $\rho$  can be interpreted as a intra-cluster correlation coefficient.

In our case, for several HHs in a cluster we do not observe  $Z_i^*$ , due to lack of attempts to contact the HHs (e.g., because an interview has been obtained already). What we observe, instead, are the total number,  $Z$  say, of HHs participating in the survey, and the total number,  $C$  say, of HHs that were attempted to be contacted. Conditionally on  $C$ , the probability that  $Z$  HHs will participate in the survey may be calculated as:

$$P(Z = z | C = c) = P\left(\sum_{i=1}^c z_i^* = z\right) = \sum_{t=0}^{4-c} P\left(\sum_{j=1}^c z_j^* = z, \sum_{i=c+1}^4 z_i^* = t\right). \quad (2)$$

From (1) it follows that:

$$P(Z = z | C = c) = \sum_{t=0}^{4-c} \binom{c}{z} \binom{4-c}{t} \frac{B(\alpha + z + t, \beta + 4 - (z + t))}{B(\alpha, \beta)}, \quad (3)$$

where  $\alpha = \pi(\rho^{-1} - 1)$  and  $\beta = (1 - \pi)(\rho^{-1} - 1)$ .

Using the above formulation, the dependence of  $\pi$  and  $\rho$  on cluster-level covariates may be modelled using generalized linear modelling ideas. For example, the logit of  $\pi$  and Fisher's z-transform of  $\rho$  may be assumed to depend on (different) linear combinations of the covariates (18). In the models presented in the paper a common value of intra-cluster correlation for all clusters was assumed.

## References

1. SCHAFER JL, KHARE M AND EZATTI-RICE TM. Multiple imputation of missing data in NHANES III. Proceedings of the Annual Research Conference 1993: 459-487. Bureau of the Census, Washington, DC.
2. QUATAERT P, VAN OYEN H, TAFFOREAU J et al. Health Interview Survey 1997. Protocol for selection of the households and the respondents. SPH/Episerie No. 12, SPH 1998, Brussels.
3. VAN OYEN H, TAFFOREAU J, HERMANS H, QUATAERT P, SCHIETTECATTE E, LEBRUN L et al. The Belgian Health Interview Survey. Arch Public Health 1997; 55: 1-14.
4. KISH L. Survey Sampling. New York: Wiley, 1995.
5. TELLIER V, DEMAREST S, LEURQUIN P, TAFFOREAU J, VAN DER HEYDEN J, VAN OYEN H. La santé de la population en Belgique et à Bruxelles. Enquête de Santé par Interview, Belgique, 1997. Brussels: Centre de Recherche Opérationnelle en Santé Publique, Institut Scientifique de la Santé Publique, 1998.
6. VAN DER HEYDEN J, DEMAREST S, LEURQUIN P, TAFFOREAU J, TELLIER V, VAN OYEN H. De gezondheid van de bevolking in het Brussels Gewest. Samenvatting. Gezondheidsenquête, België, 1997. Brussels: Centrum voor Operationeel Onderzoek in Volksgezondheid, Wetenschappelijk Instituut Volksgezondheid-Louis Pasteur, 1998.
7. VAN OYEN H, DEMAREST S, LEURQUIN P, TAFFOREAU J, TELLIER V, VAN DER HEYDEN J. De gezondheid van de bevolking in de Vlaamse Gemeenschap. Samenvatting. Gezondheidsenquête, België, 1997. Brussels: Centrum voor Operationeel Onderzoek in Volksgezondheid, Wetenschappelijk Instituut Volksgezondheid-Louis Pasteur, 1998.
8. LITTLE RJA AND RUBIN DB. Statistical Analysis with Missing Data. New York: Wiley, 1987.
9. KLEINMANN JC. Proportions with extraneous variance: single and independent samples. J Am Stat Assoc 1973; 68: 46-54.
10. WILLIAMS DA. Extra-binomial variation in logistic linear models. Appl Statist 1982; 31: 144-148.
11. RUBIN DB. Inference and missing data. Biometrika 1976; 63: 581-592.
12. MOLENBERGHS G, KENWARD MG AND LESAFFRE E. The analysis of longitudinal ordinal data with nonrandom drop-out. Biometrika 1997; 84: 33-44.
13. MOLENBERGHS G, KENWARD MG AND GOETGHEBEUR E. Sensitivity analysis for incomplete data: region of uncertainty. Submitted for publication.
14. RUBIN DB AND SCHENKER N. Multiple imputation for interval estimation from simple random samples with ignorable nonresponse. J Am Stat Assoc 1986; 81: 366-374.
15. RUBIN DB. Multiple Imputation for Nonresponse in Surveys. New York: Wiley, 1987.
16. SCHAFER JL. Analysis of Incomplete Multivariate Data. London: Chapman and Hall, 1997.
17. RUBIN DB, STERN HS AND VEHOVAR V. Handling "don't know" survey responses: the case of the Slovenian plebiscite. J Am Stat Assoc 1995; 90: 822-828.
18. MOLENBERGHS G, DECLERCK L AND AERTS M. Misspecifying the likelihood for clustered binary data. Computational Statistics and Data Analysis 1998; 26: 327-350.