

### Figure 3.4c. SAMPLE SURVEY DESIGN/EXECUTION: Managing Non-Response

Many data-based investigations encounter the difficulty that not all the data called for in the Plan are acquired – this is the general problem of *missing data*. In dealing with *human* populations, this matter is usually referred to as *non-response* and it can impose *severe* limitations on Answers.

A framework to manage non-response, as indicated in the diagram at the right, we consider the *study* population to be made up of the *respondent* population and the *non-respondent* population. The set of units selected from the study population is the *selection*, and comprises the *sample* (from the respondent population) and the *non-respondents* (from the non-respondent population). The diagram has *two* categories of symbols:

- the  $N$ s and  $n$ s refer to *numbers of elements or units*;
- the  $\bar{Y}$ s and  $\bar{y}$ s are *averages* of a response variate  $Y$  of the elements or units.

The relationships among the numbers of elements or units are:

$$\text{Study population} = \text{Respondent population} + \text{Non-respondent population}$$

$$N_s = N + N_{nr}$$

$$\text{Selection} = \text{Sample} + \text{Non-respondents}$$

$$n_s = n + n_{nr}$$

Also, suppose that a selection of  $n_s$  units is obtained from the study population by equiprobable selecting (EPS) to estimate  $\bar{Y}_s$ , the *study* population average; we have:

$$\bar{Y}_s = \frac{\text{study population total}}{\text{study population size}} = \frac{\text{respondent population total} + \text{non-respondent population total}}{\text{respondent population size} + \text{non-respondent population size}} = \frac{N \cdot \bar{Y} + N_{nr} \cdot \bar{Y}_{nr}}{N + N_{nr}} \quad \text{-----(3.4c.1)}$$

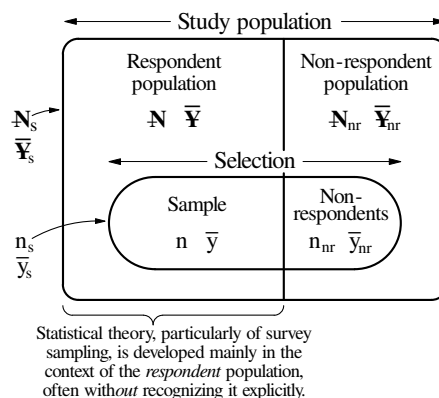
where, in the denominator:  $N + N_{nr} = N_s$ .

The direct information available to us about the response variate  $Y$  comes from the *sample* of  $n$  *respondent* population units, which yield  $\bar{Y}$  as an estimator (and  $\bar{y}$  as an estimate) of  $\bar{Y}$ ; we know from Figure 2.3 that  $\bar{Y}$  is *unbiased* for  $\bar{Y}$  – i.e.,  $E(\bar{Y}) = \bar{Y}$ . Hence, using equation (3.4c.1), the *non-responding bias*, when using  $\bar{Y}$  as an estimator of  $\bar{Y}_s$ , is given by:

$$E(\bar{Y}) - \bar{Y}_s \equiv \bar{Y} - \bar{Y}_s = \bar{Y} - \frac{N \cdot \bar{Y} + N_{nr} \cdot \bar{Y}_{nr}}{N + N_{nr}} = \frac{N_{nr}}{N + N_{nr}} (\bar{Y} - \bar{Y}_{nr}). \quad \text{-----(3.4c.2)}$$

The term:  $\frac{N_{nr}}{N + N_{nr}} = \frac{N_{nr}}{N_s}$  is the *non-response rate*; its *value* is (approximately):  $\frac{n_{nr}}{n + n_{nr}} = \frac{n_{nr}}{n_s}$ ; -----(3.4c.3)

the corresponding *response rate* is:  $\frac{N}{N + N_{nr}} = \frac{N}{N_s} = 1 - \frac{N_{nr}}{N_s}$  and its *value* is (approximately):  $\frac{n}{n + n_{nr}} = \frac{n}{n_s}$ . -----(3.4c.4)



#### 1. Dealing with Non-Response – General Considerations

Equation (3.4.2) shows that the (*unattainable*) ideal of *zero* non-responding bias, when using  $\bar{Y}$  to estimate  $\bar{Y}_s$ , occurs when:

- the non-response rate is *zero* – i.e., there are *no* non-respondents; **OR**
- $\bar{Y} = \bar{Y}_{nr}$  – i.e., the respondent and non-respondent populations have the *same* average for the response variate  $Y$ .

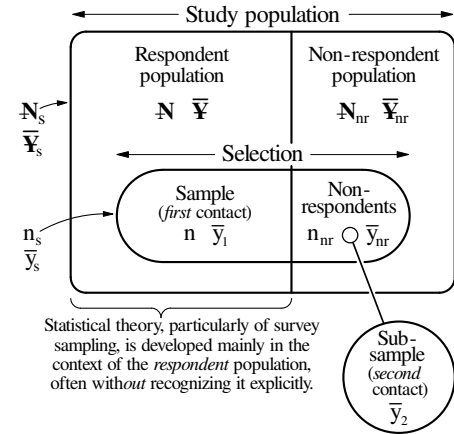
In practice, we are more concerned with the factors that *reduce* (as distinct from *eliminate*) the magnitude of the non-responding bias when using  $\bar{Y}$  to estimate  $\bar{Y}_s$ ; on the basis of our *model*, these (*two*) factors are:

- \* when  $\bar{Y}$  is *close* to  $\bar{Y}_{nr}$  – i.e., the respondent and non-respondent populations have *similar* averages for the response variate  $Y$ ;
  - this factor is *not* under the control of the investigator(s) because it depends on the distribution of values of the response variate  $Y$  over the elements (or units) which make up the respondent and the non-respondent populations;
- \* when the non-response rate is *low* – i.e., when the non-respondent population is small in size *relative* to the respondent population, or when  $N_{nr} \ll N$ ;
  - this factor *is* under the control of the investigator(s) in the sense that it is part of the Plan for a sample survey (as specified in the *sampling protocol* in the Design stage of the FDEAC cycle) – strategies for *reducing* non-response include:
    - good questionnaire design – e.g., clear, succinct, answerable questions;
    - careful interviewer selecting and training;
    - introductory letter or telephone call prior to the interview (a primary approach letter – Figure 3.2, right-hand column, first paragraph);
    - incentives for response;
    - refusal conversion techniques – e.g., randomized response (see assignment question A1 – 2);
    - follow-up reminders and call-backs;
    - use of proxy respondents – but see Figure 3.5f;
    - use of stratification.

(continued overleaf)

## 2. Dealing with Non-Response – Call-backs

One technique to reduce non-response is to *call back* (denoted below by an abbreviation *cb*) one or more times on a sub-sample of the non-respondents from the first contact, using the most experienced interviewers and/or added incentives for response. The diagram at the right (adapted from the diagram overleaf on page 3.21) illustrates the case of *one* call-back (or a *second* contact); the notation now involves an average  $\bar{y}_1$  for the sample (or *first* contact) and an average  $\bar{y}_2$  for the sub-sample of the first contact non-respondents (*i.e.*, *second* contacts or *first* call-backs). In the following symbolic development, it is assumed that *complete* response is obtained from the second contact and that samples are obtained by EPS.



Rewriting equation (3.4c.1) in appropriate notation, we have:

$$\bar{Y}_s = \frac{N}{N_s} \bar{Y} + \frac{N_{nr}}{N_s} \bar{Y}_{nr}. \quad \text{-----(3.4c.5)}$$

Remembering the Plan's responsibility to make it reasonable to treat the (real-world)  $y_s$  as the (model)  $y_s$ , we estimate  $\bar{Y}$  by  $\bar{y}_1 (= \bar{Y}_1)$  for the sample of respondents to the first contact, and  $\bar{Y}_{nr}$  by  $\bar{y}_2 (= \bar{Y}_2)$  for complete response from the sub-sample of the non-respondents (the corresponding random variables are  $\bar{Y}_1$  and  $\bar{Y}_2$ ); hence, our estimator of  $\bar{Y}_s$  is given by:

$$\bar{Y}_{cb} = \frac{N}{N_s} \bar{Y}_1 + \frac{N_{nr}}{N_s} \bar{Y}_2 \quad (\text{assuming the response rate – and the non-response rate – are known}). \quad \text{-----(3.4c.6)}$$

$$\text{Then: } E(\bar{Y}_{cb}) = \frac{N}{N_s} E(\bar{Y}_1) + \frac{N_{nr}}{N_s} E(\bar{Y}_2) = \frac{N}{N_s} \bar{Y} + \frac{N_{nr}}{N_s} \bar{Y}_{nr} = \bar{Y}_s \quad (\text{under equiprobable selecting}); \quad \text{-----(3.4c.7)}$$

*i.e.*, under EPS and if *complete* response is obtained at the *final* (here, the first) call-back,  $\bar{Y}_{cb}$  is unbiased with respect to non-response for  $\bar{Y}_s$ ; the corresponding *estimate* is  $\bar{y}_{cb}$ .

The foregoing results can be extended to *multiple* call-backs if we think of the non-respondent population as a *set* of *sub-populations*, often of progressively *decreasing* size and requiring an *increasing* level of coercion to respond; increasing coercion would be represented in practice by factors such as the experience and persistence of the interviewer, the incentive(s) offered for response, and the fact of multiple contacts by the interviewer. Each successive contact after the initial one is usually with a *sub-sample* of the non-respondents from the previous contact.

To accommodate this model of the non-respondent population into the symbolic development given above, we would extend the notation to define  $\bar{y}_r$  as the average of the responses from the *r*th contact [or (r-1)st call-back]. Ideally, the call-back process is repeated until *complete* response is obtained from the relevant sub-sample of the previous non-respondents; in practice, the process is continued until either of or some combination of:

- the sub-population of remaining non-respondents is considered too small to introduce appreciable non-responding bias into the estimate of  $\bar{Y}_s$  – this can be thought of as achieving the condition  $N_{nr} \ll N$  discussed overleaf on page 3.21;
- the use of further resources for call-backs is not considered cost effective in terms of the sample survey answers.

A final matter about our model for the use of call-backs to reduce the effect of non-responding bias when using  $\bar{Y}_{cb}$  as an estimator of  $\bar{Y}_s$  is that if the response rate,  $N/N_s$  [and the *non-response* rate,  $N_{nr}/N_s$ ] are *unknown* and are estimated by  $n/n_s$  (and  $n_{nr}/n_s$ ) in the expression (3.4c.6) for  $\bar{Y}_{cb}$ , taking the expectation [as in equation (3.4c.7)] involves  $E(N/n_s)$  and  $E(N_{nr}/n_s)$  [ $N$  and  $N_{nr}$  are the respective random variables representing  $n$  and  $n_{nr}$ , which can be shown to be  $N/N_s$  and  $N_{nr}/N_s$ , respectively, and similarly for multiple call-backs. Then, the estimator of  $\bar{Y}_s$  is:

$$\bar{Y}_{cb} = \frac{N}{n_s} \bar{Y}_1 + \frac{N_{nr}}{n_s} \bar{Y}_2 \quad (\text{if the response rate – and the non-response rate – are estimated}); \quad \text{-----(3.4c.8)}$$

$$\text{but: } E\left(\frac{N}{n_s} \bar{Y}_1\right) = E_N[E_{\bar{Y}_1|N}\left(\frac{N}{n_s} \bar{Y}_1\right)] = E_N\left[\frac{N}{n_s} E_{\bar{Y}_1|N}(\bar{Y}_1)\right] = E_N\left[\frac{N}{n_s} \bar{Y}\right] = E_N\left[\frac{N}{n_s}\right] \bar{Y} = \frac{N}{N_s} \bar{Y}; \quad \text{-----(3.4c.9)}$$

$$\text{similarly: } E\left(\frac{N_{nr}}{n_s} \bar{Y}_2\right) = \frac{N_{nr}}{N_s} \bar{Y}_{nr}, \quad \text{-----(3.4c.10)}$$

so that, as in equation (3.4c.7) above where the response rate – and the *non-response* rate – are *known*, under EPS and if *complete* response is obtained at the *final* (here, the first) call-back,  $\bar{Y}_{cb}$  is likewise *unbiased* with respect to non-response for  $\bar{Y}_s$  when the two rates are *estimated* from the selection.

**NOTES:** 1. The current context of the *study* population and (*our* terminology of) the *selection*, together with their attributes  $N_s$ ,  $\bar{Y}_s$ ,  $n_s$  and  $\bar{y}_s$ , makes equations (3.4c.8) to (3.4c.10) *different* from what we routinely encounter in the usual context of the *respondent* population and (*our* meaning of) the *sample* – in *that* context, there is seldom a reason to involve the *non-respondent* population and the non-respondents, again *unlike* the current discussion.

Here,  $n_s$  is under the control of the investigator(s), whereas  $n$  and  $n_{nr}$  are *not* and, under EPS, we model them as values of the corresponding random variables  $N$  and  $N_{nr}$ .

(continued)

**Figure 3.4c. SAMPLE SURVEY DESIGN/EXECUTION: Models for Non-Response (continued 1)**

- NOTES:** 1. ● Non-responding bias, the focus of equations (3.4c.8) to (3.4c.10) and the random variables  $N$  and  $N_{nr}$ , is peripheral to the key ideas the two diagrams at the upper right on pages 3.22 and 3.21 are intended to convey; this is why, to avoid distraction, roman, *not* italic, letters are used for  $n$  and  $n_{nr}$  in the two diagrams.
- 2 The use of conditional expectation in equations (3.4c.9) and (3.4c.10) on the lower half of the facing page 3.22, to show that the expected value of the product of two random variables is equal to the product of their expected values, reminds us of three matters:
- Equation (3.4c.9) is reminiscent of the conditional probability result for events  $A$  and  $B$ :  

$$\Pr(A \cap B) = \Pr(A) \times \Pr(B|A).$$
  - We usually associate multiplicative behaviour of expected values with *probabilistic independence*; however, such independence is a *strong* requirement in that its occurrence imposes substantial constraints on the relevant probability functions.
  - Multiplicative behaviour of expected values can occur under *weaker* conditions – here, that  $E_{\bar{Y}_1|N}(\bar{Y}_1)$  is *not* a function of  $N$  – than probabilistic independence, as illustrated in equation (3.4c.9); we are familiar with this idea from the fact that neither zero covariance nor zero correlation necessarily implies probabilistic independence.

**3. Dealing with Non-Response – Stratifying**

A second approach to try to reduce non-response is to *stratify* the study population into groups (*i.e.*, strata) within *each* of which respondents and non-respondents are *more* similar with respect to the attribute(s) of interest than are respondents and non-respondents in the study population as a whole.

A summary of some possible reasons for non-response is:

1. Non-contact:
  - inaccurate information about the physical location of a element (or unit);
  - accessibility of elements (or units) – *e.g.*, homemakers vs. university students;
  - timing of contact attempt(s) – *e.g.*, weekdays, weeknights, weekends;
  - number of contact attempts;
  - competence of interviewer – *e.g.*, locating people who have moved residence.
2. Refusal to participate in the sample survey:
  - participation is aided by interviewer's level of satisfaction in working with people, intellectual curiosity, professional identification with the sample survey's objectives, etc.
3. Item non-response (some questionnaire items not answered):
  - question is sensitive or personal;
  - respondent does not know the answer to a question.

**Blank page**