# Responsive Design – Side Effect Reduction of Prior Information on Survey Design

## A. Demnati

Independent Researcher, Ottawa, Canada, Abdellatif_Demnati@msn.com

**Abstract**

It is difficult to design a survey because prior information on response rates and the like is likely generated from a different random process than the target one governing the survey to be designed, and the impact on the properties of the estimator can be significant. We are concerned with reducing the side effect of error-prone prior information on the precision of the estimator of the parameter of interest. Nowadays, computer-assisted data collection methods provide an instant variety of observations on the target random process governing the survey under consideration. These data and paradata enable the survey producer to make decisions regarding the need for methodology-process revision, which involves consideration of both a model that represents how the target information relates to the prior information and the design that describes how the observations are obtained. We think of the prior and target information as a random process that has a joint distribution with some probability function. Then at each phase of data collection, after receiving the information that the target random process has taken specific values, we update the joint probability distribution, to revise the design specification in the course of the data collection period. A coefficient of reliability for a survey as a whole set of processes as well for a single process is also discussed.

**Key Words:** Adaptive design; Misclassification; Optimal resources allocation; Paradata; Partially observed units; Two sources of information.

## 1. Introduction

Survey or census studies start with a collection of distinct units of interest; termed population. There are multiple random variables attached to each unit, as each holds their own individual characteristics and attitudes. Each particular study targets a small subset of these random variables and involves a questionnaire to collect the data of sampled respondents in order to draw inferences at the estimation stage about the parameter of interest. Measurements on these variables of interest are intended to be collected during data collection stage from each selected unit.

At the planning stage of a survey the question of resources determination and allocation within stages of the survey design is a difficult and critical one. Survey developers must justify resources to be used, and the survey producer should review the justification to ensure the survey produces results within resources, quality and timing constraints. Efficiency is a very important issue because inefficient allocation may lead to: a) imprecise results; and b) waste of time, resources, and money. To determine optimally : i) the duration of the survey; and ii) the amount of resources and their allocation within stages of the survey design, design pre-specification requires the following: a) specification of the parameter of interest, and the associated estimator to be used; b) specification of the desired precision or the global cost; c) specification of the cost function; d) specification of the precision function; e) obtaining prior information, from the sampling frame, administrative files, or from previous surveys, needed to compute unknown quantities in formulas for both precision and cost functions; and, f) optimization of some utility function – that perhaps involves both precision and cost functions.

Suppose previous surveys suggest that the conditional probability of responding $h^{(rq)}$ (in a time period) is constant over time, where the superscript "$rq$" stands for "response to questionnaire". When the conditional response probability $h^{(rq)}$ is constant over time, then the marginal response probability over $I$ time periods is given by $\xi_{1;k}^{(rq)} = 1 - (1 - h^{(rq)})^{I}$. To reach a marginal probability of response close to 1 under constant conditional response probability, it will take around 17 time periods when $h^{(rq)} = .5$, and over 100 time periods when $h^{(rq)} = .1$. Since collecting data over such a long data collection period is time consuming, costly and the results may vary from one time period to another, the way survey sampling handles the problem of capturing information from, or estimating parameters with respect to, finite population generated from such random processes is as follows: a) selecting a random sample of units from the population; and b) increasing level of efforts in term of follow-up treatments to improve units cooperation. Sampling is based on the idea that, within a certain margin of error, one can infer something about the parameter of interest from a small sample, as long as the sample is chosen at random. Efficient follow-up requires information on the error-free

target random process governing the survey under consideration. It is difficult to pre-specify the design for certain surveys because prior information is likely generated from a different random process than the one under consideration. A naive approach simplifies the problem under the assumption that resources should be big enough to have good estimates. However, often a survey has limited budget and timing, and those in turn, in combination with the resources allocation used within stages of the survey design based on prior information, determine the achievable precision. Nowadays, computer-assisted data collection methods provide an instant variety of observations about the target random process that can be used to revise survey design during the course of its data collection. Although, previous survey designs are mostly done deterministically using prior information, there is a wide spread need for responsive design where the design is revised during the data collection period. The intent of such revision is to reduce errors attached to design pre-specification on prior information grounds. Objectives of responsive design are formulated in Groves and Heeringa (2006). Starting with an expected design based on prior information, then, cumulative collected data can be used to: 1) update information used for design specification; and, 2) revise, if necessary, specification of the design at each phase of data collection. Estimates with a certain margin of error can be obtained using few time periods, since prior information provides information about the joint probability distribution.

Design pre-specification is a special case of measurement error which refers here to the case where the prior (or error-prone) information, say $\chi$, is not necessarily identical to the target (or error-free) information, say $\psi$, of the processes underlying the finite population. We assume that the assessment of error in $\chi$ can be carried out on the basis of observations on $\psi$. We also assume that the error-prone information $\chi$ has a potential bias $b$ when used to estimate $\psi$ and that the error-free information $\psi$ has no error. Thus, the assessment of errors allows quantification of such bias. Under two random processes, we are interested in the error-free random variable $\psi$, knowing its probability density function, the probability density function of another random variable $\chi$, together with the joint probability density function with vector parameter denoted by $\lambda$. We assume that census parameter $\lambda_N$ associated with the vector model parameter $\lambda$ is defined as solution to an estimating equation (EE) of the form $\mathbf{S}(\psi, \chi; \lambda) = \sum_k \mathbf{s}(\psi_k, \chi_k; \lambda) = \mathbf{0}$, where $\sum_k$ is the sum over all the population units. It is assumed that the sampling frame has no coverage bias. It is also assumed that values of the error-prone variable are available for all units in the population, while values of the error-free variable are unknown but observable.

Once an estimate of $\lambda$, $\psi_k$, or of the parameter of interest is obtained, the question follows; what is the reliability of this estimate? In a general sense, reliability of an estimate refers to the degree to which the estimate is free from error and therefore truly measures the quantity which it is intended to measure. When reliability measures are available at all various stages of the survey process, they can serve as performance measures. Such measures enable the survey producer to make decisions regarding the need for methodology-process modification. As there is no general reliability measure that would capture all information on the impact of each stage of the survey design on the ultimate estimate, the survey producer tends to combine various measures to get a broader effect and interactions between different factors. A key step in defining reliability was the introduction of an error criterion that measures, in a probabilistic sense, the error between the desired quantity $\theta$ and an estimate $\hat{\theta}$ of it. Possible sources of error in surveys include sampling frame, sampling scheme, nonresponse, measurement, disclosure-avoidance, etc. A criterion which is commonly used in judging the performance of an estimator $\hat{\theta}$ of a quantity $\theta$ is its Mean Square Error (MSE) defined by $MSE(\hat{\theta}) = E\{(\hat{\theta} - \theta)^2\}$, where $E$ denotes total expectation under random processes involved. We can also interpret this via the MSE decomposition. For any random variable $z$, we have $E(z^2) = E\{[z - E(z)]^2\} + \{E(z)\}^2$. Applying this to $z = \hat{\theta} - \theta$ we get

$$E\{(\hat{\theta} - \theta)^2\} = E\{[(\hat{\theta} - \theta) - E(\hat{\theta} - \theta)]^2\} + \{E(\hat{\theta} - \theta)\}^2 . \tag{1.1}$$

The first term of (1.1) is the variance of $\hat{\theta} - \theta$. It is the error of the estimator due to the random processes involved. The second term of (1.1) is the square of the bias of $\hat{\theta}$, the best one can do is make this zero.

In an attempt to discuss side effect reduction of prior information on the ultimate estimator, our work below is organized as follows: in Section 2 stochastic models underlying nonresponse follow-up strategy, mode of data collection, and response behaviour are presented; in Section 3, basic ingredients for designing a survey are defined; in Section 4, design pre-specification and revision in the course of data collection is studied; in Section 5, estimator of the regression parameter as well as estimator of the parameter of interest are derived; in Section 6, surveys requiring units contactibility and eligibility identification are briefly discussed; and, in Section 7, a reliability coefficient is discussed.

## 2. Stochastic Models

To revise the survey design in discrete intervals, we divided the continuous time of the entire data collection period into a sequence of continuous time periods: 1, 2, and so on, and let's $I_{min}$ denote the minimum length of data collection period to obtain full responses. Suppose the survey limited length of duration of data collection is made up of $P_{max}$ phases, the $p^{th}$ being of size $n_p$ time periods, so that the limited duration of data collection is made up of $I_{max} = \sum_{p=1}^{P_{max}} n_p$ time periods, with $I_{max} < I_{min}$. So we shall be dealing with $N \times P_{max}$ rectangular array of phases of data collection, where $N$ is the size of the finite population. In this Section, we consider stochastic models for nonresponse follow-up strategy, mode of data collection, and response behaviour. Such models are based on the paradigm of a random experiment modeled by a probability measure on an underlying sample space; i.e., an experiment whose outcome cannot be predicted with certainty, before the experiment is run.

## 2.1 Nonresponse Follow-up Strategy Model

Suppose we have $S$ nonresponse follow-up strategies, and define a vector of strategy indicator variables as $J_{s|p;k}^{(f)} = 1$ if unit $k$ is assigned to strategy $s$ at phase $p$, and $J_{s|p;k}^{(f)} = 0$ if not, where $J_{|p;k}^{(f)} = (J_{1|p;k}^{(f)}, ..., J_{S|p;k}^{(f)})^T$ are realizations of independent distributed variables according to a multinomial distribution, $Mult_S(1, \boldsymbol{\varphi}_{|p;k}^{(f)})$, $\boldsymbol{\varphi}_{|p;k}^{(f)} = (\phi_{1|p;k}^{(f)}, ..., \phi_{S|p;k}^{(f)})^T = E_f(\boldsymbol{J}_{|p;k}^{(f)})$ is the vector of strategy probabilities with $\sum_{s=1}^{S} \phi_{s|p;k}^{(f)} = 1$, $E_f$ denotes expectation under the strategy allocation model, and the superscript "$f$" stands for "follow-up". We consider the $S^{th}$ strategy as an omitted or reference strategy. For the multinomial logistic regression model, logits of the first $S-1$ strategies are constructed with the reference strategy in the denominator

$$\log(\phi_{s|p;k}^{(f)} / \phi_{S|p;k}^{(f)}) = \boldsymbol{v}_{f|p;k}^T \boldsymbol{\eta}_{|p}, \ s = 1, ..., S-1,$$

where $\boldsymbol{v}_{f|p;k}$ is the $q_{f|p}^{(1)} \times 1$ vector of explanatory variables and $\boldsymbol{\eta}_{|p} = (\boldsymbol{\eta}_{1|p}^T, ..., \boldsymbol{\eta}_{(S-1)|p}^T)^T$ is the $q_{f|p} = q_{f|p}^{(1)}(S-1) \times 1$ unknown vector parameter to be estimated. It follows that the $S$ conditional probabilities given the vector of explanatory variables are

$$\phi_{S|p;k}^{(f)} = \{1 + \sum_{s=1}^{S-1} \exp(\boldsymbol{v}_{f|p;k}^T \boldsymbol{\eta}_{s|p})\}^{-1},$$

and for $s = 1, ..., S-1$
$$\phi_{s|p;k}^{(f)} = \phi_{S|p;k}^{(f)} \exp(\boldsymbol{v}_{f|p;k}^T \boldsymbol{\eta}_{s|p}).$$

Let's $^{(t)}c_{s|p;k}^{(f)}$ and $^{(e)}c_{s|p;k}^{(f)}$ denote respectively the error-free and error-prone cost associated with follow-up strategy $s$ ($s = 1, ..., S$) for unit $k$ during phase $p$, where the superscripts "$t$" and "$e$" stand for "error-free" and "error-prone" respectively.

## 2.2 Data Collection Models

Similarly, suppose we have $M$ modes of data collection, and define a vector of error-free data collection mode indicator variables as $^{(t)}J_{m;k}^{(dc)} = 1$ if unit $k$ uses mode $m$, and $^{(t)}J_{m;k}^{(dc)} = 0$ if not, where $^{(t)}\boldsymbol{J}_k^{(dc)} = (^{(t)}J_{1;k}^{(dc)}, ..., ^{(t)}J_{M;k}^{(dc)})^T$ are realizations of independent distributed random variables according to a multinomial distribution, $Mult_M(1, ^{(t)}\boldsymbol{\varphi}_k^{(dc)})$, $^{(t)}\boldsymbol{\varphi}_k^{(dc)} = (^{(t)}\phi_{1;k}^{(dc)}, ..., ^{(t)}\phi_{M;k}^{(dc)})^T = E_{dc}(^{(t)}\boldsymbol{J}_k^{(dc)})$ is the vector of data collection mode probabilities with $\sum_{m=1}^{M} {}^{(t)}\phi_{m;k}^{(dc)} = 1$, and $E_{dc}$ denotes expectation with respect to the data collection model. We consider the $M^{th}$ mode as a reference mode. For the multinomial logistic regression model, logits of the first $M-1$ modes are constructed with the reference mode in the denominator

$$\log(^{(t)}\phi_{m;k}^{(dc)} / ^{(t)}\phi_{M;k}^{(dc)}) = \boldsymbol{v}_{dc;k}^T \, ^{(t)}\boldsymbol{\alpha}_m, \ m = 1, ..., M-1,$$

where $\boldsymbol{v}_{dc;k}$ is the $q_{dc}^{(1)} \times 1$ vector of explanatory variables and $^{(t)}\boldsymbol{\alpha} = (^{(t)}\boldsymbol{\alpha}_1^T, ..., ^{(t)}\boldsymbol{\alpha}_{M-1}^T)^T$ is the $^{(t)}q_{dc} = q_{dc}^{(1)}(M-1) \times 1$ unknown vector parameter to be estimated. It follows that the $M$ conditional probabilities given the vector of explanatory variables are

$$^{(t)}\phi_{M;k}^{(dc)} = \{1 + \sum_{m=1}^{M-1} \exp(\boldsymbol{v}_{dc;k}^T \, ^{(t)}\boldsymbol{\alpha}_m)\}^{-1},$$

and for $m = 1, ..., M-1$
$$^{(t)}\phi_{m;k}^{(dc)} = {}^{(t)}\phi_{M;k}^{(dc)} \exp(\boldsymbol{v}_{dc;k}^T \, ^{(t)}\boldsymbol{\alpha}_m).$$

Similarly, the vector of error-prone data collection mode indicator variables $^{(e)}\boldsymbol{J}_k^{(dc)} = (^{(e)}J_{1;k}^{(dc)}, ..., ^{(e)}J_{M;k}^{(dc)})^T$ are realizations of independent distributed random variables according to a multinomial distribution, $Mult_M(1, ^{(e)}\boldsymbol{\varphi}_k)$, $^{(e)}\boldsymbol{\varphi}_k^{(dc)} = (^{(e)}\phi_{1;k}^{(dc)}, ..., ^{(e)}\phi_{M;k}^{(dc)})^T = E_{dc}(^{(e)}\boldsymbol{J}_k^{(dc)})$ is the vector of probabilities with $\sum_{m=1}^{M} {}^{(e)}\phi_{m;k}^{(dc)} = 1$. The error-prone data collection mode

for unit $k$ is characterized by the matrix $^{(e|t)}\mathbf{\Phi}_k^{(dc)}$ which consists of the conditional probability $^{(e|t)}\phi_{i|j;k}^{(dc)}$ of assigning mode $i$ given that unit $k$ should have been assigned mode $j$ defined by its components

$$^{(e|t)}\phi_{i|j;k}^{(dc)} = \Pr(^{(e)}J_{i;k}^{(dc)} = 1|^{(t)}J_{j;k}^{(dc)} = 1) .$$

The marginal distribution of $^{(e)}J_{i;k}^{(dc)}$ is:

$$^{(e)}\phi_{i;k}^{(dc)} = \Pr(^{(e)}J_{i;k}^{(dc)} = 1) = \sum_{j=1}^{M} {}^{(e|t)}\phi_{i|j;k}^{(dc)} {}^{(t)}\phi_{j;k}^{(dc)} , \quad i = 1,...,M ,$$

where $\sum_{i=1}^{M} {}^{(e|t)}\phi_{i|j;k}^{(dc)} = 1$. We consider the $M^{th}$ mode as the reference mode. For the multinomial logistic regression model, logits of the first $M-1$ conditional modes are constructed with the reference mode in the denominator

$$\log(^{(e|t)}\phi_{i|j;k}^{(dc)} / {}^{(e|t)}\phi_{M|j;k}^{(dc)}) = {}^{(|j)}\mathbf{v}_{dc;k}^{T} {}^{(|j)}\mathbf{\alpha} , \quad i = 1,...,M-1 ,$$

where $^{(|j)}\mathbf{v}_{dc;k}$ is the $^{(|j)}q \times 1$ vector of explanatory variables, $^{(|j)}\mathbf{\alpha} = (^{(1|j)}\mathbf{\alpha}^T,...,^{(M-1|j)}\mathbf{\alpha}^T)^T$ and $^{(e)}\mathbf{\alpha} = (^{(|1)}\mathbf{\alpha}^T,...,^{(|M)}\mathbf{\alpha}^T)^T$ is the $^{(e)}q_{dc} \times 1$ unknown vector parameter to be estimated, with $^{(e)}q_{dc} = \sum_{j=1}^{M} {}^{(|j)}q(M-1)$. It follows that the $M$ conditional probabilities of each mode given mode $j$ and the vector of explanatory variables are

$$^{(e|t)}\phi_{M|j;k}^{(dc)} = \{1 + \sum_{m=1}^{M-1} \exp(^{(|j)}\mathbf{v}_{dc;k}^{T} {}^{(|j)}\mathbf{\alpha}_m)\}^{-1} ,$$

and for $m = 1,...,M-1$

$$^{(e|t)}\phi_{m|j;k}^{(dc)} = {}^{(e|t)}\phi_{M|j;k}^{(dc)} \exp(^{(|j)}\mathbf{v}_{dc;k}^{T} {}^{(|j)}\mathbf{\alpha}_m) .$$

Let's $^{(t)}c_{m;k}^{(dc)}$ and $^{(e)}c_{m;k}^{(dc)}$ denote respectively the error-free and error-prone cost associated with data collection mode $m$ ( $m = 1,...,M$ ) for unit $k$ .

## 2.3 Response Models

Let $^{(a)}t^{(rq)}$ represent the discrete random variable that indicates the time period $i$ when the response occurs under random process $a$ with $a \in \{e,t,e|t\}$ for a randomly selected unit from the sample, where $e$, $t$, and $e|t$ are used to indicate error-prone, error-free and conditional random processes respectively. After $P$ phases (or equivalently $I_P = \sum_{p=1}^{P} n_p$ time periods) of data collection, each unit $k$ is observed from the $e_k^{(dc)}$ time period until the period $^{(a)}I_{P;k}$, with $^{(a)}I_{P;k} \leq I_P$, where $e_k^{(dc)}$ denotes the entry time period for unit $k$ into data collection window. Usually $e_k^{(dc)}$ equals 0. Observation of the unit could be discontinued for two reasons: 1) the unit responds; or 2) the phase $P$ of data collection ends. In the first case, $^{(a)}t_k^{(rq)} = {}^{(a)}I_{P;k}$. In the second case, it is only know that $^{(a)}t_k^{(rq)} > I_P$. Units with $^{(a)}t_k^{(rq)} > I_P$ are right-censored – it is unknown when they respond. Note that $^{(a)}t^{(rq)}$ is defined only when the unit will respond eventually using appropriate follow-up strategy. Since censoring is planned and observation is terminated at the end of data collection, the censuring mechanism is noninformative (Lagakos, 1979) in the sense that the act of censoring imparts no information about the response mechanism. The overall response indicator from the $e_k^{(dc)}$ to the $I_P$ time period, with $e_k^{(dc)} < I_P$, is given by $^{(a)}r_{P;k}^{(rq)} = 1 - \prod_{i=e_k^{(dc)}+1}^{I_{P;k}}(1 - {}^{(a)}r_{ki}^{(rq)}) = \sum_{i=e_k^{(dc)}+1}^{I_{P;k}} {}^{(a)}r_{ki}^{(rq)}$ , where $^{(a)}r_{ki}^{(rq)}$ is a sequence of response indicators defined for each unit $k$ whose values are defined as $^{(a)}r_{ki}^{(rq)} = 1$ if the unit does respond in period $i$, and $^{(a)}r_{ki}^{(rq)} = 0$ if the unit does not respond in period $i$. Because response occurrence is intrinsically conditional, Demnati (2015) characterized $^{(a)}t^{(rq)}$ by its conditional probability function – the distribution of the probability that a response will occur in each time period given that it has not already occurred in a previous time period – known as the discrete-time hazard function. Discrete-time hazard, $^{(a)}h_{ki}^{(rq)}(^{(a)}\mathbf{v}_{rq;k}, {}^{(a)}\mathbf{\beta})$, $^{(a)}h_{ki}^{(rq)}$ for short, is defined as the conditional probability that unit $k$ will respond in time period $i$, given that the unit did not respond prior to $i$:

$$^{(a)}h_{ki}^{(rq)} = \Pr(^{(a)}t_k^{(rq)} = i|^{(a)}t_k^{(rq)} \geq i) ,$$

where $^{(a)}\mathbf{v}_{rq;k}$ refers to both time-invariant and time-varying explanatory variables and $^{(a)}\mathbf{\beta}$ is the unknown $^{(a)}q_{rq} \times 1$ vector parameter to be estimated. For unit with $^{(a)}t_k^{(rq)} = i$, the probability of obtaining a response at time period $i$ could be expressed in terms of the hazard as

$$\Pr(^{(a)}t_k^{(rq)} = i) = {}^{(a)}h_{ki}^{(rq)} \prod_{j=e_k^{(dc)}+1}^{i-1}(1 - {}^{(a)}h_{kj}^{(rq)}) .$$

For units with $^{(a)}t_k^{(rq)} > i$, the probability of obtaining a response can be expressed as

$$\Pr(^{(a)}t_k^{(rq)} > i) = \prod_{j=e_k^{(dc)}+1}^{i}(1 - {}^{(a)}h_{kj}^{(rq)}) .$$

After $P$ phases of data collection, we have

$$g(^{(a)}t_k^{(rq)}) = \{\Pr(^{(a)}t_k^{(rq)} = {}^{(a)}I_{P;k})\}^{^{(a)}\delta_k^{(dc)}} \{\Pr(^{(a)}t_k^{(rq)} > {}^{(a)}I_{P;k})\}^{1 - {}^{(a)}\delta_k^{(dc)}} \equiv {}^{(a)}g_{I_P;k}^{(rq)} ,$$

where $^{(a)}\delta_k^{(dc)}=1$ if unit $k$ is uncensored (responds) and $^{(a)}\delta_k^{(dc)}=0$ if unit $k$ is censored under process $a$. When unit $k$ is censored, it is unknown when the unit responds. The joint distribution of $(^{(e)}t_k^{(rq)},{}^{(t)}t_k^{(rq)})$ is characterized by the matrix $^{(e|t)}\mathbf{G}_k^{(rq)}$ defined by its components

$$^{(e|t)}g_{i|j;k}^{(rq)}=\Pr(^{(e)}t_k^{(rq)}=i|^{(t)}t_k^{(rq)}=j)\,,\ i,j=e_k^{(dc)}+1,...,\mathrm{I_P}\,.$$

The marginal distribution of $^{(e)}t_k^{(rq)}$ is:

$$^{(e)}g_{\mathrm{I_P};k}^{(rq)}=\sum_{j=e_k^{(dc)}+1}^{\mathrm{I_P}}{}^{(e|t)}g_{\mathrm{I_P};k|j;k}^{(rq)}\,{}^{(t)}g_{j;k}^{(rq)}\,,\ \mathrm{I_{P;k}}=e_k^{(dc)}+1,...,\mathrm{I_P}\,.$$

The marginal probability of obtaining a response from the $e_k^{(dc)}$ to the $\mathrm{I_P}$ time period is given by

$$^{(a)}\xi_{\mathrm{P};k}^{(rq)}=1-\prod_{i=e_k^{(dc)}+1}^{\mathrm{I_P}}(1-{}^{(a)}h_{ki}^{(rq)})=\sum_{i=e_k^{(dc)}+1}^{\mathrm{I_P}}\Pr(^{(a)}t_{ki}^{(rq)}=i)\,.$$

## 2.4 Modeling Sample Selection Probability

For sampling, we parameterize the probability of selection as $\pi_k=\{lb+ub\times\exp(\mathbf{v}_{\wp;k}^T\varphi)\}/\{1+\exp(\mathbf{v}_{\wp;k}^T\varphi)\}$, where $lb$ and $ub$ are respectively the lower and upper bounds with $0<lb<ub\le1$, $\mathbf{v}_{\wp;k}$ is the $q_\wp\times1$ vector of explanatory variables and $\varphi$ is the $q_\wp\times1$ unknown vector parameter to be determined. Unlike the customary probability of selection $\pi_{c;k}=np_k$, where $n$ is the expected sample size and $p_k$ is a measure of size, this expression for $\pi_k$ fulfill the two criterions: : a) $0<lb\le\pi_k\le ub\le1$; and b) $\pi_k\to ub$ as $n\to N$. It is common practice to set $(lb,ub)=(0,1)$.

# 3. Basic Ingredients for Designing a Survey

In this Section, we elaborate on each step of the technical steps mentioned in the introduction. These steps are required in designing a survey. We consider examples of parameter of interest and an associated estimator, a precision function, a cost function, and an utility function.

## 3.1 Parameter of Interest

We assume that census parameter $\Theta_N(\mathbf{\psi})$ associated with the variable of interest $\mathbf{\psi}$ is defined as solution to an EE of the form

$$\mathbf{S}(\mathbf{\psi};\mathbf{\Theta})=\sum_k\mathbf{s}(\mathbf{\psi}_k;\mathbf{\Theta})-\mathbf{v}(\mathbf{\Theta})=\mathbf{0}\,,\tag{3.1}$$

where the known function $\mathbf{s}(\mathbf{\psi}_k;\mathbf{\Theta})$ is a $q_\Theta$-dimensional vector-valued function of $\mathbf{\psi}_k$ and the known function $\mathbf{v}(\mathbf{\Theta})$ allows for explicitly defined parameters. For linear and logistic regression models, $\mathbf{s}(\mathbf{\psi}_k;\mathbf{\Theta})=\mathbf{x}_k(y_k-\mu_k(\mathbf{x}_k^T\mathbf{\Theta}))$ and $\mathbf{v}(\mathbf{\Theta})=\mathbf{0}$, where $\mu_k(\mathbf{x}^T\mathbf{\Theta})=E_y(y)$, $\mathbf{x}=(x_1,...,x_{q_\Theta})^T$ is a $q_\Theta\times1$ vector of explanatory variables, $\mathbf{\Theta}=(\Theta_1,...,\Theta_{q_\Theta})^T$ is the $q_\Theta\times1$ vector of model parameter and $E_y$ denotes model expectation on the variable of interest $y$. For the special case of the finite population total $Y=\sum_k y_k$, $\mathbf{s}(\mathbf{\psi}_k;\mathbf{\Theta})=y_k$, $\mathbf{v}(\mathbf{\Theta})=\Theta_N(\mathbf{\psi})$, and $\Theta_N(\mathbf{\psi})=Y$. The finite population parameter $\Theta_N(\mathbf{\psi})$, obtained as the solution of (3.1), under the assumed ideal situation which consists of census case with complete response and without any measurement error, plays the role of a "gold standard".

## 3.2 Estimator of the Parameter of Interest

Suppose that the response probability, $\xi_{\mathrm{P};k}^{(rq)}=E_r(r_{\mathrm{P};k}^{(rq)})$, during $\mathrm{P}$ phases of data collection, is known for every unit in the population, where $r_{\mathrm{P};k}^{(rq)}$ is the overall response indicator for unit $k$, and $E_r$ denotes expectation under the response mechanism. For general sampling design with known positive inclusion probabilities, $\pi_k$, a design-response unbiased estimator of the EE defined by (3.1) is given by

$$\bar{\mathbf{S}}(\mathbf{\psi};\mathbf{\Theta})=\sum_k d_k(\wp)(r_{\mathrm{P};k}^{(rq)}/\xi_{\mathrm{P};k}^{(rq)})\mathbf{s}(\mathbf{\psi}_k;\mathbf{\Theta})-\mathbf{v}(\mathbf{\Theta})=\mathbf{0}\,,\tag{3.2}$$

where $d_k(\wp)=1_k(\wp)/\pi_k$ are the design weights, $1_k(\wp)=1(k\in\wp)$ is the sample $\wp$ membership indicator variable for unit $k$, $1(condition)$ is the truth function, i.e., $1(condition)=1$ if the $condition$ is true and $1(condition)=0$ if not, $\pi_k=E_\wp\{1_k(\wp)\}$ is the sample $\wp$ inclusion probability for unit $k$, and $E_\wp$ denotes expectation with respect to the sampling design. The solution obtained by a Newton-Raphson-type iterative method gives the estimator $\bar{\mathbf{\Theta}}_{\mathrm{P}}(\mathbf{\psi})$ of $\mathbf{\Theta}_N(\mathbf{\psi})$.

## 3.3 Derivation of the Variance Function

We assume that $P$ phases of data collection are completed, with $1 \le P \le P_{max}$, and we consider for illustration the simple estimator $\breve{\Theta}_{P;k}(\psi)$ of $\Theta_N(\psi)$ solution to (3.2). We first consider the derivation of the variance of a compact form given by

$$\hat{U} = \sum_k u_k d_k(\wp) r_{P;k}^{(rq)},$$ (3.3)

where $u_k$ is a vector of constants. We may decompose the variance of $\hat{U}$ as

$$Var(\hat{U}) = E_\wp Var_r(\hat{U}) + Var_\wp E_r(\hat{U}) \equiv \mathbf{V}_r + \mathbf{V}_\wp = \mathbf{V}(u),$$ (3.4)

where $Var_\wp$ and $Var_r$ denote variance with respect to sampling design and response mechanism respectively.

Under independent mechanism on $r_{P;k}$, the first component $\mathbf{V}_r = E_\wp Var_r(\hat{U})$ of $Var(\hat{U})$ given by (3.4) is given by

$$\mathbf{V}_r = \sum_k (1/\pi_k) u_k \xi_{P;k}^{(rq)} (1 - \xi_{P;k}^{(rq)}) u_k^T.$$ (3.4a)

The second component $\mathbf{V}_\wp = Var_\wp \{\sum_k u_k d_k(\wp) \xi_{P;k}^{(rq)}\}$ of $Var(\hat{U})$ given by (3.4) is given by

$$\mathbf{V}_\wp = \sum_k \sum_l \omega_{kl}^{-1}(1 - \omega_{kl}) u_k \xi_{P;k}^{(rq)} \xi_{P;l}^{(rq)} u_l^T,$$ (3.4b)

where $\omega_{kl} = \pi_{kl}^{-1} \pi_k \pi_l$, $\omega_{kk} = \omega_k = \pi_k$ and $\pi_{kl} = E_\wp\{1_k(\wp)1_l(\wp)\}$. The sum of (3.4a) and (3.4b) constitutes $\mathbf{V}(u) = \mathbf{V}_r + \mathbf{V}_\wp$, the variance of $\hat{U}$ given by (3.3).

It follows that the compact form given by (3.3) can be used to derive the linearization variance $Var_L\{\breve{\Theta}_P(\psi)\}$ of $\breve{\Theta}_P(\psi)$ using $u_k = \{J(\Theta_N(\psi))\}^{-1} s(\psi_k; \Theta_N(\psi))/\xi_{P;k}^{(rq)}$, where $J(\Theta) = -\partial S^T(\psi; \Theta)/\partial \Theta$. It is easily shown that, under stratified simple random sample and Poisson sampling, we can express the linearization variance $Var_L\{\breve{\Theta}_P(\psi)\}$ of $\breve{\Theta}_P(\psi)$ in the separate form as

$$Var_L\{\breve{\Theta}_P(\psi)\} = v_0(\psi) + \sum_k v_{\wp;k}(\psi)/\pi_k + \sum_k v_{\wp r;k}(\psi)/\{\pi_k \xi_{P;k}^{(rq)}\},$$

where $v_0(\psi)$, $v_{\wp;k}(\psi)$, and $v_{\wp r;k}(\psi)$ are functions independents of $\pi_k$ and $\xi_{P;k}^{(rq)}$. We denote in operator notation the estimator $\breve{\Theta}_P(\psi)$ and its linearization variance $Var_L\{\breve{\Theta}_P(\psi)\}$ by $\breve{\Theta}\{\lambda, \mathbf{A}(\psi)\}$ and $Var_L\{\breve{\Theta}\{\lambda, \mathbf{A}(\psi)\}\}$ respectively, where $\mathbf{A}(\psi)$ is an $N$-column matrix with $k^{th}$ column $\psi_k$.

## 3.4 Specification of the Cost Function
We may decompose the global cost over $P$ phases of data collection as

$$C(P) = \sum_{p=1}^{P} C_p,$$ (3.5)

with $1 \le P \le P_{max}$. The $P$ components of the global cost are

$$C_1 = c_1 + C_1^{(\wp)} + C_1^{(f)} + C_1^{(dc)},$$

and for $p = 2, ..., P$

$$C_p = c_p + C_p^{(f)} + C_p^{(dc)},$$

where $C_1^{(\wp)} = \sum_k 1_k(\wp) c_k^{(\wp)}$ is the component associated with sampling cost, $C_p^{(f)} = \sum_k 1_k(\wp)(1 - r_{p-1;k}^{(rq)}) \sum_{s=1}^{S} J_{s|p;k}^{(f)} c_{s|p;k}^{(f)}$ is the component associated with nonresponse follow-up cost for phase $p$, with $r_{0;k}^{(rq)} = 0$, and $C_p^{(dc)} = \sum_k 1_k(\wp)(1 - r_{p-1;k}^{(rq)}) \sum_{s=1}^{S} J_{s|p;k}^{(f)} \sum_{m=1}^{M} J_{m|s;k}^{(dc)} r_{p;k}^{(rq)} c_{m|p;k}^{(dc)}$ is the component associated with data collection cost for phase $p$. Here $c_p$ is the fixed cost for phase $p$, and $c_k^{(\wp)}$ is the sampling cost for unit $k$.

## 3.5 Specification of the Utility Function
Suppose guessed values $\lambda_p$ and $\mathbf{A}(\psi_p)$ are available for $\lambda$ and for $\mathbf{A}(\psi)$ respectively. To create a design in the case of one parameter of interest, we minimize the linearization variance, $Var_L\{\breve{\Theta}\{\lambda_p, \mathbf{A}(\psi_p)\}\}$ subject to constraints on the duration of the data collection and on the conditional expected global cost: $1 \le P \le P_{max}$ and $C\overline{C}\{\lambda_p, \mathbf{A}(\psi_p)\} \le C_{max}$, where $C\overline{C}\{\lambda_p, \mathbf{A}(\psi_p)\}$ denotes the conditional expectation of the cost function $C(P)$ given by (3.5), and $C_{max}$ is the survey global cost limit.

In the case of $\Lambda(>1)$ parameters of interest, let $\Theta_\kappa$ denote the parameter of interest $\kappa$, $\kappa = 1, ..., \Lambda$. To create a design, we optimize the conditional expected global cost given by $C\overline{C}\{\lambda_p, \mathbf{A}(\psi_p)\}$ subject to constraints on the duration of data collection and on $\Lambda$ variances:

$$1 \le P \le P_{max},$$

and
$$Var_L\{\breve{\Theta}_\kappa\{\lambda_p, \mathbf{A}(\psi_p)\}\} \le V_{p;\kappa}, \quad \kappa = 1,...,\Lambda,$$

where $V_{p;\kappa}$ are specified tolerances, and $Var_L\{\breve{\Theta}_\kappa\{\lambda_p, \mathbf{A}(\psi_p)\}\}$ is the linearization variance of the estimator $\breve{\Theta}_\kappa\{\lambda_p, \mathbf{A}(\psi_p)\}$ for the $\kappa^{th}$ parameter of interest $\kappa = 1,...,\Lambda$. For example, one could specify an upper limit, $\mathfrak{I}_{p;\kappa}$, on the coefficient of variation of $\breve{\Theta}_\kappa\{\lambda_p, \mathbf{A}(\psi_p)\}$ so that $V_{p;\kappa} = \{\mathfrak{I}_{p;\kappa} E\{\breve{\Theta}_\kappa\{\lambda_p, \mathbf{A}(\psi_p)\}\}\}^2$. One may repeat the optimization process with different value of $\mathfrak{I}_{p;\kappa}$ $\kappa = 1,...,\Lambda$ to obtain the desired minimum cost.

## 4. Design Pre-specification and Revision

The complete data for sampled unit $k$ is given by $(\psi_k^T, \chi_k^T)^T$, where $\psi_k = {}^{(t)}\varpi_k$, $\chi_k = {}^{(e)}\varpi_k$, ${}^{(\varepsilon)}\varpi_k = ({}^{(\varepsilon)}\mathbf{I}_k, {}^{(\varepsilon)}\delta_k^{(dc)}, {}^{(\varepsilon)}\mathbf{J}_k^{(dc)T}, {}^{(\varepsilon)}\mathbf{y}_k^T, {}^{(\varepsilon)}\mathbf{v}_k^T, {}^{(\varepsilon)}\mathbf{c}_k^T)^T$, ${}^{(\varepsilon)}\mathbf{c}_k = ({}^{(\varepsilon)}\mathbf{c}_k^{(f)T}, {}^{(\varepsilon)}\mathbf{c}_k^{(dc)T})^T$, and $\varepsilon \in \{t,e\}$. The vectors parameter associated with the marginal distributions of ${}^{(\varepsilon)}\varpi_k$ is ${}^{(\varepsilon)}\lambda = ({}^{(\varepsilon)}\alpha^T, {}^{(\varepsilon)}\beta^T, {}^{(\varepsilon)}\gamma_y^T, {}^{(\varepsilon)}\gamma_c^T)^T$, where ${}^{(\varepsilon)}\gamma_y$ is the vector parameter associated with the vector of variables of interest ${}^{(\varepsilon)}\mathbf{y}$, and ${}^{(\varepsilon)}\gamma_c$ is the vector parameter associated with the vector cost ${}^{(\varepsilon)}\mathbf{c}$. We set ${}^{(t)}\lambda_0 = {}^{(e)}\lambda$ and $\psi_{0;k} = \chi_k$ for design pre-specification. In this Section, we discuss survey design pre-specification, and the Observation-Revision-Optimization steps for design revision in the course of its data collection period.

### 4.1 Pre-specification of the Survey Design
Suppose guessed values are available for design pre-specification. Since the above global cost given by (3.5) under the models comes from what we think of as random, we consider its expectation under the sampling design and models for strategy, mode and response behaviour, given by
$$E\{C(\mathrm{P})\} = \sum_{p=1}^{\mathrm{P}} \overline{C}_p \equiv C\overline{C}\{\lambda, \mathbf{A}(\psi)\}. \tag{4.1}$$
The $\mathrm{P}$ components of the expected global cost are
$$\overline{C}_1 = c_1 + \overline{C}_1^{(\wp)} + \overline{C}_1^{(f)} + \overline{C}_1^{(dc)},$$
and for $p = 2,...,\mathrm{P}$
$$\overline{C}_p = c_p + \overline{C}_p^{(f)} + \overline{C}_p^{(dc)},$$
where $\overline{C}_1^{(\wp)} = \sum_k \pi_k c_k^{(\wp)}$, $\overline{C}_p^{(f)} = \sum_k \pi_k (1 - \xi_{p-1;k}^{(rq)}) \sum_{s=1}^S \phi_{s|p;k}^{(f)} c_{s|p;k}^{(f)}$, and $\overline{C}_p^{(dc)} = \sum_k \pi_k (1 - \xi_{p-1;k}^{(rq)}) \sum_{s=1}^S \phi_{s|p;k}^{(f)} \sum_{m=1}^M \phi_{m|s;k}^{(dc)} \xi_{p;k}^{(rq)} c_{m|p;k}^{(dc)}$, with $\xi_{0;k}^{(rq)} = 0$.

To create a pre-design using guessed values $\lambda_0$, and $\mathbf{A}(\psi_0)$, the conditional expected cost $C\overline{C}\{\lambda_0, \mathbf{A}(\psi_0)\}$, and the variance $Var_L\{\breve{\Theta}\{\lambda_0, \mathbf{A}(\psi_0)\}\}$, we determine the optimal $\mathrm{P}$, $\varphi$, and $\eta$ by optimizing the utility function. We denote the solution by $\varphi$, $\mathrm{P}_0$ and $\eta_0$, then we draw the sample $\wp$ using $\varphi$.

### 4.2 Revision of a Design in the Course of its Progress
Our method basically consists of a series of optimizations using additional information on the error-free target random process. After optimization using additional information from a certain phase of data collection is completed, the result indicates if the active design should be revised. We do not, therefore, use a fixed design, although an expected design is always specified in the survey design. Starting with an expected pre-specified design, then for $p = 1,2,...$ design update is made using the Observation-Revision-Optimization (ORO) steps:

➢ **Observation Step:** Obtain next phase $p$ of observations on the error-free process.

➢ **Revision Step:**
- **Maximization**: Update $\lambda_{p-1}$ to get $\lambda_p$ using $\mathbf{D}_p$ and the update step given in Subsection 4.3, where $\mathbf{D}_p$ denotes all observed information until the end of phase $p$ of data collection.
- **Imputation**: Impute missing values of each component $\psi$ of $\psi$ to get $\psi_{p;k} = E_\psi(\psi_k | \mathbf{D}_p, \lambda_p)$, where $E_\psi$ denotes expectation with respect to the random process governing the component $\psi$. Note that $\psi_{p;k} = \psi_k$ when item $\psi_k$ is observed.
- **Conditional Expectation of the Cost Function:** Compute the conditional expectation of the global cost to get $C\overline{C}\{\lambda_p, \mathbf{A}(\psi_p)\}$
- **Conditional Expectation of the Precision Function:** Compute the conditional expectation of the variance to get $Var_L\{\breve{\Theta}\{\lambda_p, \mathbf{A}(\psi_p)\}\}$.

➢ **Optimization Step:**

- Determine the optimal $P$, and $\eta$ conditional on $\varphi$, $J_{1;k}^{(f)},...,J_{p\text{-}1;k}^{(f)}$, $r_{p\text{-}1;k}^{(rq)}$, $\lambda_p$, and $\mathbf{A}(\psi_p)$ under the usual constraints and the following additional constraint $p \le P \le P_{max}$, using $CC\{\lambda_p, \mathbf{A}(\psi_p)\}$ and $Var_L\{\bar{\Theta}_\kappa\{\lambda_p, \mathbf{A}(\psi_p)\}\}$. The solution is denoted by $P_p$, and $\eta_p$.

The three steps are repeated until $p = P_p$.

## 4.3 Revision – Maximization Step

When the duration of data collection period is taken into account, the likelihood function of the joint distribution under census data is defined for unit $k$ as

$$L_{I_{min};k}^{(I_{min})}(\lambda) = f_{I_{min}}^{(I_{min})}(\psi_k, \chi_k),$$ 

(4.1)

where the subscript $I$ in $f_I^{(I_{min})}(\zeta)$ denotes that $\zeta$ is observed during the interval $[0,I]$. Hence when data collection period is taken into account, the census case means that $I = I_{min}$ and $d_k(\wp) = 1_k(\wp) = \pi_k = 1$. To simplify our notation we drop the superscript $I_{min}$, and write (4.1) as $L_{I_{min};k}(\lambda) = f_{I_{min}}(\psi_k, \chi_k)$. After observing $I_p$ time periods of data collection, the joint observations on $\chi$, and $\psi$ are known for respondents during the $I_p$ time periods, while only observations on $\chi$ are known for nonrespondents during the rest of the periods. Consequently, we decompose the likelihood of observed data for unit $k$ in two parts

$$L_{I_p;k}(\lambda) = f_{I_{max}}(\chi_k) f_{I_p}(\psi_k \mid \chi_k),$$ 

(4.2)

in which case the log-likelihood is given by

$$\ell_{I_p;k}(\lambda) = \log f_{I_{max}}(\chi_k) + \log f_{I_p}(\psi_k \mid \chi_k),$$ 

(4.3)

where $f_{I_{max}}(\chi_k) = \int f_{I_{max}}(\chi_k, \psi_k) d\psi_k$, and $f_{I_p}(\psi_k \mid \chi_k) = f_{I_p}(\chi_k, \psi_k)/f_{I_p}(\chi_k)$. Note that $f_{I_{max}}(\chi_k) f_{I_p}(\psi_k \mid \chi_k) \to f_{I_{max}}(\psi_k, \chi_k)$ as $I_p \to I_{max}$. In arriving to (4.2), we decomposed $\chi$ in two parts: $\chi = (\chi_{[e_k^{(dc)}, I_p]}^T, \chi_{[I_p+1, I_{max}]}^T)^T$, where $\chi_{[e_k^{(dc)}, I_p]}$ denotes observation during the interval $[e_k^{(dc)}, I_p]$, while $\chi_{[I_p+1, I_{max}]}$ denotes observation during the interval $[I_p + 1, I_{max}]$. The joint distribution is given by

$$f(\chi, \psi) = f(\chi_{[e_k^{(dc)}, I_p]}, \chi_{[I_p+1, I_{max}]}, \psi) = f(\chi_{[e_k^{(dc)}, I_p]}) f(\chi_{[I_p+1, I_{max}]} \mid \chi_{[e_k^{(dc)}, I_p]}) f(\psi \mid \chi_{[e_k^{(dc)}, I_p]}) = f(\chi) f(\psi \mid \chi_{[e_k^{(dc)}, I_p]}).$$

Taking the derivatives of (4.3) and adjusting for unequal probability of selection, we get the weighted EE

$$\hat{\mathbf{S}}_{I_p}(\psi, \chi; \lambda) = \sum_k d_k(\wp)\{\mathbf{s}_{I_{max}}(\chi_k; \lambda) + \mathbf{s}_{I_p}(\psi_k; \lambda \mid \chi_k)\} = \mathbf{0},$$ 

(4.4)

where 

$$\mathbf{s}_{I_{max}}(\chi_k; \lambda) = \partial \log f_{I_{max}}(\chi_k)/\partial\lambda \text{ and } \mathbf{s}_{I_p}(\psi_k; \lambda \mid \chi_k) = \partial \log f_{I_p}(\psi_k \mid \chi_k)/\partial\lambda.$$

The estimator $\hat{\lambda}_p$ of $\lambda$ is obtained using the following update step.

**Update Step of $\lambda$ :** Starting with a guessed value, $\lambda^{(0)} = \lambda_{p\text{-}1}$, then for $b = 1,2,...$ updates are made using

$$\lambda^{(b)} = \lambda^{(b-1)} + \{\hat{\mathbf{J}}_{I_p}(\lambda^{(b-1)})\}^{-1}\hat{\mathbf{S}}_{I_p}(\psi, \chi; \lambda^{(b-1)}),$$

where $\hat{\mathbf{J}}_{I_p}(\lambda) = -\partial\hat{\mathbf{S}}_{I_p}^T(\psi, \chi; \lambda)/\partial\lambda$. The solution to (4.4) gives the estimator $\hat{\lambda}_p$ of $\lambda$.

To update only the vector parameter $\beta$ associated with the response mechanism after observing $p$ phases of data collection, we set $\chi_k = {}^{(e)}t_k^{(rq)}$, $\psi_k = {}^{(t)}t_k^{(rq)}$, $\lambda = \beta$,

$$f_1(\chi_k) = \prod_{i=1}^I \{\sum_{j=1}^I {}^{(e|t)}g_{i|j;k}^{(rq)(t)} g_{j;k}^{(rq)}\}^{(e)t_{i;k}^{(rq)}},$$

and 

$$f_{I_p}(\psi_k, \chi_k) = \prod_{i=1}^{I_p} \prod_{j=1}^{I_p} \{{}^{(e|t)}g_{i|j;k}^{(rq)(t)} g_{j;k}^{(rq)}\}^{(e)t_{i;k}^{(rq)(t)}t_{j;k}^{(rq)}},$$

where ${}^{(a)}t_{i;k}^{(rq)} = 1({}^{(a)}t_k^{(rq)} = i)$ for $i = 1,..., I_{max}$.

After observing the first phase of data collection, the joint distribution is given by $f(\psi, \chi) = f_{I_{max}}(\chi) f_{I_1}(\psi \mid \chi)$, and the conditional distribution of $\psi$ given $\chi$ is given by

$$f_{I_1}(\psi \mid \chi) = \frac{f_{I_1}(\chi \mid \psi)}{f_{I_1}(\chi)} f_{I_1}(\psi),$$

where the factor $\{f_{I_1}(\chi \mid \psi)/f_{I_1}(\chi)\}$ represents the impact of the error-prone information $\chi$ on the distribution of the target error-free information $\psi$. After observing two phases of data collection, we may decompose the conditional distribution as

$$f_{I_2}(\psi \mid \chi) = f_{I_1}(\psi \mid \chi) \times \frac{f(\chi_{[I_1+1,I_2]} \mid \psi_{[e_k^{dc}+1,I_1]}, \psi_{[I_1+1,I_2]})}{f(\chi_{[I_1+1,I_2]} \mid \psi_{[e_k^{dc}+1,I_1]})} f(\psi_{[I_1+1,I_2]} \mid \psi_{[e_k^{dc}+1,I_1]}),$$

where the factor $\{f(\chi_{[I_1+1,I_2]} \mid \psi_{[e_k^{dc}+1,I_1]}, \psi_{[I_1+1,I_2]})\}/\{f(\chi_{[I_1+1,I_2]} \mid \psi_{[e_k^{dc}+1,I_1]})\}$ represents the partial impact of the information $\chi_{[I_1+1,I_2]}$ during phase 2 on the conditional distribution of $\psi_{[I_1+1,I_2]}$ during phase 2 given observation $\psi_{[e_k^{dc}+1,I_1]}$ during phase 1, provided $e_k^{dc} < I_1$.

## 4.4 Revision – Imputation Step

Our interest here is in estimating the variable $\psi$ by an estimator $\theta$ using the well known mean square approach. One may estimate the unknown quantity $\psi$ as follows: a) set the lost function to the square error $SE(\psi,\theta) = (\psi - \theta)^2$; then, b) pick estimate $\theta$ to minimize $E_\psi\{SE(\theta,\psi)\}$. If all of the available information on $\psi$ is summarised in its distribution $f(\psi)$, then we need to solve $\min_\theta E_\psi\{SE(\theta,\psi)\} = \min_\theta \int (\psi - \theta)^2 f(\psi) d\psi$. Differentiating with respect to $\theta$ gives the optimal estimate $\theta = E_\psi(\psi)$. Thus, the case of no additional information the minimum MSE estimator is simply the expected value $\theta = E_\psi(\psi)$. When we have additional information in the form of observed information $\chi$ that is related somehow to $\psi$. We could use that information to get a better estimation than its mean. A simple model assumes that $\chi = \psi + b$, where $b$ is a random variable that represents error. Suppose that $\psi$ and $b$ are independent, $\psi$ has a normal distribution with mean $\mu_\psi$ and variance $\sigma_\psi^2$, and $b$ has a normal distribution with mean 0 and variance $\sigma_b^2$. Then the conditional expectation of $\psi$ given $\chi$ is $E_\psi(\psi \mid \chi) = \mu_\psi + \sigma_\psi^2(\chi - \mu_\psi)/(\sigma_\psi^2 + \sigma_b^2)$, and the conditional variance of $\psi$ given $\chi$ is given by $Var_\psi(\psi \mid \chi) = \sigma_\psi^2\{1 - \sigma_\psi^2/(\sigma_\psi^2 + \sigma_b^2)\}$. When the variance of the error is smaller, $\psi$ becomes closer to $\chi$, i.e. $Var_\psi(\psi \mid \chi)$ becomes closer to 0. Hence, the conditional variance of $\psi$ given $\chi$ decreases as $\chi$ becomes closer to $\psi$. More generally, if the form of the joint distribution $f(\chi,\psi)$ is known, then our estimator of $\psi$ given $\chi$ is of course some function of $\chi$, say $\theta(\chi)$ and our aim now is to minimize $\min_{\theta(\chi)} E_\psi[SE\{\psi,\theta(\chi) \mid \chi\}] = \min_{\theta(\chi)} \int \{\psi - \theta(\chi)\}^2 f(\psi \mid \chi) d\psi$. Exactly the same calculations as in the case of no additional information then show that $\theta(\chi) = E_\psi(\psi \mid \chi)$, the conditional expectation of $\psi$ given $\chi$. The updated probability distribution of $\psi$ is the conditional probability distribution of $\psi$ given $\chi$. If $\psi$ and $\chi$ are independent, then all conditional expectation of $\psi$ are independent of $\chi$, and coincide with the unconditional expectation $E_\psi(\psi)$.

## 4.5 Revision – Cost Conditional Expectation Step

Consider for simplicity the case of $P = 2$. In this case, the global cost given by (3.5) reduces to $C(P) = C_1 + C_2$. After observing the first phase of data collection, the first part $C_1 = c_1 + C_1^{(\wp)} + C_1^{(f)} + C_1^{(dc)}$ is known, while the second part $C_2 = c_2 + C_2^{(f)} + C_2^{(dc)}$ is partially unknown; and the resulting conditional expected cost is given by

$$C\overline{C}\{\lambda_1, \mathbf{A}(\psi_1)\} = C_1 + C\overline{C}_2,$$

with
$$C\overline{C}_2 = c_2(\lambda_1) + C\overline{C}_2^{(f)} + C\overline{C}_2^{(dc)},$$

where $C\overline{C}_2^{(f)} = \sum_k 1_k(\wp)(1 - r_{1;k}^{(rq)}) \sum_{s=1}^S \phi_{s|2;k}^{(f)}(\lambda_1) c_{s|2;k}^{(f)}(\mathbf{A}(\psi_1))$, and $C\overline{C}_2^{(dc)} = \sum_k 1_k(\wp)(1 - r_{1;k}^{(rq)}) \sum_{s=1}^S \phi_{s|2;k}^{(f)}(\lambda_1) \sum_{m=1}^M \phi_{m|s;k}^{(dc)}(\lambda_1) \xi_{2;k}^{(rq)}(\lambda_1) c_{m|2;k}^{(dc)}(\mathbf{A}(\psi_1))$. Note that the remaining cost is given by $C_{\max} - C_1$.

## 5. Estimation

## 5.1 Estimation of the Regression Parameter

Once data collection is completed, observed values of the target information in combination with prior information values are used in the EE given by (4.4) at the estimation stage to get estimate of regression parameter. However it may happens that prior information is ignored in the estimation stage, and only observed values of the target information are used. In this case, the likelihood for unit $k$ is given by

$$L_{P;k}(\lambda) = f_P(\psi_k),$$

Taking the derivatives of the log-likelihood, and adjusting for sampling unequal probabilities, we get the weighted EE

$$\hat{\mathbf{S}}_P(\psi;\lambda) = \sum_k d_k(\wp) \mathbf{s}(\psi_k;\lambda) = \mathbf{0}, \tag{5.1}$$

where $\mathbf{s}(\psi_k;\lambda) = \partial \log f_P(\psi_k)/\partial \lambda$. Starting with a guessed value, $\lambda^{(0)}$, then for $b = 1,2,\dots$ updates are made using

$$\lambda^{(b)} = \lambda^{(b-1)} + \{\hat{\mathbf{J}}_{\mathrm{P}}(\mathbf{\psi};\lambda^{(b-1)})\}^{-1}\hat{\mathbf{S}}_{\mathrm{P}}(\mathbf{\psi};\lambda^{(b-1)}) \,,$$

where $\hat{\mathbf{J}}_{\mathrm{P}}(\mathbf{\psi};\lambda) = -\partial\hat{\mathbf{S}}_{\mathrm{P}}^{T}(\mathbf{\psi};\lambda)/\partial\lambda$. The solution of (5.1) obtained by a Newton-Raphson-type iterative method gives the estimator $\hat{\lambda}$ of $^{(t)}\lambda$. Let's consider our case which consists on models for strategy, mode and response behaviour. Since, the mode of data collection is known for a unit that responds at any time of data collection, while it is unknown for a unit that is censored, the likelihood for unit $k$ may be decomposed as

$$L_{\mathrm{P};k}(\lambda) = f(\mathbf{J}_{k}^{(f)})L_{R|\mathbf{J}_{k}^{(f)};k}L_{M|\mathbf{J}_{k}^{(f)};k} \,,$$

with $\qquad L_{R|\mathbf{J}_{k}^{(f)};k} = \{\prod_{m=1}^{M}[\phi_{m|\mathbf{J}_{k}^{(f)};k}^{(dc)} f_{m|\mathbf{J}_{k}^{(f)}}(t_{k}^{(rq)})]^{J_{m|s;k}^{(dc)}}\}^{\delta_{k}^{(dc)}}$ and $L_{M|\mathbf{J}_{k}^{(f)};k} = \{\sum_{m=1}^{M}\phi_{m|\mathbf{J}_{k}^{(f)};k}^{(dc)} f_{m|\mathbf{J}_{k}^{(f)}}(t_{k}^{(rq)})\}^{(1-\delta_{k}^{(dc)})}$,

where $\lambda = (\mathbf{\eta}^{T},\mathbf{\alpha}^{T},\mathbf{\beta}^{T})^{T}$, $\mathbf{J}_{k}^{(f)} = (\mathbf{J}_{1|k}^{(f)T},...,\mathbf{J}_{\mathrm{P}|k}^{(f)T})^{T}$ is the stochastic process representing the evolution of nonresponse follow-up strategies over data collection phases, $\mathbf{J}_{\mathrm{p}|k}^{(f)}$ indicates which follow-up strategy is assigned to unit $k$ at phase p, and $f_{m|\mathbf{J}_{k}^{(f)}}(t_{k}^{(rq)})$ is $f(t_{k}^{(rq)})$ for mode $m$ of data collection under follow-up stochastic process $\mathbf{J}_{k}^{(f)}$. It remains to specify the joint probability distribution of the random process $\mathbf{J}_{k}^{(f)} = (\mathbf{J}_{1|k}^{(f)T},...,\mathbf{J}_{\mathrm{P}|k}^{(f)T})^{T}$. It maybe informative in our context to use Markov chains model with $s$ states. Markov chains model is characterized by – the probability of the next follow-up strategy depends only on the current follow-up strategy and not on the sequence of follow-up strategies that proceed it – the matrix which consists of the conditional probability defined by its components

$$\Pr(J_{i|\mathrm{p}+1;k}^{(f)} = 1 \mid J_{j|\mathrm{p};k}^{(f)} = 1) \,, \; i,j = 1,...,S \,, \text{ and } \; \mathrm{p} = 1,...,\mathrm{P}-1 \,.$$

An excellent source of information on Marcov Chains estimation is the volume by Brémaud (1998).

## 5.2 Estimation of the Parameter of Interest

Since the probability of response is unknown, estimated response probability $\hat{\xi}_{\mathrm{P};k} = \xi_{\mathrm{P};k}(\hat{\lambda})$ is used in the EE given by (3.2) to get

$$\hat{\mathbf{S}}(\mathbf{\psi};\mathbf{\Theta}) = \sum_{k} d_{k}(\wp)(r_{\mathrm{P};k}^{(rq)}/\hat{\xi}_{\mathrm{P};k}^{(rq)})\mathbf{s}(\mathbf{\psi}_{k};\mathbf{\Theta}) - \mathbf{v}(\mathbf{\Theta}) = \mathbf{0} \,, \tag{5.2}$$

and the solution gives the estimator $\hat{\mathbf{\Theta}}_{\mathrm{P}}(\mathbf{\psi})$ of $\mathbf{\Theta}_{N}(\mathbf{\psi})$. As noted by Rosenbaum (1987) and others, estimator $\hat{\mathbf{\Theta}}_{\mathrm{P}}(\mathbf{\psi})$ using the estimated response probability can be more efficient than estimator $\breve{\mathbf{\Theta}}_{\mathrm{P}}(\mathbf{\psi})$ using the true response probability.

If prior information values are used, then the likelihood for unit $k$ in given by

$$L_{\mathrm{P};k}(\mathbf{\Theta}) = \{f_{\mathrm{I}_{\max}}(\mathbf{\chi}_{k})\}^{1-1_{k}(\wp)n_{\mathrm{P};k}} f_{\mathrm{I}_{\mathrm{P}}}(\mathbf{\psi}_{k},\mathbf{\chi}_{k})^{1_{k}(\wp)n_{\mathrm{P};k}} \,,$$

in which case the weighed EE is given by

$$\hat{\mathbf{S}}_{\mathrm{P}}(\mathbf{\psi},\mathbf{\chi};\mathbf{\Theta}) = \sum_{k}\mathbf{s}_{\mathrm{I}_{\max}}(\mathbf{\chi}_{k};\mathbf{\Theta}) + \sum_{k} d_{k}(\wp)(r_{\mathrm{P};k}/\hat{\xi}_{\mathrm{P};k})\mathbf{s}_{\mathrm{I}_{\mathrm{P}}}(\mathbf{\psi}_{k};\mathbf{\Theta} \mid \mathbf{\chi}_{k}) = \mathbf{0} \,,$$

where $\qquad \mathbf{s}_{\mathrm{I}_{\max}}(\mathbf{\chi}_{k};\mathbf{\Theta}) = \partial\log f_{\mathrm{I}_{\max}}(\mathbf{\chi}_{k})/\partial\mathbf{\Theta}$ and $\mathbf{s}_{\mathrm{I}_{\mathrm{P}}}(\mathbf{\psi}_{k};\mathbf{\Theta} \mid \mathbf{\chi}_{k}) = \partial\log f_{\mathrm{I}_{\mathrm{P}}}(\mathbf{\psi}_{k} \mid \mathbf{\chi}_{k})/\partial\mathbf{\Theta}$.

For the simple model defined in Section 4.4, where the vectors $(\chi_{k},\psi_{k})^{T}$ are realizations of independent distributed random variables according to a bivariate Normal distribution,

$$\begin{pmatrix}\chi_{k}\\\psi_{k}\end{pmatrix} \sim N_{2}\left(\begin{pmatrix}\mu_{\psi}\\\mu_{\psi}\end{pmatrix},\begin{pmatrix}\sigma_{\psi}^{2}+\sigma_{b}^{2} & \sigma_{\psi}^{2}\\\sigma_{\psi}^{2} & \sigma_{\psi}^{2}\end{pmatrix}\right),$$

we have $\mathbf{s}_{\mathrm{I}_{\max}}(\chi_{k};\mathbf{\Theta}) = (\chi_{k}-\mu_{\psi})\rho_{\psi\chi}^{2}/\sigma_{\psi}^{2}$, and $\mathbf{s}_{\mathrm{I}_{\mathrm{P}}}(\psi_{k};\mathbf{\Theta}\mid\chi_{k}) = \{\psi_{k}-\mu_{\psi}-\rho_{\psi\chi}^{2}(\chi_{k}-\mu_{\psi})\}/\{(1-\rho_{\psi\chi}^{2})\sigma_{\psi}^{2}\}$, provided $\rho_{\psi\chi}^{2}<1$, where $\rho_{\psi\chi}^{2} = \sigma_{\psi}^{2}/(\sigma_{\psi}^{2}+\sigma_{b}^{2})$. If $\rho_{\psi\chi} = 0$, then the EE reduces to $\hat{\mathbf{S}}_{\mathrm{P}}(\mathbf{\psi},\mathbf{\chi};\mathbf{\Theta}) = \sum_{k} d_{k}(\wp)(r_{\mathrm{P};k}^{(rq)}/\hat{\xi}_{\mathrm{P};k}^{(rq)})(\psi_{k}-\mu_{\psi})/\sigma_{\psi}^{2} = 0$.

## 6. Quick Look at Surveys Requiring Contactibility and Eligibility Identification

The process of determining the contactibility and the eligibility status of each sampled unit is another stage of the survey design where the side effect of prior information can be reduced during its phase of follow-up to establish a first contact. In the following, we sketch the necessitate ingredients.

## 6.1 Contactibility Models
Let $^{(a)}t^{(fc)}$ represent the discrete random variable that indicates the time period $i$ when the first contact occurs under random process $a$, where the superscript "$fc$" stands for "first contact". The process of follow-up to reach unit $k$ at

specified time periods is conducted until some period $^{(a)}I_k^{(fc)}$, with $^{(a)}I_k^{(fc)} \leq I_{max}^{(fc)}$, where $I_{max}^{(fc)}$ is the maximum time periods for follow-up to get a first contact, and $I_{max}^{(fc)} < I_{max}$. The process of reaching a first contact could be discontinued for two reasons: 1) the unit is contacted for the first time; or 2) the period to reach sampled unit ends. In the first case, $^{(a)}t_k^{(fc)} = ^{(a)}I_k^{(fc)}$. In the second case, it is only known that $^{(a)}t_k^{(fc)} > I_{max}^{(fc)}$. Units with $^{(a)}t_k^{(fc)} > I_{max}^{(fc)}$ are right-censored – it is unknown whether they are contactable during survey data collection period. The overall contactibility indicator over $I_{max}^{(fc)}$ time periods is given by $^{(a)}r_k^{(fc)} = 1 - \prod_{i=1}^{I_k^{(fc)}} (1 - ^{(a)}r_{ki}^{(fc)})^{z_{ki}^{(fc)}}$, where $^{(a)}r_{ki}^{(fc)}$ is a sequence of contactibility indicators defined for each unit $k$ whose values are defined as $^{(a)}r_{ki}^{(fc)} = 1$ if the unit is contacted in time period $i$, and $^{(a)}r_{ki}^{(fc)} = 0$ if not, and $^{(a)}z_{ki}^{(fc)}$ denote the 0/1 variable indicating whether a unit $k$ is to be contacted at time period $i$. Because contactibility occurrence is intrinsically conditional, we characterized $^{(a)}t^{(fc)}$ by its discrete-time hazard function

$$^{(a)}h_{ki}^{(fc)} = \Pr(^{(a)}t_k^{(fc)} = i | ^{(a)}t_k^{(fc)} \geq i),$$

where $^{(a)}h_{ki}^{(fc)} = ^{(a)}h_{ki}^{(fc)}(^{(a)}\boldsymbol{v}_{fc;k}, ^{(a)}\boldsymbol{\beta}^{(fc)})$, $^{(a)}\boldsymbol{v}_{fc;k}$ refers to both time-invariant and time-varying explanatory variables and $^{(a)}\boldsymbol{\beta}^{(fc)}$ is the unknown $^{(a)}q_{fc} \times 1$ vector parameter to be estimated. For unit with $^{(a)}t_k^{(fc)} = i$, the probability of obtaining a first contact at time period $i$ could be expressed in terms of the hazard as

$$\Pr(^{(a)}t_k^{(fc)} = i) = ^{(a)}z_{ki}^{(fc)(a)}h_{ki}^{(fc)} \prod_{j=1}^{i-1}(1 - ^{(a)}h_{kj}^{(fc)})^{^{(a)}z_{ki}^{(fc)}}. \tag{6.1}$$

For units with $^{(a)}t_k^{(fc)} > i$, the probability of obtaining a response can be expressed as

$$\Pr(^{(a)}t_k^{(fc)} > i) = \prod_{j=1}^{i}(1 - ^{(a)}h_{kj}^{(fc)})^{^{(a)}z_{kj}^{(fc)}}.$$

We have

$$g(^{(a)}t_k^{(fc)} = ^{(a)}I_k^{(fc)}) = \{\Pr(^{(a)}t_k^{(fc)} = ^{(a)}I_k^{(fc)})\}^{^{(a)}\delta_k^{(fc)}} \{\Pr(^{(a)}t_k^{(fc)} > ^{(a)}I_k^{(fc)})\}^{1 - ^{(a)}\delta_k^{(fc)}},$$

where $^{(a)}\delta_k^{(fc)} = 1$ if unit $k$ is uncensored (contacted) and $^{(a)}\delta_k^{(fc)} = 0$ if unit $k$ is censored under process $a$. When unit $k$ is censored, either unit $k$ will be reachable at some future time period $^{(a)}t_k^{(fc)} > I_{max}^{(fc)}$ or the unit will not be contactable during survey period. The joint distribution of $(^{(e)}t_k^{(fc)}, ^{(t)}t_k^{(fc)})$ is characterized by the matrix $^{(e|t)}\mathbf{G}_k^{(fc)}$ defined by its components

$$^{(e|t)}g_{i|j;k}^{(fc)} = \Pr(^{(e)}t_k^{(fc)} = i | ^{(t)}t_k^{(fc)} = j), \quad i,j = 1,...,I_{max}^{(fc)}.$$

The marginal probability of obtaining a first contact after $I_{max}^{(fc)}$ time periods is given by

$$\Pr(^{(a)}r_k^{(fc)} = 1) = 1 - \prod_{i=1}^{I_{max}^{(fc)}}(1 - ^{(a)}h_{ki}^{(fc)})^{^{(a)}z_{ki}}.$$

## 6.2 Parameter of Interest and Associated Estimator

We now define the census parameter for eligible subpopulation as the solution to

$$\mathbf{S}(\boldsymbol{\psi}; \boldsymbol{\Theta}) = \sum_k J_k^{(el)} \mathbf{s}(\boldsymbol{\psi}_k; \boldsymbol{\Theta}) - \mathbf{v}(\boldsymbol{\Theta}) = \mathbf{0}, \tag{6.2}$$

where $J_k^{(el)}$ denotes the eligibility indicator for unit $k$, i.e., $J_k^{(el)} = 1$ if unit $k$ is eligible for the given survey and $J_k^{(el)} = 0$ if unit $k$ is not eligible. Suppose that the probability, $\xi_k = E(\Delta_k)$, of getting a response to questionnaire during P phases of follow-ups to both establish a contact and get response to questionnaire is known for every eligible unit in the population, where $\Delta_k = r_k^{(fc)} r_k^{(rq)}$, $r_k^{(fc)}$ is the overall contactibility indicator, and $r_k^{(rq)}$ is the overall response indicator for eligible unit $k$. A design-response unbiased estimator of the EE defined by (6.2) is given by

$$\breve{\mathbf{S}}(\boldsymbol{\psi}; \boldsymbol{\Theta}) = \sum_k J_k^{(el)} d_k(\wp)(\Delta_k / \xi_k) \mathbf{s}(\boldsymbol{\psi}_k; \boldsymbol{\Theta}) - \mathbf{v}(\boldsymbol{\Theta}) = \mathbf{0}, \tag{6.3}$$

where the probability of getting response to questionnaire for eligible unit $k$ is given by

$$\xi_k = \sum_{i=1}^{I_{max}^{(fc)}} \Pr(t_k^{(fc)} = i)\{1 - \prod_{j=i+1}^{I_{max}}(1 - h_{kj}^{(rq)})\},$$

and $\Pr(t_k^{(fc)} = i)$ is given by (6.1).

The variance of the estimator obtained as the solution to (6.3) can be derived along the lines of Section 3.3 by replacing $r_{P;k}^{(rq)}$ and $\xi_{P;k}^{(rq)}$ by $\Delta_k$ and $\xi_k$ respectively.

## 6.3 Specification of the Cost Function

We may decompose the global cost over P phases of follow-ups and data collection as

$$C(\mathrm{P}) = \sum_{p=1}^{P} C_p,$$

with $1 \leq \mathrm{P} \leq \mathrm{P}_{max}$. The P components of the global cost are

$$C_1 = c_1 + C_1^{(\wp)} + C_1^{(fc)} + C_1^{(f)} + C_1^{(dc)} \, ,$$

and for $p = 2, ..., P$
$$C_p = c_p + C_p^{(fc)} + C_p^{(f)} + C_p^{(dc)} \, ,$$

where the extra component associated with the establishment of a first contact is $C_p^{(fc)} = \sum_k 1_k(\wp)(1 - r_{p-1;k}^{(fc)}) \sum_{i=I_{p-1}+1}^{I_p} (1 - r_{i-1;k}^{(fc)}) z_{ki}^{(fc)} \{c_{ki}^{(ffc)} + r_{ki}^{(fc)} c_{ki}^{(rfc)}\}$ for $1 \le p \le P_{max}^{(fc)}$ and $C_p^{(fc)} = 0$ for $P_{max}^{(fc)} < p \le P_{max}$, with $r_{0;k}^{(fc)} = 0$ and $I_0 = 0$, $c_{ki}^{(ffc)}$ is the cost associated with follow-up to get a first contact at time period $i$, and $c_{ki}^{(rfc)}$ is the cost associated with first contact data collection at time period $i$. Then the two remaining components of the global cost are adjusted as follows: $C_p^{(f)} = \sum_k 1_k(\wp) r_{p-1;k}^{(fc)} J_k^{(el)} (1 - r_{p-1;k}^{(rq)}) \sum_{s=1}^{S} J_{s|p;k}^{(f)} c_{s|p;k}^{(f)}$, and $C_p^{(dc)} = \sum_k 1_k(\wp) r_{p-1;k}^{(fc)} J_k^{(el)} (1 - r_{p-1;k}^{(rq)}) \sum_{s=1}^{S} J_{s|p;k}^{(f)} \sum_{m=1}^{M} J_{m|s;k}^{(dc)} r_{p;k}^{(rq)} c_{m|p;k}^{(dc)}$ for $1 \le p \le P_{max}$.

## 6.4 Estimation of the Regression Parameter

The eligibility indicator $J_k^{(el)}$, with $J_k^{(el)} \in \{0,1\}$, is known for a unit that have been contacted during any time of the contactibility period, while it is unknown for a unit that is censored, where the superscript "$el$" stands for "eligible". Consequently the likelihood for the census observed data associated with (4.1) may be decomposed as

$$L_k^{(all)} = f(z_k) L_{R|z_k;k}^{(all)} L_{M|z_k;k}^{(all)} \, ,$$

with $L_{R|z_k;k}^{(all)} = \{\prod_{i=0}^{1} \phi_{i;k}^{(el)} f(t_k^{(fc)} | J_{i;k}^{(el)}) L_k(\lambda | J_{i;k}^{(el)})\}^{J_{i;k}^{(fc)}}\}^{\delta_k^{(fc)}}$ and $L_{M|z_k;k}^{(all)} = \{\sum_{i=0}^{1} [\phi_{i;k}^{(el)} f(t_k^{(fc)} | J_{i;k}^{(el)})]^{J_{i;k}^{(fc)}}\}^{1-\delta_k^{(fc)}}$,

where $J_{i;k}^{(el)} = 1(J_k^{(el)} = i)$, $\phi_{i;k}^{(el)} = \Pr(J_k^{(el)} = i)$, $z_k$ is the stochastic process representing the evolution of follow-up to get a first contact, $L_k(\lambda | J_{i;k}^{(el)}) = L_k(\lambda)$ for illegible units, $L_k(\lambda | J_{i;k}^{(el)}) = 1$ for ineligible units, and $L_k(\lambda)$ is given by (4.1). In likelihood components that incorporate a contactibility piece, appropriate conditioning would be assumed.

# 7. Coefficient of Reliability

To start our discussion on the reliability coefficient, suppose interest is in estimating the variable $\psi$ by an estimator $\theta$ rather than evaluating the reliability of an estimator. In the case of no additional information, the minimum MSE of $\theta = E(\psi)$ is the variance of $\psi$, namely $Var(\psi) = E\{[\psi - E(\psi)]^2\}$; while in the case in which we have additional information, the min MSE of $\theta(\chi) = E(\psi | \chi)$ is also the variance, $Var(\psi | \chi) = E\{[\psi - E(\psi | \chi)]^2 | \chi\}$, but of the conditional density $f(\psi | \chi)$. This conditional variance characterizes the spread of $\psi$ about its conditional expectation $E(\psi | \chi)$ for a given value of $\chi$. If $\chi$ and $\psi$ are independent, then $Var(\psi | \chi) = Var(\psi)$. The remaining relative error (or the missed information) of $\psi$ based on the knowledge of $\chi$ is given by $Var(\psi | \chi) / Var(\psi)$; so that the proportion of knowledge (or the attained information) about $\psi$ obtained after observing $\chi$ constitutes our coefficient of reliability given by

$$K\{\psi; \chi\} = 1 - \frac{Var(\psi | \chi)}{Var(\psi)} \, . \tag{7.1}$$

If $Var(\psi | \chi) = Var(\psi)$ then $K\{\psi; \chi\} = 0$, and if $Var(\psi | \chi) = 0$ then $K\{\psi; \chi\} = 1$.

Suppose that $\psi$ and $\chi$ have a correlation coefficient $\rho_{\psi\chi}$, $\psi$ has a normal distribution with mean $\mu_\psi$ and variance $\sigma_\psi^2$, and $\chi$ has a normal distribution with mean $\mu_\chi$ and variance $\sigma_\chi^2$. Then the conditional distribution has

$$E(\psi | \chi) = \mu_\psi + \rho_{\psi\chi} \sigma_\psi (\chi - \mu_\chi) / \sigma_\chi \, , \tag{7.2}$$

and
$$Var(\psi | \chi) = \sigma_\psi^2 (1 - \rho_{\psi\chi}^2) \, .$$

As shown by Goldberger (1962), the linear estimator given by (7.2) is the best linear unbiased predictor of $\psi$ under the general linear model. In this case under the normality assumption, the coefficient of reliability of $\psi$ based on the knowledge $\chi$ given by (7.1) reduces to

$$K\{\psi; \chi\} = \rho_{\psi\chi}^2 = \left(\frac{Cov(\psi, \chi)}{\sigma_\psi \sigma_\chi}\right)^2 \, . \tag{7.3}$$

Tenenbein (1970) introduced the square of the correlation coefficient given by (7.3) as a measure of reliability between the error-prone and error-free classification variables to measure the strength of the relationship between the true and fallible classifications; i.e., it measures how well the true classification can be predicted from the fallible classification on a given sampling unit. Expression (7.3) gives a convenient way to derive the coefficient of reliability: It is reasonable in practice to replace conditional variance, which depends on the joint distribution, with correlation which can be calculated more easily. That being said, conditional independence is more meaningful and preferable than zero-correlation.

For the survey case, where the estimator $\hat{\Theta}_P(\psi)$ of our gold standard $\Theta_N(\psi)$ is obtained as the solution of (5.2), the coefficient of reliability under the normality assumption is

$$K\{\Theta_N(\psi);\hat{\Theta}_P(\psi)\} = \rho^2_{\Theta_N(\psi)\hat{\Theta}_P(\psi)} = \frac{\{Cov(\Theta_N(\psi),\hat{\Theta}_P(\psi))\}^2}{Var(\Theta_N(\psi))Var(\hat{\Theta}_P(\psi))} \ .$$

Suppose $\Theta_N(\psi) = \sum_k y_k$ and $\hat{\Theta}_P(\psi) = \sum_k F_k y_k$, where $F_k$ is the frame membership indicator for unit $k$. We have under independent observations, $Var(\Theta_N(\psi)) = N\sigma_y^2$, $Var(\hat{\Theta}_P(\psi)) = N_F\sigma_y^2$, $Cov(\Theta_N(\psi),\hat{\Theta}_P(\psi)) = N_F\sigma_y^2$, and $K\{\Theta_N(\psi);\hat{\Theta}_P(\psi)\} = N_F/N$, where $N_F = \sum_k F_k$ and $\sigma_y^2$ denotes the variance of the variable $y$. If $\hat{\Theta}_P(\psi) = \sum_k F_k v_k$ is instead used to estimate $\Theta_N(\psi)$, then $K\{\Theta_N(\psi);\hat{\Theta}_P(\psi)\} = \rho^2_{yv}(N_F/N)$, where $\rho_{yv}$ denotes the correlation coefficient between $y$ and $v$. If $\hat{\Theta}_P(\psi) = \sum_k d_k(\wp)F_k v_k$ under simple random sampling is used to estimate $\Theta_N(\psi)$, then $K\{\Theta_N(\psi);\hat{\Theta}_P(\psi)\} = \rho^2_{yv}(n_F/N)$, where $Var(\hat{\Theta}_P(\psi)) = N_F\sigma_v^2 + N_F\sigma_v^2(N_F - n_F)/n_F$, $\sigma_v^2$ denotes the variance of the variable $v$, and $n_F$ is the sample size. Finally, if $\hat{\Theta}_P(\psi) = \sum_k d_k(\wp)F_k v_k$ under Bernoulli sampling is used to estimate $\Theta_N(\psi)$, then $K\{\Theta_N(\psi);\hat{\Theta}_P(\psi)\} = \rho^2_{yv}(n_F/N)\{1 + (1-\pi_F)(\mu_v/\sigma_v)^2\}^{-1}$, where $\pi_F = n_F/N_F$, $n_F$ is the expected sample size, and $\mu_v$ is the mean of the variable $v$.

Suppose now a third set of information $\aleph$ was available previously to $\chi$, with $E(\aleph) = \mu_\aleph$, $Var(\aleph) = \sigma_\aleph^2$, $\rho(\psi,\aleph) = \rho_{\psi\aleph}$, and $\rho(\chi,\aleph) = \rho_{2\aleph}$. Then, the conditional variance-covariance matrix of $(\psi,\chi)$ given $\aleph$ under the normality assumption is given by

$$\Sigma_{|\aleph} = \begin{bmatrix} \sigma_\psi^2(1-\rho_{\psi\aleph}^2) & \sigma_\psi\sigma_\chi(\rho_{\psi\chi} - \rho_{\psi\aleph}\rho_{\chi\aleph}) \\ \sigma_\psi\sigma_\chi(\rho_{\psi\chi} - \rho_{\psi\aleph}\rho_{\chi\aleph}) & \sigma_\chi^2(1-\rho_{\chi\aleph}^2) \end{bmatrix},$$

and the resulting partial coefficient of reliability after removing the effect of $\aleph$ from each variable is

$$K\{\psi;\chi \mid \aleph\} = \frac{\{\rho_{\psi\chi} - \rho_{\psi\aleph}\rho_{\chi\aleph}\}^2}{(1-\rho_{\psi\aleph}^2)(1-\rho_{\chi\aleph}^2)},$$

provided $|\rho_{\psi\aleph}| < 1$ and $|\rho_{\chi\aleph}| < 1$.

The quality of an estimate is usually characterized by MSE. However, coefficients of reliability for each survey process as well as for a survey as a whole set of processes when supplied with MSE enhance the quality of information on a) the survey results; b) the comparisons between surveys; and c) the contribution of the given survey as addition to prior information.

## Concluding Remarks

We formulated an optimization problem for designing a survey, and we identified steps for its revision in the course of the data collection period. We considered the error-prone and error-free information as a random variable with a joint distribution with some probability function. Then, we updated the joint probability distribution after observing some of realizations of the error-free random process at each phase of data collection, to revise the design specification in the course of the data collection period. The proposed approach makes full use of error-prone information while requiring only few observations from the error-free and expensive random process. Since revision of a design indicates when a design is nearly "optimal", and how the error-free random process varies from the error-prone random process, the revision of the design has an important role to play in survey quality and cost. In the case of insufficient sample size on prior information grounds, our approach can be extended to increase the sample in the course of data collection. Details are omitted to focus on the response mechanism. A reliability coefficient for a survey as a whole set of processes, as well as for a single process, is also discussed. Such a coefficient when supplied with MSE enhances information on a) the survey results; b) the comparisons between surveys; and c) the contribution of the given survey as addition to prior information.

## Acknowledgement

# References

Brémaud, P. 1998. *Markov chains: Gibbs fields, Monte Carlo simulation, and Queues*. Springer.

Demnati, A. 2015. Linearization variance estimators for mixed-mode survey data when response indicators are modeled as discrete-time survival. In *Proceeding of Federal Committee on Statistical Methodology Research Conference.*

Goldberger, A.S. 1962. Best linear unbiased prediction in the generalized linear regression model. *Journal of the American Statistical Association*, **57**, pp. 369-375.

Groves, R.M. and S.G. Heeringa. 2006. Responsive design for household surveys: tools for actively controlling survey errors and costs. *Journal of the Royal Statistical Society*, Series A, **169**, pp. 439-457.

Lagakos, S.W. 1979. General right censoring and its impact on the analysis of survival data. *Biometrics*, **35**, pp. 139–156.

Rosenbaum, P. R. 1987. Model-based direct adjustment. *Journal of the American Statistical Association*, **82**, pp. 387–394.

Tenenbein, A. 1970. A double sampling scheme for estimating from binomial data with misclassifications. *Journal of the American Statistical Association*, **331**, pp. 1350–1361.