

Jean Opsomer  
Westat

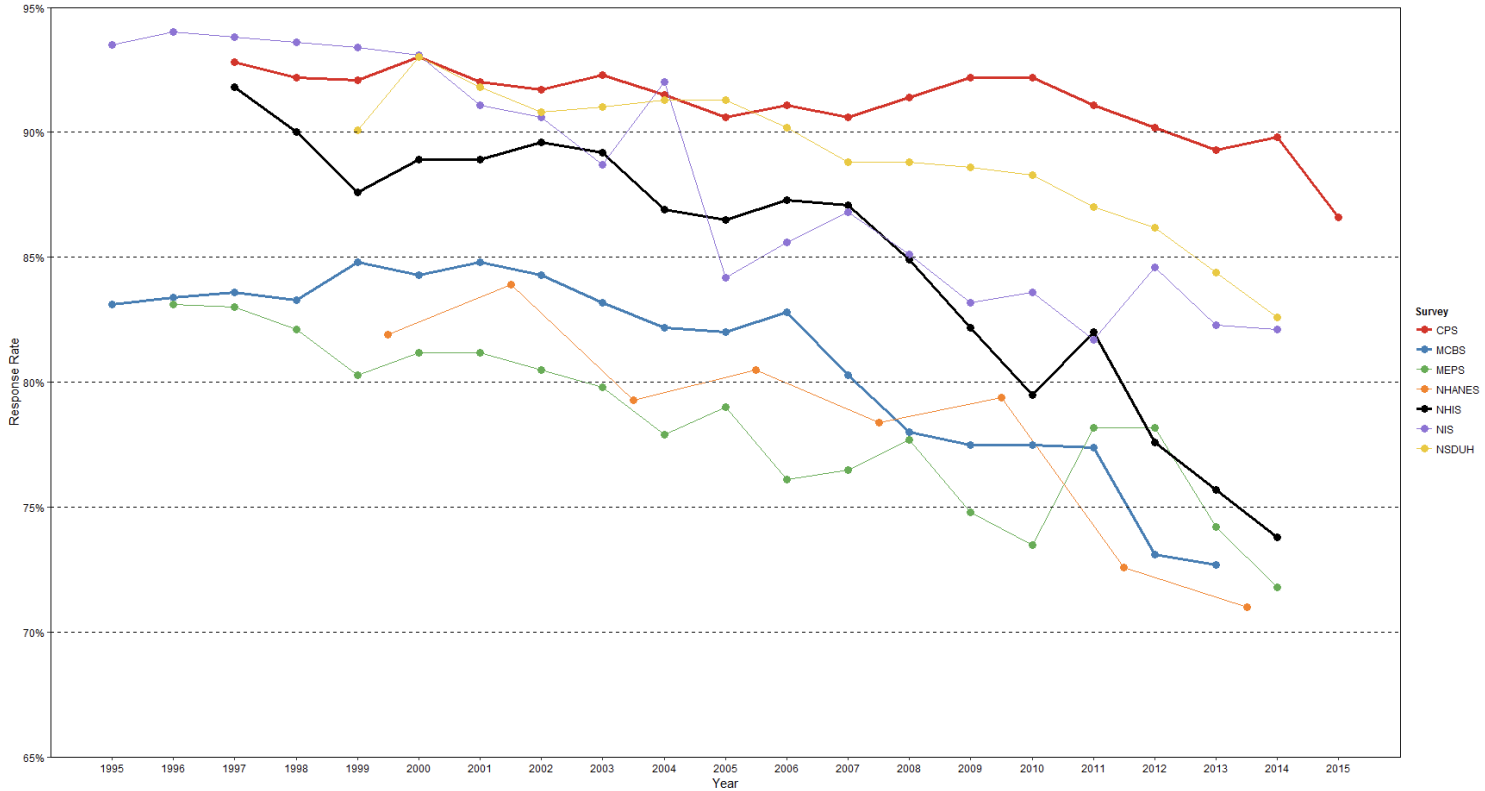
August 20, 2018

1. Design-based survey estimation and inference
2. Constructing estimators
3. Weighting by observed response probabilities
4. Case study: constrained estimation of response probabilities
5. Conclusions

- “Typical” social science survey:
  - large-scale data collection effort conducted on behalf of government agency, using complex multi-stage design
  - output: summary tables and/or weighted datasets
- Key concept: target of inference is specific finite population, e.g. all infants born in US hospitals in 2018, *not* characteristics of a model
- This traditionally leads to *design-based inference*: population treated as fixed but unknown, only randomness comes from sampling design

- Design-based inference is conceptually attractive
  1. assumption-free inference, because design is known
  2. model-free tools available to quantify sampling variability
  3. enables access to high quality datasets for analysis
    - variables available in their original form
    - analyses do not have to be pre-specified
- But:
  1. estimators based on design often inefficient
  2. high nonresponse “breaks” known design assumption

# Growing nonresponse issue in official surveys



(Czajka and Beyler, 2016)

- Nonresponse is seen as important practical and research issue in human population surveys
- Tourangeau and Plewes (2013), *Nonresponse in Social Science Surveys: A Research Agenda*, National Academies Press
- Important on-going research on reducing nonresponse, especially in adaptive multi-mode approaches
- Nevertheless, nonresponse rates are generally expected to continue to increase

- Model-based inference: build model for target variables; once estimated, allows full set of model-based techniques including prediction of population quantities of interest
  1. maximize efficiency, subject only to inherent variability of data (and skill of modeler)
  2. bypass “nuisance” random processes: sampling design, response mechanism
- But:
  1. labor-intensive
  2. sample selection effect can invalidate results

- Thriving area of research within survey statistics, including by many SMURF participants
- Options:
  - ignore
  - apply model-assisted ideas and rely on relationships between variables to correct for nonresponse
  - explicit modeling of response mechanism
  - double-robust approaches
  - etc
- Theory is well understood (still room for improvement!)



- Conceptually, nonresponse is treated as “add-on” to design randomness
- We continue to appeal to classical design properties to justify design-based (weighted) approach to survey inference to users of survey data
- Is this counter-productive?

- “Generalized design-based” (?) relies on combination of design *and models* to account for selection process of obtaining data
  - sampling design, response mechanism, other selection steps (e.g. response-driven sampling)
  - “selection probability” is longer assumption-free, but does not claim to be
- Key aspects:
  1. finite population still treated as fixed target of inference
  2. avoids modeling of survey variables to extent possible

- Modeling selection process: design known, but other components need to be modeled
  - access to paradata, frame data and confidential unit-level data can lead to better models and creation of weights that can be released
  - survey specialists focus on developing and fitting selection models, no need to be subject-matter specialist
  - weights are presented as result of careful modeling, instead of modified design inclusion probabilities (similar to output in other areas of statistics)

- Modeling data selection process instead of data can still result in inefficient inference, since model is “generic” w.r.t. survey variables
- Improving efficiency
  - model-assisted approaches continue to apply
  - weighting by empirical response probabilities
- In both cases, efficiency gains depend on relationship between model variables and survey variables

- Could in principle be handled as a selection problem and modeled as such
- But: “Swiss cheese” nonresponse makes this often not practical
- Approaches:
  - explicit modeling (e.g. multiple imputation, regression imputation)
  - implicit modeling (e.g. hierarchical and/or fractional hot-deck imputation)

- Goal of imputation: create pseudo-data that look like original data
- Pro: allows survey users to continue using data as if obtained under selection-only approach
- Cons:
  - requires modeling of survey variables
  - can increase variability of estimators
- Might be preferable to leave this to subject-matter analysts?

- Finite population:  $U = \{1, 2, \dots, k, \dots, N\}$
- Survey variables

$y_k$  = target variables (unknown outside sample, fixed)

$x_k$  = auxiliary variables (known, fixed)

- Target population parameters: totals, means, proportions, e.g.

$$T_y = \sum_{k \in U} y_k$$

- Sample:  $s \subset U$ , obtained by selection mechanism  $p(s)$

- Sample membership indicator (random)

$$I_k = \begin{cases} 1 & \text{if } k \in s \\ 0 & \text{otherwise} \end{cases}$$

- Selection probabilities

$$p_k = \Pr(k \in s) = \mathbf{E}(I_k)$$

$$p_{kl} = \Pr(k, l \in s) = \mathbf{E}(I_k I_l)$$

- traditional:  $p_k, p_{kl}$  known
- if generalized design-based: well-defined quantities, to be estimated/predicted



- Specifications of selection probability (unit nonresponse case)

$$p_k = \pi_k r(x_k)$$

$$p_{kl} = \pi_{kl} r(x_k) r(x_l)$$

- $\pi_k, \pi_{kl}$  are “pure” design inclusion probabilities
- $r(x) = r(x; \theta)$  is unknown function of auxiliary variable(s)
  - \* usually:  $x$  is multivariate and categorical
  - \* usually:  $r(\cdot)$  is parametric

- Inverse-probability weighting estimator

$$\hat{T}_y = \sum_{k \in s} \hat{w}_k y_k = \sum_{k \in U} \frac{I_k}{\hat{p}_k} y_k = \sum_{k \in U} \frac{I_k}{\pi_k r(x_k; \hat{\theta})} y_k$$

- $\hat{T}_y$  does not behave like “oracle” 2-phase estimator

$$\tilde{T}_y = \sum_{k \in s} w_k y_k = \sum_{k \in U} \frac{I_k}{\pi_k r(x_k; \theta)} y_k$$

- model dependent
- no longer exactly unbiased, even if model is correct
- often includes additional variance terms

- Response homogeneity group (RHG) model
  - common nonresponse adjustment in practice
  - flexible and efficient “all-purpose” approach, as approximation to more complicated underlying model
  - closely related to post-stratification
- Revisit efficiency of RHG (Särndal *et al*, 1992, Ch. 15.6)

- Selection process

- sample  $s$  drawn according to sampling design  $p(s)$
- conditional on  $s$ , units respond independently with unknown probabilities that are equal within groups  $s_g$  ( $s = \cup s_g$ )

- Selection probabilities

$$p_k = \pi_k \theta_g \quad \text{for all } k \in s_g$$
$$p_{kl} = \pi_{kl} \theta_g \theta_{g'} \quad \text{for all } k \in s_g, l \in s_{g'}$$

- Groups  $s_g$  can be sample-dependent and  $\theta_g$  are unknown parameters

---

- Notation

- $R_k = 1$  if unit  $k$  responds, 0 otherwise

- $n_g = \sum_{s_g} 1$ : sample size in  $s_g$

- $m_g = \sum_{s_g} R_k$ : respondent sample size in  $s_g$

- $r_g =$  subset of respondents in  $s_g$

- If  $\theta_g$  known, classical 2-phase estimator

$$\tilde{T}_y = \sum_{g=1}^G \sum_{r_g} \frac{1}{\pi_k \theta_g} y_k$$

- Properties

$$E(\tilde{T}_y) = T_y$$

$$\text{Var}(\tilde{T}_y) = \sum \sum_U \Delta_{kl} \frac{y_k y_l}{\pi_k \pi_l} + E_p \left( \sum_{g=1}^G \frac{1 - \theta_g}{\theta_g} \sum_{s_g} \frac{y_k^2}{\pi_k^2} \right)$$

$$(\Delta_{kl} = \pi_{kl} - \pi_k \pi_l)$$

- RHG estimator

$$\widehat{T}_y = \sum_{g=1}^G \sum_{r_g} \frac{1}{\pi_k \widehat{\theta}_g} y_k = \sum_{g=1}^G \sum_{r_g} \frac{1}{\pi_k \frac{m_g}{n_g}} y_k$$

- Properties

$$\begin{aligned} \mathbf{E}(\widehat{T}_y) &= T_y \\ \text{Var}(\widehat{T}_y) &\approx \sum \sum_U \Delta_{kl} \frac{y_k y_l}{\pi_k \pi_l} \\ &\quad + \mathbf{E}_p \left( \sum_{g=1}^G \frac{1 - \theta_g}{\theta_g} \sum_{s_g} \left( \frac{y_k}{\pi_k} - \frac{\sum_{s_g} y_k / \pi_k}{n_g} \right)^2 \right) \end{aligned}$$

- Compare

$$\text{Var}(\widehat{T}_y) \approx \sum \sum_U \Delta_{kl} \frac{y_k y_l}{\pi_k \pi_l} + \mathbf{E}_p \left( \sum_{g=1}^G \frac{1 - \theta_g}{\theta_g} \sum_{s_g} \left( \frac{y_k}{\pi_k} - \frac{\sum_{s_g} y_k / \pi_k}{n_g} \right)^2 \right)$$

$$\text{Var}(\widetilde{T}_y) = \sum \sum_U \Delta_{kl} \frac{y_k y_l}{\pi_k \pi_l} + \mathbf{E}_p \left( \sum_{g=1}^G \frac{1 - \theta_g}{\theta_g} \sum_{s_g} \frac{y_k^2}{\pi_k^2} \right)$$



- Using observed response probabilities  $\hat{\theta}_g = m_g/n_g$  is equivalent to ratio-type estimator
  - efficiency gains relative to Horvitz-Thompson estimator
- Gains depend on:
  - correctness of response model
  - homogeneity of  $y_k/\pi_k$  within groups
- Gains can offset efficiency losses due to (modest) model departures

- We consider estimator of RHG type

$$\widehat{T}_y = \sum_{g=1}^G \sum_{r_g} \frac{1}{\pi_k \widehat{\theta}_g} y_k = \sum_{g=1}^G \sum_{r_g} \frac{1}{\pi_k \frac{m_g}{n_g}} y_k$$

with  $r_g, s_g$  defined by values of ordinal variable  $x$

- Assume response probability monotone in  $x$ :

$$x_k \leq x_l \Rightarrow r(x_k) \leq r(x_l)$$

and for simplicity, rewrite as

$$\theta_1 \leq \dots \leq \theta_G$$

- Can be set up as design-weighted or unweighted problem; consider unweighted here
- Estimators  $\hat{\theta}_1^c, \dots, \hat{\theta}_G^c$  are solution to

$$\min \sum_{g=1}^G \sum_{s_g} n_g (R_k - \theta_g)^2 \quad \text{subject to } \theta_1 \leq \dots \leq \theta_G$$

with  $R_k = 1$  if unit  $k$  responds, 0 otherwise

- If minimizer satisfies constraint,

$$\hat{\theta}_g^c = \frac{m_g}{n_g} = \hat{\theta}_g$$

- If constraint is binding,

$$\hat{\theta}_g^c = \frac{m_{g_1:g_2}}{n_{g_1:g_2}} = \hat{\theta}_{g_1:g_2}$$

with  $g_1 \leq g \leq g_2$

- In general,

$$\hat{\theta}_g^c = \max_{g_1 \leq g} \min_{g \leq g_2} \frac{m_{g_1:g_2}}{n_{g_1:g_2}}$$

and  $\hat{\theta}_{g_1}^c = \dots = \hat{\theta}_{g_2}^c$  (Brunk, 1955)

→ automatic determination of response homogeneity groups,  
by pooling neighboring groups

- Estimator

$$\widehat{T}_y^c = \sum_{g=1}^G \sum_{r_g} \frac{1}{\pi_k \widehat{\theta}_g^c} y_k = \sum_{g'=1}^{G_s^*} \sum_{r'_g} \frac{1}{\pi_k \widehat{\theta}_{g'}^c} y_k$$

with  $G_s^*$  sample-dependent, determined by pooling

- We study its theoretical properties

- classical design-based asymptotic framework ( $N \rightarrow \infty$ , sequence of designs  $p_N$ , asymptotic normality, etc)
- assuming constrained RHG model holds in population

1. Response probability estimator  $\hat{\theta}_g^c$  is consistent for  $\theta_g$  w.r.t. design and response model
2. RHG estimator  $\hat{T}_y^c$  is consistent for  $T_y$  w.r.t. design and response model
3. Let  $\hat{V}_s^*$  = linearized variance estimate treating the pooled groups  $\{r_{g'}, g' = 1, \dots, G_s^*\}$  as fixed. Then,

$$\frac{\hat{T}_y^c - T_y}{\sqrt{\hat{V}_s^*}} \rightarrow \mathcal{N}(0, 1)$$

- Population:  $N = 10,000$ , 5 equal-sized groups  $U_g$  with
$$y_k \sim \mathcal{N}(1 + g, 1) \quad \text{for } k \in U_g$$

- Sampling design: SRSWOR with  $n = 400$

- Response mechanism: RHG with

$$R_k \sim \text{Ber}(\theta_g) \quad \text{for } k \in U_g$$

and we consider different specifications of  $\theta_1, \dots, \theta_5$

- 10,000 replications

- Estimators of population mean

$\bar{Y}$  = unconstrained RHG estimator

$\bar{Y}^c$  = constrained RHG estimator

$\bar{Y}_{HT}$  = Horvitz-Thompson estimator, true  $\theta_g$

$\bar{Y}_{HA}$  = Hájek (ratio) estimator, true  $\theta_g$

with

$$\bar{Y}_{HT} = \frac{\sum_g \sum_{s_g} y_k / (\pi_k \theta_g)}{N}$$

$$\bar{Y}_{HA} = \frac{\sum_g \sum_{s_g} y_k / (\pi_k \theta_g)}{\sum_g \sum_{s_g} 1 / (\pi_k \theta_g)}$$



- Scenario 1: high response rate, monotone

$$(\theta_1, \dots, \theta_5) = (0.5, 0.6, 0.7, 0.8, 0.9)$$

- Scenario 2: medium response rate, monotone

$$(\theta_1, \dots, \theta_5) = (0.3, 0.4, 0.5, 0.6, 0.7)$$

- Scenario 3: low response rate, monotone

$$(\theta_1, \dots, \theta_5) = (0.2, 0.25, 0.3, 0.35, 0.4)$$

- Scenario 4: equal-probability (monotone)

$$(\theta_1, \dots, \theta_5) = (0.5, 0.5, 0.5, 0.5, 0.5)$$

$$(\theta_1, \dots, \theta_5) = (0.5, 0.6, 0.7, 0.8, 0.9)$$

	Rel. Bias (%)	Scaled MSE
$\bar{Y}$	-0.014	—
$\bar{Y}^c$	-0.21	1.04
$\bar{Y}^{HT}$	-0.034	8.52
$\bar{Y}^{HA}$	-0.009	3.08

$$(\theta_1, \dots, \theta_5) = (0.3, 0.4, 0.5, 0.6, 0.7)$$

	Rel. Bias (%)	Scaled MSE
$\bar{Y}$	0.011	—
$\bar{Y}^c$	-0.247	1.06
$\bar{Y}^{HT}$	-0.034	5.82
$\bar{Y}^{HA}$	0.106	2.90

$$(\theta_1, \dots, \theta_5) = (0.2, 0.25, 0.3, 0.35, 0.4)$$

	Rel. Bias (%)	Scaled MSE
$\bar{Y}$	-0.004	—
$\bar{Y}^c$	-0.586	1.12
$\bar{Y}^{HT}$	-0.011	11.00
$\bar{Y}^{HA}$	0.155	2.97

$$(\theta_1, \dots, \theta_5) = (0.5, 0.5, 0.5, 0.5, 0.5)$$

	Rel. Bias (%)	Scaled MSE
$\bar{Y}$	-0.001	—
$\bar{Y}^c$	-1.772	2.51
$\bar{Y}^{HT}$	0.005	11.07
$\bar{Y}^{HA}$	-0.012	2.97

- Applying RHG estimation at smallest scale possible appears to be most efficient
- Using external knowledge about response probabilities not sufficient to offset this
- Not shown: effects disappear if  $y_k$  *iid* across RHG groups

- Generalized design-based inference
  - corresponds to current “best practice” in survey organizations, but is hidden behind nominal design-based approach
  - claim: should be explicitly recognized and advocated
- RHG (and PS) provide good all-purpose approach for constructing efficient estimators in social surveys, since most variables are categorical

CONTACT: [JeanOpsomer@westat.com](mailto:JeanOpsomer@westat.com)