Answers to criticisms of using design-based methods for inference

> David Binder Statistics Canada (retired) November 7, 2008

Quebec Inter-University Centre for Social Statistics

OUTLINE

- 1. What do we mean by design-based or model-based inference?
- 2. Example (informative sampling)
- 3. The rationale for design-based inference under the model-design-based randomization
- 4. Some problematic cases
- 5. An artificial example
- 6. Modeler's criticisms
- 7. Closing remarks

1. What do we Mean by Design-based or Model-based Inference?

We assume we have observed a sample obtained from a complex survey design, where the units sampled are taken from a finite population.

Most important are:

- 1. The target population(s) vs. the survey population
- 2. The randomization distribution
 - a) Model-based framework
 - b) Design-based framework
 - c) Model-design-based framework

1. What do we mean by designbased or model-based inference?

The target population is the population about which the researcher wishes to make conclusions.

The survey population consists of all the units that are eligible for selection through the frame and survey design being used.

Finite Target Population vs. Survey Population Quantities of interest are finite population quantities



The differences between the survey population and the target population has led to some criticisms of design-based methods; for example, for cutoff sampling.



The Randomization Distribution

• For statistical inference you are interested not only what is <u>actually observed</u> in the selected sample, but also what <u>could have been observed</u> had other samples been selected under the particular randomization framework

• Of primary interest is the distribution of the estimates under hypothetical random repetitions

• The distribution of the estimates depends on whether or not a statistical model is presumed to have generated the values of a finite population (*§*-randomization)

• The distribution of the estimates may or may not be affected by the sample design (*p*-randomization)







Informative Sampling Chambers (2003, 2004)

• Broadly speaking – sampling is informative if distributions of population and observed sample values are different –More precisely, sampling is informative (non-informative) given some information if the two corresponding probabilities are unequal (equal).

• Informative/non-informative status depends on what is being conditioned on (e.g. design variables, sampling process outcome, response/nonresponse status).



Population values of Y are *iid* ~ one parameter exponential

 $f(y;\theta) = \theta \exp(-\theta y)$

Target parameter of interest: $\mu = E[y] = \theta^{-1}$

Sample selected with known inclusion probabilities: $\pi_i = \frac{nY_i}{N\overline{Y}}$

Complete response, no other auxiliary information

Value \overline{Y} is deducible from sample values of Y and their inclusion probabilities

$$\hat{\mu}_{ML} = \overline{Y}; \quad Var_{\xi p}[\hat{\mu}_{ML}] = \frac{\mu^2}{N}$$

Comment:

Full Information MLE becomes difficult in secondary analysis under informative sampling because much of the information desired for this modelling exercise is not available.

Pseudo-Likelihood Approach Design-Based Methodology

$$\hat{\mu}_{PL} = \frac{\sum I_i \pi_i^{-1} Y_i}{\sum I_i \pi_i^{-1}}; \quad Var_{\xi p} \left[\hat{\mu}_{PL}\right] = o(n^{-1})$$

Clearly $\hat{\mu}_{PL}$ is suboptimal by a long shot!. We know the inclusion probabilities π_i , so we know \overline{Y} , which is the ML estimator of μ based on the entire population!

The PL estimator is approximately unbiased in large populations

Sample Likelihood Approach

Approximates distribution of sample values $\{Y_i; I_i = 1\}$ as a function of the population distribution and the sampling weights.

Chambers showed that SL estimator is approximately unbiased under size-biased sampling for the oneparameter exponential model in large populations, and that the variance reduction over the PL estimator could be very substantial!

In general, in terms of efficiency, ML dominates SL, which in turn dominates PL.

Both ML and SL require that the informative sampling mechanism be modelled.

Both SL and PL require values of the actual sample inclusion probabilities.

3. The Rationale for Design-based Inference under the Model-design-based Randomization

We use as a motivating example the standard linear regression model.

$$Y_i = \mathbf{x}'_i \mathbf{0} + \varepsilon_i$$
, where $\varepsilon_i \sim IN(0, \sigma^2)$.

If the complete finite population could be observed, the usual estimating equations for estimating θ are given by

$$\mathbf{U}_{N}(\mathbf{\theta}_{N}) = \sum_{i=1}^{N} \mathbf{x}_{i}(Y_{i} - \mathbf{x}_{i}'\mathbf{\theta}_{N}) = \sum_{i=1}^{N} \mathbf{u}_{i}(\mathbf{\theta}_{N}) = \mathbf{0},$$

where θ_N is the finite population parameter associated with model parameter θ .

It is important to note that

$$\mathbf{E}_{\xi}\mathbf{u}_{i}(\widetilde{\boldsymbol{\theta}}) = \mathbf{E}_{\xi}[\mathbf{x}_{i}(Y_{i} - \mathbf{x}_{i}'\widetilde{\boldsymbol{\theta}})] = \mathbf{x}_{i}\mathbf{x}_{i}'(\boldsymbol{\theta} - \widetilde{\boldsymbol{\theta}}),$$

so that when $\widetilde{\boldsymbol{\theta}} = \boldsymbol{\theta}$, we have $\mathbf{E}_{\xi}\mathbf{u}_{i}(\boldsymbol{\theta}) = \mathbf{0}$.

However, if the true model is

 $Y_i = \mathbf{x'}_i \,\mathbf{\theta} + \mathbf{z'}_i \,\mathbf{\gamma} + \mathcal{E}_i,$

then $\mathbf{E}_{\xi}\mathbf{u}_{i}(\widetilde{\mathbf{\theta}}) = \mathbf{x}_{i}\mathbf{x}_{i}'(\mathbf{\theta} - \widetilde{\mathbf{\theta}}) + \mathbf{x}_{i}\mathbf{z}_{i}'\boldsymbol{\gamma},$

so that when $\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}$, we have $\mathbf{E}_{\xi} \mathbf{u}_{i}(\boldsymbol{\theta}) = \mathbf{x}_{i} \mathbf{z}_{i}^{\prime} \boldsymbol{\gamma}$.

We now define the model-based and the design-based estimators for θ . The model-based maximum likelihood estimator $\hat{\theta}$ is the solution to

$$\hat{\mathbf{U}}(\hat{\boldsymbol{\theta}}) = \sum_{i=1}^{N} I_i \mathbf{x}_i (Y_i - \mathbf{x}'_i \hat{\boldsymbol{\theta}}) = \sum_{i=1}^{N} I_i \mathbf{u}_i (\hat{\boldsymbol{\theta}}) = \mathbf{0}.$$

On the other hand, the design-based the *pseudo*maximum likelihood estimator $\hat{\theta}_n$ is the solution to

$$\hat{\mathbf{U}}_{p}(\hat{\boldsymbol{\theta}}_{p}) = \sum_{i=1}^{N} \pi_{i}^{-1} I_{i} \mathbf{x}_{i} (Y_{i} - \mathbf{x}_{i}^{\prime} \hat{\boldsymbol{\theta}}_{p}) = \sum_{i=1}^{N} \pi_{i}^{-1} I_{i} \mathbf{u}_{i} (\hat{\boldsymbol{\theta}}_{p}) = \mathbf{0}.$$

(For simplicity here, we consider only a single-stage sample design with no weight adjustments.)

This simple regression example can be easily extended to more general estimating equations, covering a wide range of estimators; for example,

- Generalized Linear Models
- Quasi-likelihood
- Generalized Estimating Equations for longitudinal data
- M-estimators

As well, Cox PH models can be adapted to this approach.

To derive variances for large sample situations, we write the "kernel" of the estimating equations as

$$\mathbf{u}_{i}(\widetilde{\boldsymbol{\theta}}) = \mathbf{u}_{i}(\boldsymbol{\theta}) + \left[\left. \frac{\partial \mathbf{u}_{i}(\widetilde{\boldsymbol{\theta}})}{\partial \widetilde{\boldsymbol{\theta}}} \right|_{\widetilde{\boldsymbol{\theta}}=\boldsymbol{\theta}} \right] (\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}) + \mathbf{r}_{i}(\widetilde{\boldsymbol{\theta}})$$

For the linear regression case, this is simply

$$\mathbf{u}_{i}(\widetilde{\boldsymbol{\Theta}}) = \mathbf{x}_{i}(Y_{i} - \mathbf{x}_{i}'\widetilde{\boldsymbol{\Theta}}) = \mathbf{x}_{i}(Y_{i} - \mathbf{x}_{i}'\boldsymbol{\Theta}) - \mathbf{x}_{i}\mathbf{x}_{i}'(\widetilde{\boldsymbol{\Theta}} - \boldsymbol{\Theta}).$$

We now consider the properties of θ_N the finite population parameter associated with model parameter θ . Since

$$\mathbf{0} = \sum_{i=1}^{N} \mathbf{u}_{i}(\mathbf{\theta}_{N}) = \sum_{i=1}^{N} \left[\mathbf{x}_{i}(Y_{i} - \mathbf{x}_{i}'\mathbf{\theta}) - \mathbf{x}_{i}\mathbf{x}_{i}'(\mathbf{\theta}_{N} - \mathbf{\theta}) \right]$$

we have

$$\boldsymbol{\theta}_{N} - \boldsymbol{\theta} = \left[\sum_{i=1}^{N} \mathbf{x}_{i} \mathbf{x}_{i}'\right]^{-1} \sum_{i=1}^{N} \left[\mathbf{x}_{i} (Y_{i} - \mathbf{x}_{i}' \boldsymbol{\theta})\right] = \mathbf{S}_{xx}^{-1} \sum_{i=1}^{N} \mathbf{u}_{i} (\boldsymbol{\theta}),$$

so that under the model,

$$\mathbf{E}_{\xi}[\mathbf{\theta}_{N} - \mathbf{\theta}] = \mathbf{0},$$

and

$$\mathbf{V}_{\xi}[\mathbf{\theta}_{N}-\mathbf{\theta}] = \mathbf{S}_{xx}^{-1} \left[\sum_{i,j=1}^{N} \sigma_{ij} \mathbf{x}_{i} \mathbf{x}_{i}' \right]^{-1} \mathbf{S}_{xx}^{-1} \stackrel{?}{=} \sigma^{2} \mathbf{S}_{xx}^{-1}.$$

These standard results are very well known. We now apply a similar technique to the design-based and model-based estimators. For $\hat{\theta}_{p}$, we have

$$\mathbf{0} = \sum_{i=1}^{N} \pi_i^{-1} I_i \Big[\mathbf{x}_i (Y_i - \mathbf{x}'_i \mathbf{\theta}_N) - \mathbf{x}_i \mathbf{x}'_i (\hat{\mathbf{\theta}}_p - \mathbf{\theta}_N) \Big]$$

so that

$$\hat{\boldsymbol{\theta}}_p - \boldsymbol{\theta}_N = \hat{\mathbf{S}}_{xx}^{-1} \sum_{i=1}^N \pi_i^{-1} I_i [\mathbf{x}_i (Y_i - \mathbf{x}'_i \boldsymbol{\theta}_N)].$$

Therefore, $\mathbf{E}_{p}[\hat{\boldsymbol{\theta}}_{p} - \boldsymbol{\theta}_{N}] \rightarrow \mathbf{0}$ and $\mathbf{V}_{p}[\hat{\boldsymbol{\theta}}_{p} - \boldsymbol{\theta}_{N}] \rightarrow \mathbf{S}_{xx}^{-1} \left\{ \sum_{i,j=1}^{N} \left(\frac{\pi_{ij} - \pi_{i}\pi_{j}}{\pi_{i}\pi_{j}} \right) \mathbf{u}_{i}\mathbf{u}_{j}' \right\} \mathbf{S}_{xx}^{-1}.$

Under the model-design-based randomization framework, we have

$$\mathbf{E}_{\xi p}[\hat{\boldsymbol{\theta}}_{p} - \boldsymbol{\theta}] = \mathbf{E}_{\xi p}[(\hat{\boldsymbol{\theta}}_{p} - \boldsymbol{\theta}_{N}) + (\boldsymbol{\theta}_{N} - \boldsymbol{\theta})] \rightarrow \mathbf{0}$$

and

$$\begin{aligned} \mathbf{V}_{\xi p}[\hat{\boldsymbol{\theta}}_{p}-\boldsymbol{\theta}] &= \mathbf{V}_{\xi p}[(\hat{\boldsymbol{\theta}}_{p}-\boldsymbol{\theta}_{N})+(\boldsymbol{\theta}_{N}-\boldsymbol{\theta})] \\ &= \mathbf{E}_{\xi}\mathbf{V}_{p}[\hat{\boldsymbol{\theta}}_{p}-\boldsymbol{\theta}_{N}]+\mathbf{V}_{\xi}[\boldsymbol{\theta}_{N}-\boldsymbol{\theta}] \\ &\rightarrow \mathbf{S}_{xx}^{-1}\mathbf{E}_{\xi}\left\{\sum_{i,j=1}^{N}\left(\frac{\pi_{ij}-\pi_{i}\pi_{j}}{\pi_{i}\pi_{j}}\right)\mathbf{u}_{i}\mathbf{u}_{j}'\right\}\mathbf{S}_{xx}^{-1}+\mathbf{S}_{xx}^{-1}\sum_{i,j=1}^{N}\left[\sigma_{ij}\mathbf{x}_{i}\mathbf{x}_{j}'\right]\mathbf{S}_{xx}^{-1}.\end{aligned}$$

Main result #1:

If
$$\mathbf{S}_{xx}^{-1} \sum_{i,j=1}^{N} \left[\sigma_{ij} \mathbf{x}_{i} \mathbf{x}_{j}' \right] \mathbf{S}_{xx}^{-1} = O(N^{-1})$$

and $\mathbf{S}_{xx}^{-1} \mathbf{E}_{\xi} \left\{ \sum_{i,j=1}^{N} \left(\frac{\pi_{ij} - \pi_{i} \pi_{j}}{\pi_{i} \pi_{j}} \right) \mathbf{u}_{i} \mathbf{u}_{j}' \right\} \mathbf{S}_{xx}^{-1} = O(n^{-1})$

then $\mathbf{V}_p[\hat{\mathbf{\theta}}_p - \mathbf{\theta}_N]$ is model-design-unbiased for $\mathbf{V}_{\xi p}[\hat{\mathbf{\theta}}_p - \mathbf{\theta}]$ when n/N is negligible.

Main result: #2

Since, if the sampling is non-informative,

$$\mathbf{V}_{\xi p}[\hat{\boldsymbol{\theta}}_{p} - \boldsymbol{\theta}] = \mathbf{V}_{\xi}[\hat{\boldsymbol{\theta}}_{p} - \boldsymbol{\theta}]$$

Then, under similar conditions, for large sample sizes, we have

$$\mathbf{V}_{\xi}[\hat{\boldsymbol{\theta}}_{p}-\boldsymbol{\theta}] = \mathbf{V}_{\xi p}[\hat{\boldsymbol{\theta}}_{p}-\boldsymbol{\theta}] \approx \mathbf{V}_{p}[\hat{\boldsymbol{\theta}}_{p}-\boldsymbol{\theta}_{N}].$$

For the standard regression case, under non-informative sampling,

$$\mathbf{V}_{\xi}[\hat{\mathbf{\theta}}_{p} - \mathbf{\theta}] \rightarrow \sigma^{2} \mathbf{S}_{xx}^{-1} \Big(\mathbf{E}_{\xi} \sum_{i=1}^{N} \pi_{i}^{-1} \mathbf{X}_{i} \mathbf{X}_{i}' \Big) \mathbf{S}_{xx}^{-1}.$$

This is *not* what standard software estimates under WLS regression, even if using normalized weights.

(Note the distinction among probability weights, analytic weights and frequency weights.)

Applying the same techniques to the model-based estimator $\hat{\theta}$ does not give the same conclusions unless the sampling is non-informative

Main result #3

When the sampling is informative, the estimator may not be model-design unbiased, and the model-based variance can also be a biased estimate of $V_{\xi_p}[\hat{\theta} - \theta]$.

	Assumed first-phase model is valid and sampling is	Assumed first-phase model is misspecified or the sampling is
	ignorable	nonignorable
Model- based estimator	 Asymptotically unbiased Efficient Valid variance estimates Valid inferences May be best 	1.May be inconsistent2.Variance estimates may be invalid3.Inferences may be invalid
Design- based estimator	 Asymptotically unbiased May be inefficient Valid variance estimates Valid inferences 	If the mean of the estimating equation is zero under the model: 1. Asymptotically unbiased 2. Valid variance estimates 3. Valid inferences

4. Some Problematic Cases

- Non-negligible sampling fractions
- Rare characteristics
- Small sample sizes
- Models that include random effects
- Population-based case control studies
- Integrating data from more than one survey
- Event history analysis

The following data have been generated randomly using 100 independent and identically distributed trials (Bernoulli trials) taking the value (M or F) The "unknown" parameter is $\theta = \Pr(F) = .51$.

M F F M F F F M M M F F M M M M M M F MF F F MF F M F F F F F M F M M F MM M MM F MM MM MM F F M M F F M F MF MMFMFFMFMMMM Μ M F F M MMMFFFMFMMFFFMF F F M MM

We have 46F's and 54M's, so $\theta_N = .46$.

We now consider a finite consisting of these 100 observations as a starting point. The observations have been arranged into households using a non-random process. The results are as follows:

(M M M M F),(F M),(M M F),(F F F F M), (F M),(F M),(M F),(M F),(M F),(M F), (F M),(F F F F F F F M),(F M),(M F),(F M), (M M M M F),(M M M M M M F),(F M),(M F), (F M),(F M),(F M),(M F),(M F),(F M),(F M), (M M M M M F),(F M),(M M M F),(F F M), (F M),(M F),(F F M),(F F F M),(M M)

We now take a cluster sample by selecting *m* households from the 35 households at random, with replacement, using selection probabilities proportional to household size, and then enumerating all persons in the selected households.

It turns out that $E_p[\hat{\theta}] = .445$, so that the design bias is only -.015.

However, for large samples, we have

$$E_{p}\hat{V}_{\xi}[\hat{\theta}_{p}] = .087/m$$
$$E_{p}\hat{V}_{p}[\hat{\theta}_{p}] = .049/m$$
$$E_{p}\hat{V}_{\xi}[\hat{\theta}] = .068/m$$
$$E_{p}\hat{V}_{\xi}[\hat{\theta}] = .114/m$$

Obviously, the sampling plan is informative.

We see from this simple example that the impact of informative sampling on model-based methods can be substantial.

Hoem, Jan M. (1989)

"... sample-based analyses of individual longitudinal behavior can normally do well without sampling weights. Instead of worrying about such weights, it pays to concentrate on the modelling of behavior and on drawing inference about features of the model. One should not feel confined to finite population totals and means, finite population regression coefficients, and other finite population statistics. Also, some of the claims about the good properties of conventional weighting seem exaggerated."

Hoem, Jan M. (1989)

• Use of weights is not appropriate when modeling probabilistic models of human behaviour because we are not estimating a finite population quantity.

• Regard sampling mechanism as part of the total model of the "random experiment" that produces the survey data, with normal consequences for the statistical analysis.

Hoem, Jan M. (1989) – cont'd

• Informativeness is not relevant for the issue of robustness against behavioural model misspecification

• The finite population quantity associated with the model parameter may not always be a useful descriptive statistic.

Fienberg (1989)

• Sampling weights, as they are usually constructed, are at best irrelevant to a likelihood-based approach to statistical inference.

• However, weights may be appropriate for outcome-based sampling.

Breckling et al. (1994) and others have mentioned the problem of using design-based methods for estimating finite population quantities with cut-off sampling.

- this is a clear-cut example of where the survey population and the target populations differ.

Sugden (1998)

There is a danger of using purely asymptotic arguments to justify robustness. The asymptotic scheme requires an infinite sequence of samples and populations, usually replicating the given population and sampling scheme somewhat arbitrarily and incapable of verification. Instead of limit theorems, I would like to see approximations to the bias for use with moderate sample sizes.

Little (2004)

The design-based approach to survey inference has a number of strengths that make it popular with practitioners. It automatically takes into account features of the survey design and provides reliable inferences in large samples, without the need for strong modeling assumptions.

Little (2004) – cont'd

On the other hand, it is essentially asymptotic, and hence yields limited guidance for small-sample adjustments. Unlike models, which lead to efficient inferences based on likelihood or Bayesian principles, the design-based approach is not prescriptive for the choice of estimator. It lacks a theory for optimal estimation, and the estimates that it yields are potentially inefficient.

Little (2004) - cont'd

• Modelling provides a unified approach to survey inference, aligned with mainline statistics approaches in other application areas, such as econometrics.

- In large samples and with uninformative prior distributions, results can parallel those from design-based inference.
- The Bayesian approach is well equipped to handle complex design features.

• The Bayesian approach may yield better inferences for small-sample problems where exact frequentist solutions are not available, by propagating error in estimating parameters.

Little (2004) - cont'd

•The Bayesian approach allows prior information to be incorporated, when appropriate.

•The Bayesian approach avoids the ambiguities in the choice of reference distribution and has useful features of coherency not shared by frequentist approaches, such as satisfying the likelihood principle.

• Likelihood-based approaches like Bayes or maximum likelihood have the property of large-sample efficiency, and hence match or outperform design-based inferences if the model is correctly specified.

The Question of Nonresponse Modelling

To this point, we have not considered the complexities which arise when making adjustments for non-response. Many of the criticisms of using design-based weights in making survey inferences are related to modeling for non-response.

As Kim and Kim (2007) note that the issue of whether or not to use the weights when modeling non-response propensities is not clear-cut, and that when one considers the estimates with the estimated non-response adjustment incorporated neither using the weights nor ignoring may be optimal. This is clearly an area of future research.

As we mentioned earlier, when successful at using models to incorporated the reasons for informative sampling, Maximum Likelihood or Sample Likelihood methods will usually outperform the Pseudo-likelihood approach using designbased methods.

When sample sizes (or number of psu's) are small, a modelbased or even a Bayesian approach can be preferred to a design-based approach even if the model is not quite correct.

Models that include random effects are definitely problematic, even for large total sample sizes.

The use of the original weights when the weights are very variable may not be best, even when interested in finite population quantities. Some weight modifications may be suitable; for example in population-based case control studies, or in integrating data from more than one survey. Users may also find that a few observations are too influential due to the large weights, and it would make sense to make adjustments.

Weighted and unweighted point estimates may or may not be similar. If they are not, think about modifying the model to incorporate the fact that the sampling is informative, so that the model better explains the sampling distribution of the data.

Even if the point estimates are similar, consider modifying the model to account for reasons why the sampling could be informative.

Care must be exercised when modifying a model to include design variables. One needs to ensure that the model is being interpreted as was originally intended.

Standard errors based on a design-based approach may tend to be more robust in cases where the informativeness of the sample has not been fully incorporated into the model.

The design-based approach will still give correct inferences for the parameters of the model used to generate the finite population, when the assumed model for the finite population is true, even if the sampling is informative.

What do we mean by robustness? The design-based approach guards against misspecification of the model error structure. It is still important to get the best functional form for the model expectation - such as linearity and no missing variables.

FINAL WORDS

'all models are wrong, some are useful'- George Box

DISCUSSIONS WILL CONTINUE FOR A LONG TIME!

Contact information: David Binder <u>dbinder49@hotmail.com</u>