**Routledge**
Taylor & Francis Group

Check for updates

# Abstract: An Evaluation of Planned Missing Data Designs in Large Surveys

Dan Su

University of Wisconsin-Madison

Planned missing data designs in large surveys can efficiently reduce respondents' burden and lower the cost associated with data collection, without cutting down on the questionnaire items. If the missing data are not appropriately planned, descriptive and potential causal parameter estimates will be biased. This paper implemented two studies which use simulated data based on the PISA 2006 data. The two simulation studies investigate how the bias in marginal means, correlations and regression coefficients depend on the chosen planned missing data designs and related characteristics. The first study shows that, for a fixed sample size, the extent of bias depends on the three major properties of design and data: overlap percentages (i.e., the portion of cases where two variables are jointly observed), missing percentages, and distributions of variables. The second study applies the properties of design to illustrate that an optimal incomplete block design that ensures overlap can be a better choice than a multiple-form design given the same amount of missingness. In large surveys, the available planned missing data designs have not yet been systematically investigated with respect to the missing percentage, overlap percentage, distribution of variables, and sample sizes. The findings will guide researchers in choosing a design with optimally arranged items in forms. With appropriately chosen planned missing data designs and multiple imputations, the number of survey questions can be increased without harming the validity of parameter estimates.

The results of the first simulation show that the estimates of means are unbiased for large survey data (i.e., sample sizes exceeding 1000 cases) even when overlap is zero and the missing percentage is high (e.g., 80%). However, the recovery of correlations and regression coefficients requires sufficient overlap. The bias in correlations is negligibly small when there is 20% or more overlap for continuous data. Regarding regression coefficients, the bias is negligibly small when overlap exceeds 20% for multivariate normal data. For skewed and categorical data, the coefficients are estimated with larger bias and less reliability. The second simulation study compares a two-form design, a three-form design, and two optimal block designs with 50 and 33% missingness. The results show that all designs recover the means of variables without bias. The biases in correlations are negligibly small for all designs except for the two-form design which has no overlap across forms. With the same amount of missingness, the optimal block design with 50% missingness largely reduces the bias in regression coefficients due to its 20% overlap compared to the two-form design. Furthermore, the three-form design and the optimal block design with 30% missingness produce negligibly small bias and more reliable estimates compared to the other two designs.

## Article information

**CONTACT** Dan Su ✉ dsu4@wisc.edu 🖥 Graduate School, University of Wisconsin-Madison, Madison, WI 53706.