



Research Report

Department of Statistics

No. 2003:12

**Multiple Imputation with Double
Samples: A Simulation Study**

Boris Lorenc

Multiple Imputation with Double Samples: A Simulation Study

Boris Lorenc

Department of Statistics, Stockholm University

SE-106 91 Stockholm, Sweden

E-mail: boris.lorenc@stat.su.se

Abstract

In the double samples procedure for web surveys, introduced by Terhanian, study variables are collected only from the web panel. The problem of obtaining valid population estimates for these variables, treated by weighting on the propensity score in Terhanian's work, may be recast in different terms: of imputing values of the study variables missing in the sample from the general population. Specifically, multiple imputation is suggested, as it has the advantage of easily obtainable measures of uncertainty regarding the point estimates—they come as a byproduct of the multiple imputation procedure.

In this study, some multiple imputation techniques were compared in the double samples setting, using artificial populations with known distributions. Amongst the applied techniques was the propensity score technique for multiple imputation, enabling implicit comparison with the previously obtained simulation results on the propensity score weighting. The results indicated that, in the simulated setting, all the imputation techniques except the propensity score gave nearly unbiased estimates. The accompanying estimators of variance performed as intended, but were somewhat conservative in some circumstances. Also studied was the performance of the estimators in the situations when the assumptions pertaining to the propensity score weighting did not hold.

Introduction

Building on the work of Rosenbaum and Rubin, who introduced the propensity score as a means for estimating causal effects in observational studies (Rosenbaum and Rubin, 1983, 1984), Terhanian used the propensity score weighting in a situation when double samples are taken from one and the same population (e.g. Terhanian, Marcus, Bremer, and Smith, 2001). Besides surveying a web panel, Terhanian collected auxiliary information¹ from

¹Age, sex, etc., but even some attitudinal variables. In general, what in the survey literature are known as *auxiliary* (or sometimes *background*) variables are in regression analysis referred to as *independent* variables, in the biomedical research as *covariates*, and in the econometric literature as *conditioning* variables.

an explicitly drawn sample from the whole population. He used the auxiliary data from both samples to estimate the propensity scores for being in the web panel, and weighted the values of the study variables observed only in the panel by stratifying them on this estimated propensity score. It is worth noting that as the weights do not depend on the study variables, they may be reused on new surveys of the panel as long as it is believed that the panel or the population have not changed sufficiently enough to necessitate new generation of the weights.

In theory, the propensity score weighting produces unbiased estimates. In practice, it typically removes a significant proportion of bias of the estimates, provided that the assumptions pertaining to the technique hold. For instance, about 90% of the original bias was removed in an analytically explored situation (Lorenz, 2003a). The percentage reduction in bias was there a function of only the number of classes into which the distribution of the propensity scores was stratified. In a simulation study (Lorenz, 2003b), with the assumptions holding, the bias reduction was again about 90% across a range of factors and levels. In practical applications on real data, the technique was also reported to function well (e.g. Terhanian et al., 2001).

But, the fact that the propensity score weighting produces point estimates with a residual bias might be unsettling. Furthermore, the suggested estimator of variance of the adjusted point estimates (Rosenbaum and Rubin, 1984) is approximate, not taking into account the uncertainty stemming from estimation of the propensity scores by use of a model instead of exactly knowing the true propensity scores. In the aforementioned simulation study, whether the approximate variance estimates underestimated the true variance of the adjusted point estimator or not depended on the covariance structure of the variables involved: in half of the studied covariance structures, the estimated variance was on mark or nearly so. But, even when the variance was estimated approximately correctly, confidence intervals built around these biased point estimates gave confidence levels constantly below the targeted ones: by a few percent on average for that half of the covariance structures where the variance estimates were correct, and by more than 50% on average for the other half.

With these concerns, one may be led to consider alternatives. The structure of data for the double samples procedure is represented in Table 1. Denotation of the variables is similar to the one commonly met in the survey literature: \mathbf{X} denominates multivariate auxiliary information about the participants (sex, age, etc., obtained from current survey or external sources like register data), \mathbf{Y} denominates the study variables particular to just the current survey, and Z is an indicator variable. Specific to the double samples procedure as applied by Terhanian, \mathbf{X} includes behavioural and attitudinal variables. Z indicates observations collected using the *restricted* sample, drawn from a subset of the population (e.g. a sample among all the web users ready to participate in web panels).² As conceived by Terhanian, the procedure does not require collection of \mathbf{Y} variables from the *unrestricted* sample, drawn from the whole population³.

²It may be noted that there is a positive probability of a unit being present in both samples simultaneously. In such a case, it appears twice in the data: once among the n units without the \mathbf{y} values and with $z = 0$, and once among the k units with a complete set of observations and with $z = 1$.

³The unrestricted sample s is in practice often a random-digit-dialing sample, for which the data are collected through a telephone interview. In general, collection method for s is more expensive than that for r , the web panel, which justifies the whole setup.

Table 1: Data matrix for the double samples procedure: an unrestricted sample s of n units drawn with known inclusion probabilities π from the whole population but with data missing on \mathbf{Y} , and a restricted sample r of k units drawn with unknown inclusion probabilities from a subset of the population but having complete information. $Z=1$ indicates that a vector of observations $(-, \mathbf{x}, \mathbf{y})$ is collected using the restricted sample r and $Z=0$ indicates that a vector $(\pi, \mathbf{x}, -)$ is collected using the unrestricted sample s .

Observations	Variables:	π	X_1	\cdots	X_p	Y_1	\cdots	Y_q	Z
1		π_1	$X_{1,1}$	\cdots	$X_{p,1}$	—	\cdots	—	0
2		π_2	$X_{1,2}$	\cdots	$X_{p,2}$	—	\cdots	—	0
\vdots		\vdots	\vdots	\ddots	\vdots	\vdots	\ddots	\vdots	\vdots
n		π_n	$X_{1,n}$	\cdots	$X_{p,n}$	—	\cdots	—	0
$n+1$		—	$X_{1,n+1}$	\cdots	$X_{p,n+1}$	$Y_{1,n+1}$	\cdots	$Y_{q,n+1}$	1
$n+2$		—	$X_{1,n+2}$	\cdots	$X_{p,n+2}$	$Y_{1,n+2}$	\cdots	$Y_{q,n+2}$	1
\vdots		\vdots	\vdots	\ddots	\vdots	\vdots	\ddots	\vdots	\vdots
$n+k$		—	$X_{1,n+k}$	\cdots	$X_{p,n+k}$	$Y_{1,n+k}$	\cdots	$Y_{q,n+k}$	1

Another way of looking at the data in Table 1 is of a data matrix with missing values. Missing are the values for the study variables, \mathbf{Y} , for the units in the unrestricted sample. A way of obtaining an unbiased estimate of the parameter of interest for the general population would be through imputing the missing \mathbf{Y} values. Using a mildly simplifying assumption that the unrestricted sample is a simple random sample from the population, an unbiased imputation of the missing values would yield an unbiased estimate of the population mean for the study variable, the parameter we are seeking to estimate.

So, there exist reasons to consider imputation as a serious alternative to the propensity score weighting in the situations requiring double samples:

- propensity score weighting reduces bias but does not remove it completely: in practice, due to a limited amount of data and a limited number of strata, about 90% of the bias is removed in the most favourable conditions; many of the multiple imputation techniques are asymptotically unbiased under assumptions no stronger than those for the propensity score weighting,
- to generate weights, the propensity score technique uses only information in \mathbf{X} , so weights once generated are the same irrespective of which \mathbf{Y} variable they are to be applied to; in imputation, what missing values are imputed is dependent even on the existing \mathbf{Y} values, in addition to the \mathbf{X} values, indicating better use of the available information,
- the propensity score weighting gives an estimate of the variability of the point estimate only conditional on the model chosen to estimate the propensity scores, while uncertainty concerning this choice is left out; in contrast, multiple imputation seems to be of particular use for this specific problem due to its possibility to

provide information about uncertainty regarding the imputed values and thereby uncertainty regarding the estimate of the parameter of interest,

- there are no off-the-shelf statistical programs that do the propensity score weighting; multiple imputation is by now a readily available tool for treatment of missing data in several statistical packages and dedicated programs and routines⁴.

In the present study, a simulation is used to demonstrate that indeed multiple imputation is a useful alternative to the propensity score weighting in the double samples situation. The study compares a number of procedures for multiple imputation, among them the propensity score as a technique for multiple imputation (SAS Institute Inc., 2001). Using the same population model as in a simulation study of the performance of the propensity score weighting (Lorenç, 2003b), the present investigation includes an implicit comparison of the two techniques in the double samples setup. Section 1 gives the population model and some theoretical background, Section 2 gives methodological details of the simulation study, while Section 3 presents the results. In Section 4 some concluding remarks are given.

1 Population model and background theory

1.1 Population model

A multivariate normal distribution served as a model for the population, following the use of the same kind of model in some related analytical investigations (e.g. Cochran, 1968; Cochran and Rubin, 1973). In at least a number of practical situations it may be reasonable to consider the normal model as applicable.

The following multivariate normal model was used:

$$(X_1, X_2, Y, V) \sim N(\mathbf{0}, \Sigma),$$

where

$$\Sigma = \begin{bmatrix} 1 & \rho_{12} & \rho_{13} & \rho_{.4} \\ \rho_{12} & 1 & \rho_{23} & \rho_{.4} \\ \rho_{13} & \rho_{23} & 1 & \rho_{.4} \\ \rho_{.4} & \rho_{.4} & \rho_{.4} & 1 \end{bmatrix}. \quad (1)$$

This model defined a population, while inclusion into the subset was defined through either $Z = I_{V < X_2}$ or $Z = I_{\max(V, 0) < X_2}$. The variables in the model were given the following meanings, not uncommon in the survey literature:

⁴There are even other alternatives for the double samples setting worth considering. For instance, GREG estimation, calibration and the use of kernel estimation techniques. These techniques, however, will not be discussed here.

X_1 , an auxiliary variable,
 X_2 , another auxiliary variable, also involved in defining
the subset—“the participation variable”,
 Y , the study variable,
 V , another variable involved in defining the subset.

Two samples drawn using simple random sampling were conceived, an unrestricted sample from the complete population, denoted by s and of size n , and a restricted sample from the subset given by $Z = 1$, denoted by r and of size k .

Each of the variables in the model (1) is by itself a standard normal, so expected values of their population means are zero. Likewise are expected values of the variables’ means in the unrestricted sample s zero, because s is a simple random sample from the population.

1.2 Estimation goal

It is desired to estimate correctly, in the double samples setup (c.f. Table 1), the mean, \bar{Y} , of the study variable Y in the population.

Using the unadjusted mean in the restricted sample r to estimate the mean of Y in the whole population yields biased estimates. When $\rho_{.4} = 0$, conforming to the pair of assumptions known as “strongly ignorable treatment assignment” (given under the heading “The propensity score approach”, below)⁵, the means of both the auxiliary variables and the study variable are biased with respect to their corresponding means in the population: for the auxiliary variables, $E_r(\bar{X}_2) = \pi^{-\frac{1}{2}} \approx 0.564$, $E_r(\bar{X}_1) = \rho_{12} \times \pi^{-\frac{1}{2}}$, and for the study variable, $E_r(\bar{Y}) = \rho_{23} \times \pi^{-\frac{1}{2}}$. These estimates differ from zero whenever the correlation coefficients between the participation variable X_2 and the corresponding variables are not zero, so in the present study they are biased.

1.3 Theory for the approaches

The two mentioned techniques, the propensity score weighting and multiple imputation, may be used for correcting the aforementioned bias. Their theoretical background is now presented in more detail.

1.3.1 The propensity score approach

Terhanian (Terhanian et al., 2001) suggested using the propensity score weighting to reduce the bias of the estimates obtained using only web panels (i.e., in the present study, using only the r sample).

The propensity score (Rosenbaum and Rubin, 1983), denoted by $e(\mathbf{x})$, is a function of the auxiliary variables. It is defined as the conditional probability that a unit with the properties \mathbf{x} is included in the restricted sample r , $e(\mathbf{x}) = \Pr(Z = 1 | \mathbf{X} = \mathbf{x})$.

⁵Without this assumption, the mean may be either more biased or less biased, depending on the sign and magnitude of $\rho_{.4}$; the effect of this and some other factors was explored in (Lorenz, 2003b).

Let *strongly ignorable treatment assignment* (SITA) denote fulfillment of the following two conditions: (i) independence of the study variable and the group assignment conditional on the auxiliary information, $(Y \perp Z) | \mathbf{X}$, and (ii) a positive probability at every level of the propensity score for every unit in the population to be assigned to any of the groups, $0 < e(\mathbf{x}) \leq 1$. Then, when the SITA assumptions hold, weighting of the observed Y values in the restricted sample r by conditioning on the propensity score yields, in theory, unbiased estimates of the population means of the study variables (Lorenc, 2003a).

More practically, the technique consists of the following steps:

1. collecting complete data—the auxiliary variables and the study variables—from the restricted sample (e.g., a web panel) and collecting the auxiliary variables from the unrestricted sample (e.g., a random sample drawn from the target population),
2. given the whole set of auxiliary information (from the unrestricted and the restricted samples) but not the sample membership indicator, estimating for each unit the probability of being a panel member (this magnitude is known as the estimated propensity score); a common way of estimating this probability is by building a logistic regression model,
3. estimating the distribution of the propensity score in the target population by considering the distribution of the estimated propensity score in the unrestricted sample only; in particular, identifying cutoff points for stratification: usually equidistant cutoff points are chosen and 5 intervals are used, in which case the cutoff points would be the 20th, 40th, 60th, and 80th percentile of the estimated propensity score distribution in the population,
4. classifying the units in the restricted sample (panel) into appropriate strata based on their individual estimated propensity score values,
5. for each stratum, building a mean of the study variable values of the panelists in that stratum; then, weighting the strata means appropriately together to produce the final, adjusted estimate for that study variable; in the case of equidistant intervals the weighting amounts to calculating the arithmetic mean of the strata means.

Justification for the procedure and its details were given in (Lorenc, 2003a).

The identified shortcomings of the propensity score weighting, discussed above, include: (a) it might leave a residual bias, presumably in practical applications of the order of 10% of the original bias or more, provided the assumptions pertaining to the technique hold, and (b) it is difficult to produce confidence intervals for the point estimates that would have a desired, predetermined confidence level.

1.3.2 The multiple imputation approach

Values of the study variables for the unrestricted sample s —the questions that in fact by design were not posed to the respondents in that sample—may be even viewed as missing values. Then, if they could be perfectly imputed (replaced by correct values), the situation would have been the standard one from the usual sampling theory: with simple

random sampling, the s sample's mean would be the estimate of the population mean, and the only uncertainty—that stemming from taking a sample instead of performing a census—would be estimated based on the variance of the Y values in the s sample. Unfortunately, the missing values cannot be imputed exactly.

Prior to introduction of the multiple imputation procedure, imputations resulted in a single value being imputed, not reflecting the degree of uncertainty regarding appropriateness of the imputed value. In the multiple imputation approach (Rubin, 1987), imputation for a single missing value is performed several times using a model that includes stochastic elements, creating each time a quasi-complete data set. While none of the imputed values by itself purports to represent what just i unit's y_i value would be, taken together the imputed values both represent an estimate of the missing value and reflect the uncertainty regarding this estimate. Thus, the values imputed to a sample as a whole allow for building a sample point estimate and also an estimate of its variance. In a second step, a generalization is made from this sample to the population as a whole, which is straightforward with for instance simple random sampling.

Let the imputation of missing Y values in a data set with $i = 1, 2, \dots, n$ units be performed $j = 1, 2, \dots, m$ times. If the i unit's Y value, y_i , was observed, then it is left unchanged; if it is missing, a value, y_{ij}^* , according to a model is set in place of the missing value, where y_{ij}^* may differ from the next imputation, $y_{i(j+1)}^*$. Each round of imputations yields a quasi-complete data set, for which two statistics may be calculated: a point estimate for that data set, $\bar{y}_j^* = \frac{1}{n} \sum_{i=1}^n y_{ij}^*$, and a corresponding variance estimate, $\hat{V}_j^*(\bar{y}_j^*) = \frac{1}{n(n-1)} \sum_{i=1}^n (y_{ij}^* - \bar{y}_j^*)^2$.

After all the m imputation rounds have been performed, assuming that the n units were drawn by a simple random sampling procedure, the population point estimator can be calculated,

$$\hat{Y}_{MI} = \bar{y}^* = \frac{\sum_{j=1}^m \bar{y}_j^*}{m}, \quad (2)$$

and the estimator of variance for this estimator,

$$\hat{V}(\hat{Y}_{MI}) = \frac{\sum_{j=1}^m \hat{V}_j^*(\bar{y}_j^*)}{m} + \left(1 + \frac{1}{m}\right) B, \quad (3)$$

where B is the between-sets variance of the point estimator,

$$B = \frac{\sum_{j=1}^m (\bar{y}_j^* - \bar{y}^*)^2}{m}$$

In words, the variance comprises of two components, the first being the mean of the individual data sets' point estimate variances, and the other a slightly inflated between-sets variance of the point estimates. Theoretical underpinning for the technique was given by Rubin (1987).

2 Method

The study was performed as an experiment with a number of factors, with the primary aim to investigate the bias reducing performance of multiple imputation in a double samples

Table 2: Denotations for the covariance structures used in the study.

Covariance structure	ρ_{12}	ρ_{13}	ρ_{23}
1	low	low	low
2	low	low	high
3	low	high	low
4	high	low	low
5	low	high	high
6	high	low	high
7	high	high	low
8	high	high	high

setup, varying the relevant conditions. Several multiple imputation procedures were included, amongst them the propensity score, enabling comparison of the two alternative adjustment approaches.

In order to verify that the results regarding the propensity score imputation would agree with the ones from an earlier simulation regarding the propensity score weighting (Lorenc, 2003b), some of the factors pertaining to that study were applied here. These concerned in part the inner workings of the propensity scores technique (e.g. sample sizes, ratio of the sample sizes), and in part breakage of the assumptions supporting the propensity scores technique.

2.1 Factors

A brief motivation for inclusion of the factors follows.

2.1.1 Covariance structure

A bias in estimators arises in general due to correlation between the study variables and the variables causing the unwelcome event (for instance, nonresponse). Indeed, with non-response independent with the study variable, the data observed on the respondents are perfectly valid for a point estimate. Correspondingly, all the adjustment methods attempt to use the correlation of the relevant variables to correct for the bias. Thus, correlation of the variables relevant for the situation at hand is a factor of great importance for performance of the adjustment techniques. Furthermore, in real situations, the covariance structure is not known (had it been known, no survey would have been needed!), but is assumed instead to be such and such. The effect of eventual misspecifications may be of interest.

The covariance matrix in (1), setting V aside for the moment, produces 8 different models when each of ρ_{12} , ρ_{13} , and ρ_{23} is held on one of the two positive levels, “high” and “low”. Varying the covariance structure in this way gave the opportunity to investigate the efficiency of the multiple imputation techniques under the “high” and “low” levels of correlation between each of the covariates and the response (Table 2).

As this reduced, 3×3 , covariance matrix needs to be positive definite, a pair of the lowest and the highest values of the three ρ 's that in all the 8 combinations produced

a positive definite matrix was by trial and error determined to be $\rho_{\cdot,low} = .22$ and $\rho_{\cdot,high} = .78$.

2.1.2 Sample size

Sample size of about 1100 is by tradition used in surveys if an estimate of a population proportion is to be given with a 3 percent bound of error with a 95% confidence. With increased uncertainties due to imputation, it was hypothesized that samples of that size perhaps would not suffice for achieving sufficiently precise results. So, this factor had two levels, $n_{low} = 1000$ and $n_{high} = 5000$, where the latter level, by comparison, would result in a 1.4 percent bound of error.

2.1.3 Ratio of the sample sizes

It is reasonable to assume that the two samples, the unrestricted and the restricted one, will not be of the same size in practice, for instance that one of them was originally drawn and the values recorded for a different purpose. While, in general, more data ought to improve precision of the estimates, in some circumstances ratio of the sample sizes might be of influence. This factor is thus included in the experiment, using three levels (unrestricted sample's size is in the denominator): 1/2, 2/2, and 3/2.

2.1.4 Multiple imputation methods

Five methods of multiple imputation were used, that is, all those existing in an experimental version of PROC MI included in SAS 8.2 (the latest version of SAS available at the time of performing the present study): expectation maximization (EM), Markov Chain Monte Carlo (MCMC) with initial mean and covariance estimates obtained by EM, MCMC with initial mean and covariance estimates obtained by bootstrapping, regression, and propensity score. (For the details concerning application of these techniques to multiple imputation c.f. SAS Institute Inc., 2001).

2.1.5 SITA violation #1

When $\rho_{.4} = 0$, the SITA assumption of the conditional independence $(Y \perp Z) | \mathbf{X}$ holds. In other cases, it is violated. The impact of setting $\rho_{.4}$ to a particular value different from zero on the performance of the multiple imputation adjustment was explored in the experiment. In order not to inflate the number of factors, V 's correlation with the other variables, $\rho_{.4}$, was the same across the variables within a condition. The value of $\rho_{.4}$ was set to $-.175$ in order to ensure comparability with the corresponding level in (Lorenc, 2003b).

For this factor, consisting of two levels, the reference "SITA violation #1" is used in what follows.

2.1.6 Inclusion of all relevant variables

Efficiency of the adjustment methods is contingent on inclusion of all relevant information among the observed data. In the propensity score approach for instance, this requirement is expressed through the assumption $(Y \perp Z) | \mathbf{X}$, effectively stating that all information

regarding sample inclusion indicator Z ought to have been gathered into the auxiliary variables \mathbf{X} .

This factor was a variant of “SITA violation #1”, the difference being somewhat a conceptual one: here, a variable existed that we ought to have observed but failed to do so while, in the previous case ($(Y \angle Z) \mid \mathbf{X}$, where \angle denotes dependence), the nature of the phenomenon was such that Y and Z were tangled and could not be untangled by conditioning.

In real conditions, the requirement is difficult to verify and presumably not strictly fulfilled. Whether it is of higher importance to measure “the background information” variable X_1 , or “the participation” variable X_2 , or both, is investigated by varying this factor.

2.1.7 SITA violation #2

When assignment to the subset is set by $Z = I_{U < X_2}$, the SITA assumption $0 < e(\mathbf{x}) \leq 1$ holds. In that case, $e(\mathbf{x}) \equiv \Phi(x_2)$, the cumulative distribution function of the standard normal variable X_2 , which is never strictly 0. But, setting, for instance, $Z = I_{\max(V, 0) < X_2}$ violates the above assumption—in words, units with x_2 less than 0 have no chance of appearing in the restricted sample r , and vice versa.

It may be shown that, in the case $Z = I_{\max(V, 0) < X_2}$ which violates the SITA assumption, the regression line of $E(Y|X_2)$ is in the present model nevertheless the same for both samples. Thus, the regression techniques are expected not to be affected by this violation but the propensity score technique is expected to be affected.

For this factor, consisting of two levels, the reference “SITA violation #2” is used below.

2.2 Summary of the studied factors

The following factors were thus included in the study:

1. Covariance structure [denoted COVSTR in the Tables and Figures]: 8 levels (the 8 models presented in Table 2),
2. Sample sizes [SSIZE]: 2 levels (“low”, $n_{low} = 1000$, and “high”, $n_{high} = 5000$, for s sample),
3. Ratio of k , the size of the sample r , to n , the size of the sample s [KNRATIO]: 3 levels ($1/2$, $2/2$, and $3/2$, giving the restricted sample’s sizes $k_{low} = \{500, 1000, 1500\}$ for the “SSIZE low” condition and $k_{high} = \{2500, 5000, 7500\}$ for the “SSIZE high” condition),
4. Observed variables [OBSERVED]: 3 levels (only X_1 observed, only X_2 observed, both X_1 and X_2 observed),
5. Method of multiple imputation [METHOD]: 5 levels (EM, MCMC/EM, MCMC/BOOT, REG, and PROP),
6. SITA violation #1 [SITAVIO1]: 2 levels (“N”, $\rho_{.4, no} = 0$, and “Y”, $\rho_{.4, yes} = -.175$),

7. SITA violation #2 [SITAVIO2]: 2 levels (“N”, $\forall i : 0 < e(\mathbf{x}_i) < 1$, and “Y”, $\exists i : e(x_i) = 0$).

2.3 Procedure

For each combination of the levels of the all the factors except METHOD and OBSERVED, $b = 1000$ independent trials were run⁶, where a trial consisted of generating a simulated population given in (1) of size $N = 50000$ with the required properties, taking an unrestricted sample s and a restricted sample r , performing the multiple imputations, and calculating the required statistics (given below under this same heading) from them.

As comparison of the multiple imputation techniques and of the effects of observing differing amount of information were of interest, the required statistics for the levels of the factors METHOD and OBSERVED were calculated on the same sets of data. To each pair of drawn samples, the five methods of multiple imputation were applied (i.e. the METHOD factor), and, within each method, multiple imputation was performed for the partial variable observation (only X_1 —the background variable, only X_2 —the participation variable) and the complete variable observation (both X_1 and X_2) (i.e. the OBSERVED factor). Based on the multiple imputations, a point estimate and an estimate of its variance were calculated using the expressions in (2) and (3).

Number of imputations of Y values into the unrestricted sample was set to always give 50000 imputed observations, the size of the population, and was thus $m = 50$ for n_{low} and $m = 10$ for n_{high} .

The experimental PROC MI of SAS 8.2 was used throughout. For the propensity score imputation technique, the default number of strata, $L = 5$, was used.

For each of the draws, two statistics were recorded:

1. bias of the point estimate, the difference between the estimator and the estimand (i.e. $\hat{Y}_{MI} - \bar{Y}$, where \bar{Y} is the population mean of Y in the current population), and
2. whether the estimand was within the nominal 95% confidence interval computed using the estimated variance in (3), that is

$$CI = \hat{Y}_{MI} \pm 1.96 \sqrt{\hat{V}(\hat{Y}_{MI})}.$$

These two statistics enabled derivation of two summary statistics for each combination of the experimental levels:

- mean bias across the b trials (*MeanBias* in the reported Figures and Tables), and

⁶It took about two and a half days for an average computer at our department to perform the 1000 runs at one such combination of the levels. Of the 192 combinations, 96 were actually run, cutting out the levels `KNRATIO=2/2` and `KNRATIO=3/2` after the first round of simulations, as explained below.

As each run consisted of doing the multiple imputation 15 times (5 METHOD levels \times 3 OBSERVED levels), it took on average a quarter of a minute for a single multiple imputation. The MCMC techniques were though much more time consuming than the other ones.

- empirical confidence level: proportion of confidence interval “hits”—the mean of the statistics in item 2 above across the b trials (*Clevel*).

When percentage reduction in bias for the summary statistic *MeanBias* is presented, it was calculated using

$$prb\left(\hat{\theta}_{\{\cdot\}}\right) = 100 \left(1 - \frac{\left| \frac{1}{b} \sum_{j=1}^b \hat{\theta}_{\{\cdot\}} - \theta \right|}{\left| \hat{\theta}_r - \theta \right|} \right),$$

where $\hat{\theta}_{\{\cdot\}}$, $\hat{\theta}_r$, and θ are the estimator adjusted using the technique and under the circumstances $\{\cdot\}$, the unadjusted estimator (based on the r sample only), and the estimand, respectively. Thus, *prb* was calculated from the summary data, and not for each generated population separately.

The statistics *MeanBias* and *Clevel* are reported as results in the next section.

3 Results

The results of the simulation are presented in tabular and graphical form. The main table⁷ of results consists of percentages reduction in bias and empirical confidence levels of the multiple imputation adjusted (MI-adjusted, for short) estimator under conformance and the deviations from the assumptions. Second-order interaction plots of the studied factors are added with the aim to give the reader an impression about the individual contributions of the studied factors on the simulation statistics, as well as about the contributions of their interactions. Two additional kinds of tables, containing more detailed information, also exist: ANOVA tables for each of the summary statistics, up to second order effects, and tables of means of the first and second order effects, across all the levels partaking in the current analysis. These tables—too large and detailed to constitute a part of the text—are given in the Appendix.

Amongst the factors that showed to have a dominating effect on the observed simulation statistics were those related to violations of the assumptions for the propensity score technique. In order to give a clear picture of the contributions of all the factors investigated, first presented is the case where all the assumptions held. Investigated there were the effects of covariance structure, sample size, ratio of the samples’ sizes and method of imputation. Then, keeping constant a factor of lesser significance (*KNRATIO*), the three deviations from the perfect situation were introduced.

There are 8 cases all in all (including the one where all the assumptions held), as the Table 3 illustrates. The results are presented in this order.

Within cases, the results are presented first for the point estimation (i.e., the simulation statistics *MeanBias*), followed by those regarding confidence levels for the point estimation (i.e., the statistics *Clevel*).

⁷Table 4 on p. 16.

Table 3: The eight cases of assumption violations.

Case	Only X_1 observed	$\rho_{.4} \neq 0$	$Z = I_{\max(V,0) < X_2}$
0	no	no	no
1	yes	no	no
2	no	yes	no
3	no	no	yes
4	yes	yes	no
5	yes	no	yes
6	no	yes	yes
7	yes	yes	yes

3.1 Case 0: All the assumptions held

In the situation where all the assumptions held and the complete information was observed, the factor with the dominating effect was METHOD: the impact of a sole imputation technique, the propensity score, overrode all the other effects for both *MeanBias* and *Clevel* (Figures 1 and 2). In order not to obscure the effects of the other factors, the analysis was split into two: the main analysis was performed without the METHOD=PROP level, while separately the performance of the propensity score as a technique for multiple imputation was compared to that of the propensity score weighting (p. 27).

With the propensity score excluded, the multiple imputation corrected practically all of the bias due to observation of the Y values only in the restricted sample r . For none of the factors and levels conforming to the SITA assumptions was *MeanBias* larger than .002 (Figure 3), giving a reduction in bias of at least 98.5% for any particular combination of the levels. Across all the levels, the mean percent reduction in bias was 99.8%. Likewise, all the confidence intervals were on at least the nominal 95% level (Figure 4), more than half of them though somewhat conservative, actually achieving a 98 – 99% confidence level.

Because the factor METHOD, with the propensity score technique excluded, showed no impact in the ANOVA decompositions for the simulation statistics (Tables I and II in the Appendix), the results were in the following taken across all the four remaining levels of this factor.

The other factors did have a significant impact in the analysis of variance of both simulation statistics, *MeanBias* and *Clevel*. Except for COVSTR, their effects were straightforward.

For *MeanBias*, increase of SSIZE decreased the bias of the MI-adjusted point estimator, and the increase of KNRATIO decreased the bias of the estimator (Figure 3). When evaluating the effect of COVSTR, it ought to be recalled that the 8 covariance structures had two levels of the original bias. This bias was a function of ρ_{23} , the correlation between the participation variable and the study variable, which itself had two levels, $\rho_{23,low} = .22$ and $\rho_{23,high} = .78$, giving the bias ($\rho_{23} \times \pi^{-\frac{1}{2}}$) of either .124 or .440. The results indicate that percentage reduction in bias with MI-adjustment was very good (Table 4), but was somewhat lower for the structures with low ρ_{23} (i.e. 1, 3, 4, and 7), 99.7% on average, than for those with high ρ_{23} , 99.9% on average.

With respect to the statistic *Clevel*, the factors had the following effect. Increase in

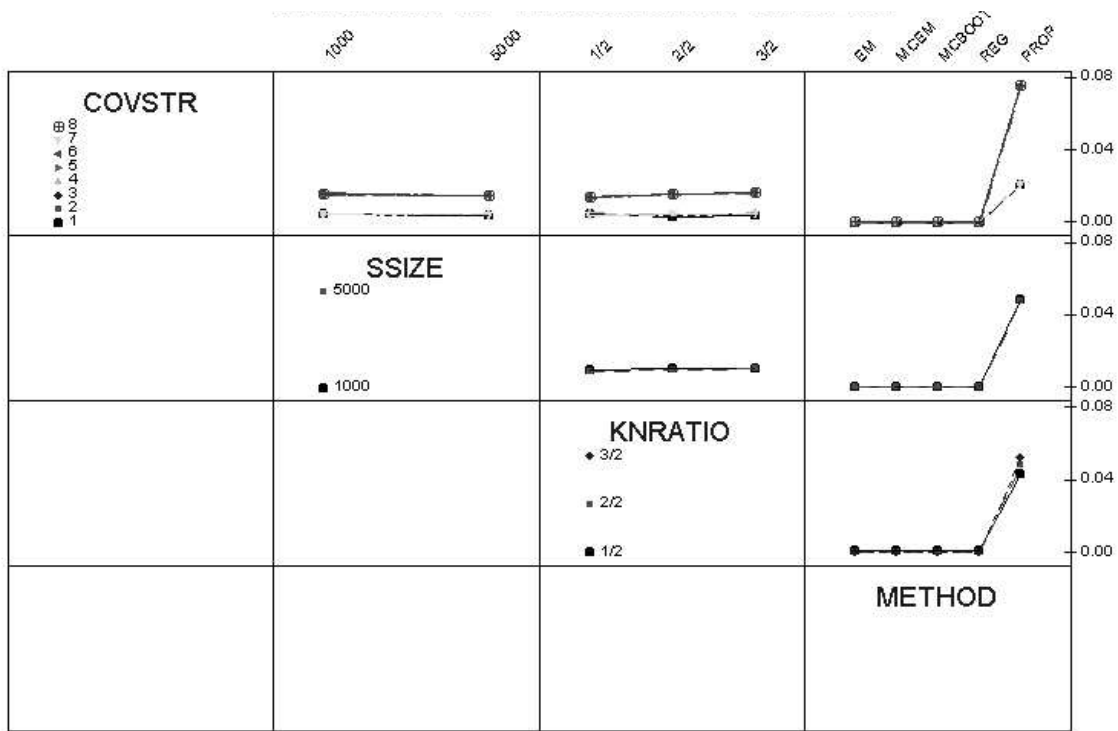


Figure 1: Case 0, interaction plot for *MeanBias*, with METHOD=PROP included. (Identification of COVSTR levels given in the text.)

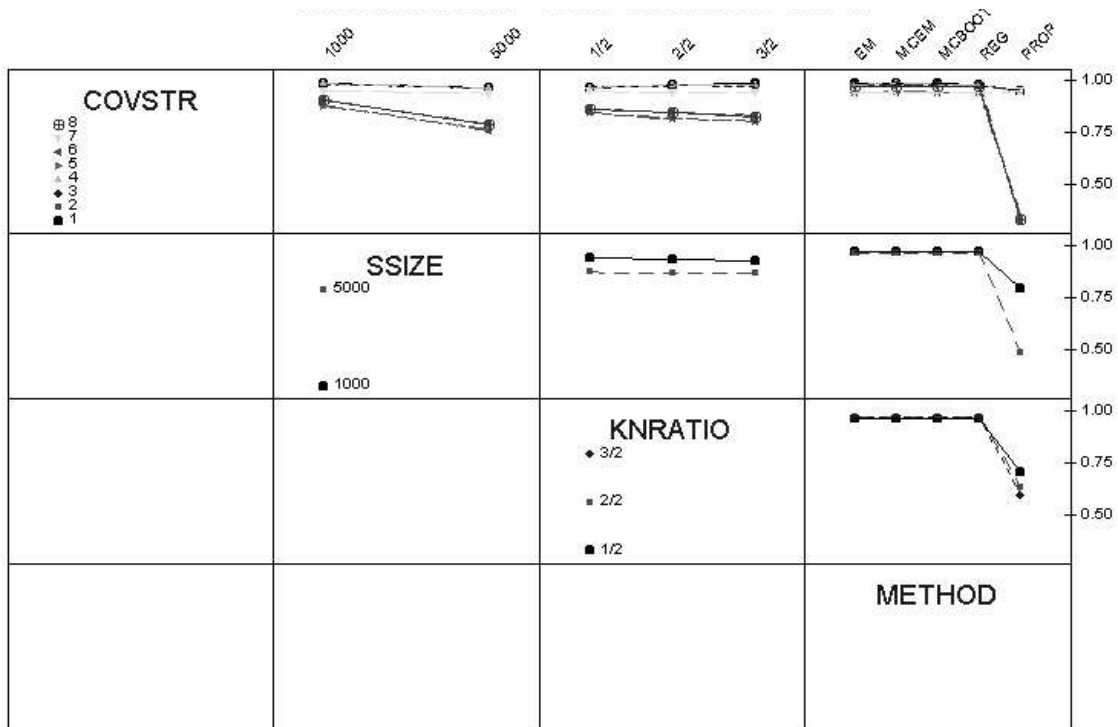


Figure 2: Case 0, interaction plot for *Clevel*, with METHOD=PROP included. (Identification of COVSTR levels given in the text.)

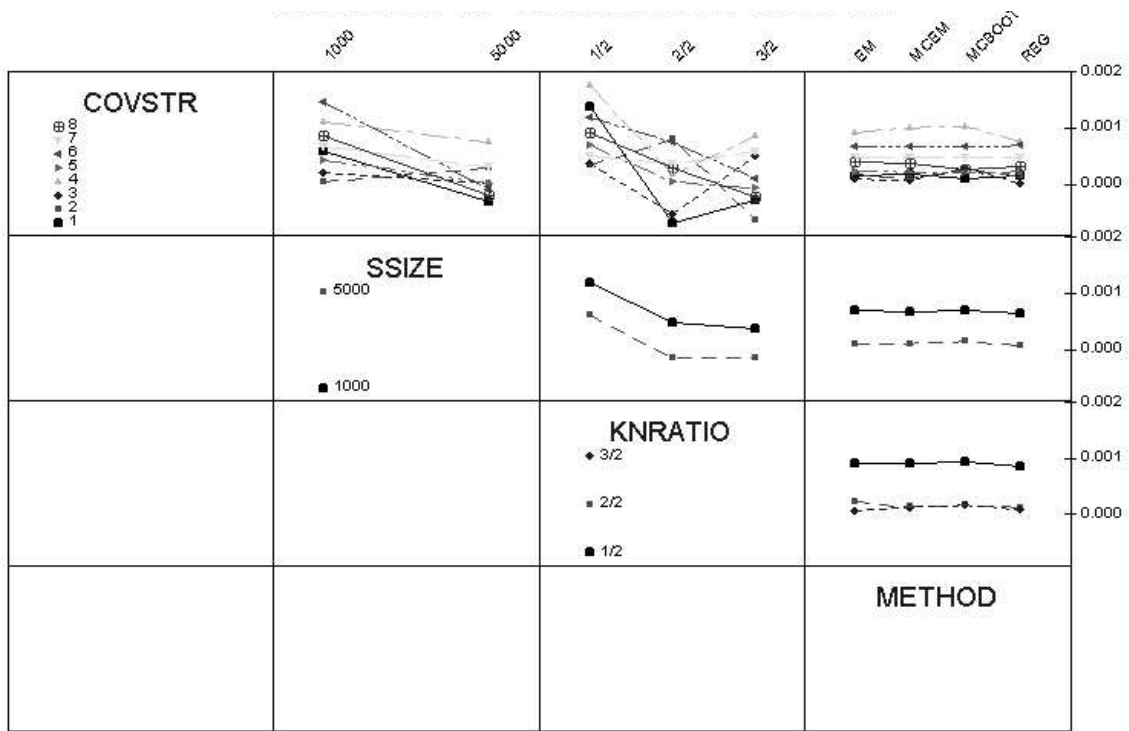


Figure 3: Case 0, interaction plot for *MeanBias*, with METHOD=PROP excluded. (Identification of COVSTR levels given in the text.)

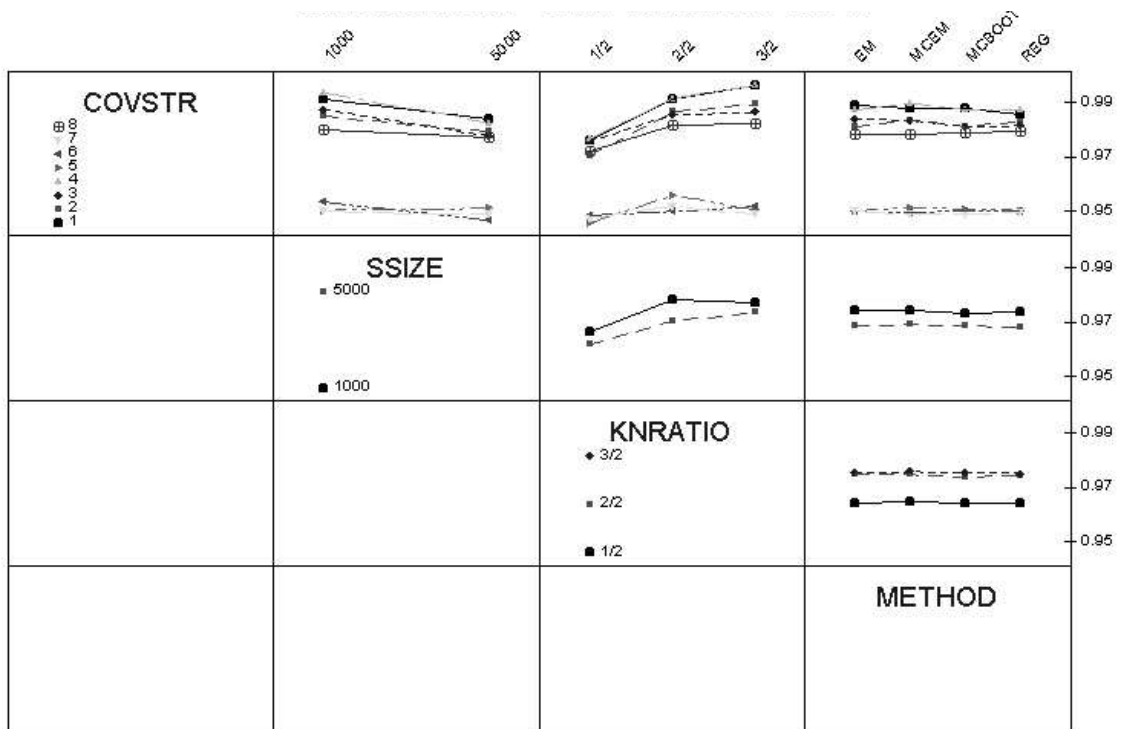


Figure 4: Case 0, interaction plot for *Clevel*, with METHOD=PROP excluded. (Identification of COVSTR levels given in the text.)

Table 4: The unadjusted estimators \bar{Y}_r (the means across all drawn populations), and the adjusted estimators $\hat{Y}_{MI,}$ with their corresponding percentages reduction in bias (*prb*) and empirical confidence levels of the nominal 95 percent confidence intervals (*Clevel*) for the 8 treated cases and, within each, for the 8 covariance structures across all the levels not defining a case.

Note: the table was set up the way that would enable an easy visual comparison with the corresponding results in (Lorenc, 2003b), which had to accomodate an additional level of SITAVIO1, $\rho_{.4} = .175$. The results of the neutral case (no assumption violations) are in the latter table placed approximately in its middle, just below the two componenet tables regarding only SITAVIO2 (i.e., with the indices 3 and 5). Here, the table starts with these two componenet tables. Also, what are here denoted as Cases 2, 4, 6, and 7 are in the other table more correctly denoted as Cases $-2, -4, -6,$ and -7 as they refer to a negative $\rho_{.4}$.

COVSTR	ρ_{12}	ρ_{13}	ρ_{23}	\bar{Y}_r	$\hat{Y}_{MI,3}$	<i>prb</i>	<i>Clevel</i>	$\hat{Y}_{MI,5}$	<i>prb</i>	<i>Clevel</i>
1	.22	.22	.22	.200	.001	100	.966	.120	40	.252
2	.22	.22	.78	.708	.000	100	.964	.511	28	0
3	.22	.78	.22	.200	.001	99	.961	.034	83	.768
4	.78	.22	.22	.200	-.001	100	.960	.049	75	.782
5	.22	.78	.78	.708	.000	100	.963	.424	40	0
6	.78	.22	.78	.708	.000	100	.958	.608	14	0
7	.78	.78	.22	.200	.000	100	.956	-.389	-95	0
8	.78	.78	.78	.708	.000	100	.961	.172	76	.054
COVSTR	ρ_{12}	ρ_{13}	ρ_{23}	\bar{Y}_r	$\hat{Y}_{MI,0}$	<i>prb</i>	<i>Clevel</i>	$\hat{Y}_{MI,1}$	<i>prb</i>	<i>Clevel</i>
1	.22	.22	.22	.124	.000	100	.988	.099	20	.342
2	.22	.22	.78	.440	.000	100	.982	.419	5	0
3	.22	.78	.22	.124	.000	100	.983	.028	78	.838
4	.78	.22	.22	.124	.001	99	.988	.035	71	.878
5	.22	.78	.78	.440	.000	100	.951	.347	21	0
6	.78	.22	.78	.440	.001	100	.950	.424	4	0
7	.78	.78	.22	.124	.000	100	.949	-.272	-119	0
8	.78	.78	.78	.440	.000	100	.979	.121	73	.107
COVSTR	ρ_{12}	ρ_{13}	ρ_{23}	\bar{Y}_r	$\hat{Y}_{MI,2}$	<i>prb</i>	<i>Clevel</i>	$\hat{Y}_{MI,4}$	<i>prb</i>	<i>Clevel</i>
1	.22	.22	.22	.206	.104	49	.390	.171	17	0.092
2	.22	.22	.78	.497	.030	94	.838	.648	-30	0
3	.22	.78	.22	.206	.069	66	.609	.049	76	0.594
4	.78	.22	.22	.206	.112	46	.361	.097	53	0.493
5	.22	.78	.78	.497	-.004	99	.637	.526	-6	0
6	.78	.22	.78	.497	.048	90	.525	.852	-71	0
7	.78	.78	.22	.206	.089	57	.286	-.515	-150	0
8	.78	.78	.78	.497	.025	95	.859	.241	52	0
COVSTR	ρ_{12}	ρ_{13}	ρ_{23}	\bar{Y}_r	$\hat{Y}_{MI,6}$	<i>prb</i>	<i>Clevel</i>	$\hat{Y}_{MI,7}$	<i>prb</i>	<i>Clevel</i>
1	.22	.22	.22	.245	.080	67	.651	.200	18	.022
2	.22	.22	.78	.719	.022	97	.903	.694	4	0
3	.22	.78	.22	.245	.055	78	.760	.058	77	.490
4	.78	.22	.22	.245	.086	65	.626	.132	46	.300
5	.22	.78	.78	.719	-.002	100	.749	.551	23	0
6	.78	.22	.78	.719	.038	95	.621	.912	-27	0
7	.78	.78	.22	.245	.069	72	.477	-.523	-114	0
8	.78	.78	.78	.719	.020	97	.910	.257	64	0

SSIZE moved the empirical confidence level somewhat towards the nominal one: from .974 to .969, while increase in KNRATIO moved the empirical confidence level away from the nominal one: from .964 to .974 to .975 (Figure 4). There was a significant interaction between the factors (Table II in the Appendix). The structures in COVSTR fell into two groups with respect to *Clevel*: a smaller one (the structures 5, 6, and 7) with the nominal confidence level practically coinciding with the empirical one, 95%, and a larger one with the nominal confidence level too conservative with respect to the empirical one, 98–99%. The three covariance structures in the former group are characterised by having exactly two of the three correlation coefficients (ρ_{12} , ρ_{13} , ρ_{23}) high⁸.

The length of the simulations necessitated that one of the factors be held at a constant level. The factor chosen for this was KNRATIO, with its *worst* performing level kept. Regarding the choice of the factor, both COVSTR and SSIZE seemed indispensable; regarding the choice of the level, its implication was that the actual results could not be *worse* than those reported, seen in the frame of reference of the originally conceived experiment. Missing though will still be the interactions between the levels of the factor KNRATIO and of those of the other factors. At the level KNRATIO=1/3, the MI-adjusted point estimator did have a residual bias (about .0009 across all the levels), and it is from this point of departure that the rest of the analysis was performed. In practical applications the situation is probably much worse than this.

The results reported thus far were obtained using complete information, that is, both X_1 and X_2 were available for the multiple imputation procedures. Availability of only X_2 (i.e. exclusion of X_1) does not represent a violation of the studied assumptions as $X_1 \perp V$ when $\rho_{.4} = 0$ (which was the case thus far), why even $(Y \perp Z) | X_2$ instead of $(Y \perp Z) | \mathbf{X}$ holds. In continuing the analysis for Case 0 (no assumption violations), the factor OBSERVED replaced the previous KNRATIO.

The factor OBSERVED interacted with the factor COVSTR in the following way. A change from observing only X_2 to observing both auxiliary variables influenced *MeanBias* (Figure 5) of four of the covariance structures (5–8), converging it towards the common value—which, it may be recalled, is about .0009 at the current level of KNRATIO. A change from observing only X_2 to observing both X_1 and X_2 took also the *Clevel* (Figure 6) of three of these four covariance structures—all but number 8—from the common one of .97 to the nominal .95. For the other covariance structures, neither *MeanBias* nor *Clevel* were affected by the factor OBSERVED. And, the other interactions of this factor seem to be just the consequence of its aforementioned interactions with COVSTR.

3.2 Case 1: Participation variable not observed

Next in the analysis, also the level of observing only X_1 was included. Observing only X_1 amounts to observing incomplete information: the participation variable X_2 is required instead in order to ‘explain’ Z . It might be recalled that the original bias was due to the correlation ρ_{23} between the participation variable and the study variable. So, insofar

⁸These three matrices are close to singular, with the consequence that three of the four variables practically determine the remaining one. I am grateful to Daniel Thorburn for pointing this out.

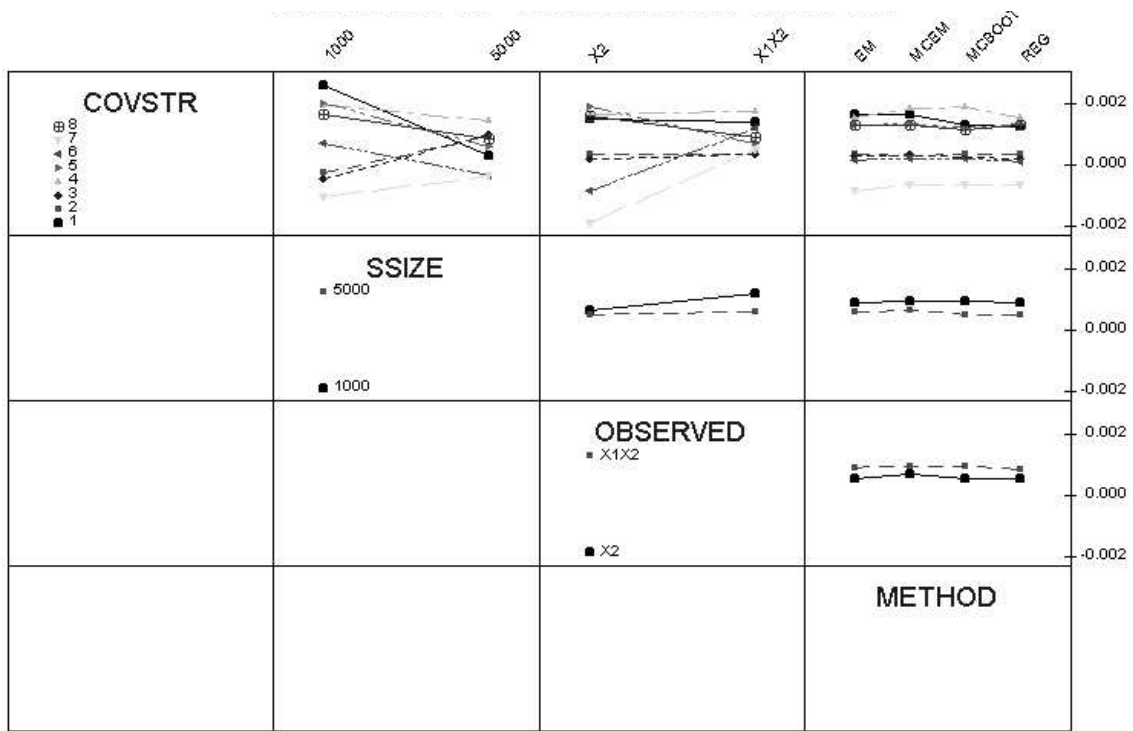


Figure 5: Case 0, interaction plot for *MeanBias*—the factor KNRATIO replaced by OBSERVED. (Identification of COVSTR levels given in the text.)

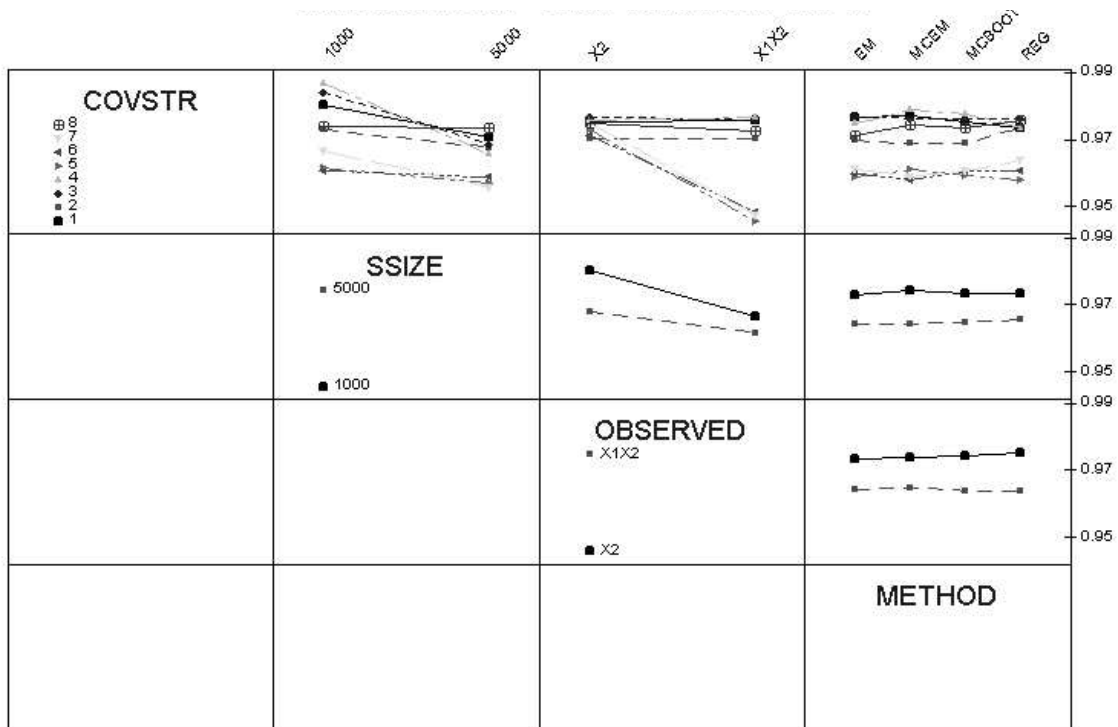


Figure 6: Case 0, interaction plot for *Clevel*—the factor KNRATIO replaced by OBSERVED. (Identification of COVSTR levels given in the text.)

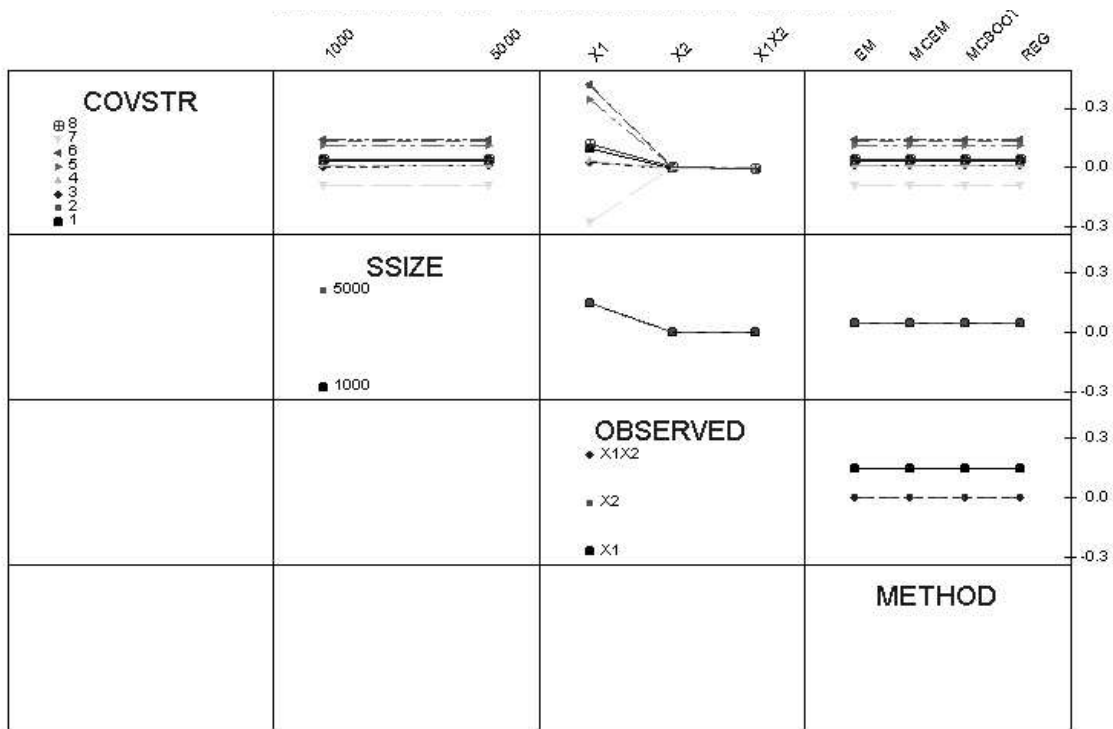


Figure 7: Case 1, interaction plot for *MeanBias*. (Identification of COVSTR levels given in the text.)

as the observed X_1 and X_2 or X_1 and Y would be correlated, it might be expected that observation of X_1 would help in correcting the bias of the unadjusted estimator \bar{Y}_r .

Observing only X_1 (i.e. OBSERVED= X_1) had a large effect on both the point estimator (Figure 7) and on its estimated nominal 95% confidence level derived from the point estimate and its estimated variance (Figure 8). For the point estimator, the dominating effects in the ANOVA decomposition (Table VII in the Appendix) were those of OBSERVED and COVSTR, as well as of their interaction. It is notable that neither SSIZE nor METHOD had a significant effect.

When the original bias of \bar{Y}_r as the estimator of \bar{Y} was high, observing an auxiliary variable highly correlated with only either X_2 or Y contributed little to the bias reducing power of the MI-adjusted estimate (rows 2, 5, and 6 vs. row 8 corresponding to the estimator $\hat{Y}_{MI,1}$ in Table 4). Both correlations needed to be strong in order to achieve a larger—here 73%—reduction in bias. When the original bias was low, it took exactly one of the correlations of the observed X_1 with X_2 or Y to be high to achieve this same level of reduction in bias (row 1 vs. rows 3 and 4 corresponding to the estimator $\hat{Y}_{MI,1}$ of Table 4). That both were high though became detrimental to the efficiency of the adjustment, by overadjusting in the negative direction thus doubling the original bias (row 7 corresponding to the estimator $\hat{Y}_{MI,1}$ of Table 4). An analogous effect when observing only X_1 was noticed in the related simulation study of the efficiency of the propensity score weighting (Lorenc, 2003b).

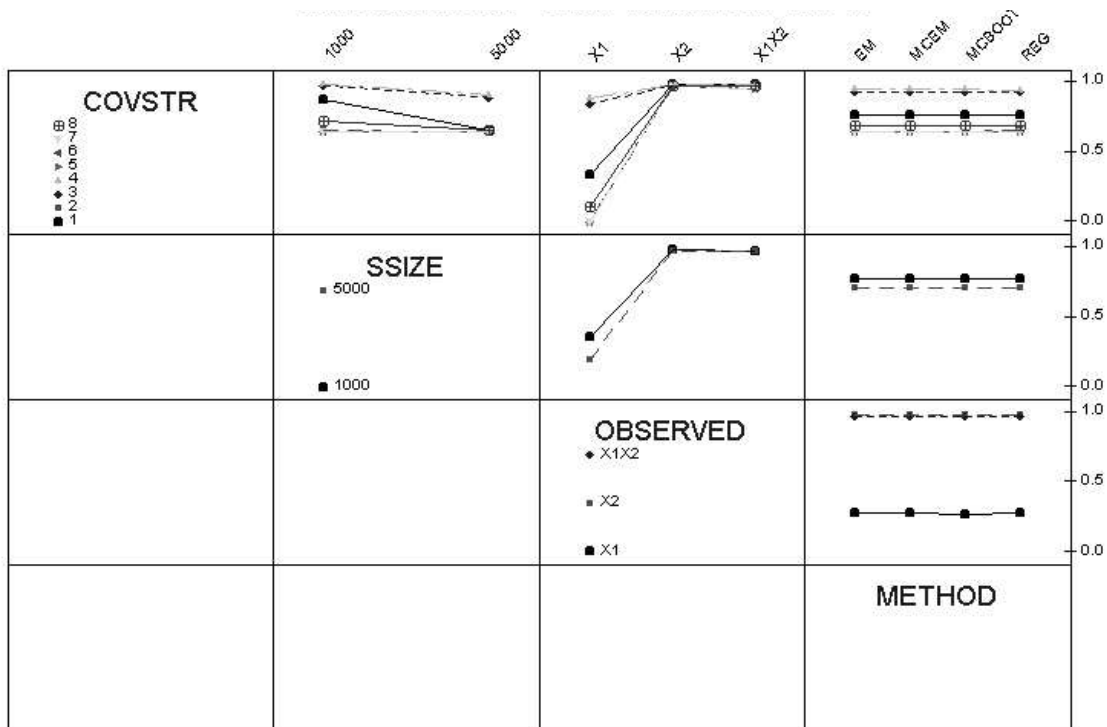


Figure 8: Case 1, interaction plot for *Clevel*. (Identification of COVSTR levels given in the text.)

It can be also noted, from Figure 7, that under the studied conditions and other things being equal, it was by far more important to observe an auxiliary variable strongly related to participation (here, X_2) than an auxiliary variable not strongly related to participation (here, X_1).

Performance of *Clevel*, with only X_1 observed, followed the performance of *MeanBias* in the same situation. The covariance structures for which the point adjustment had little or no effect (or an adverse effect) were also completely off the mark with respect to the calculated confidence intervals (the structures 2 and 5–7) or much below the targeted 95% level (the structures 8 and 1). The empirical *Clevel* of the adjusted estimator was for the remaining two covariance structures (3 and 4) closer to but still below the nominal level. All the first order effects but METHOD had a highly significant contribution in the ANOVA decomposition (Table VIII in the Appendix).

3.3 Case 2: Y and Z correlated after all relevant information observed (SITA violation #1)

The dependence between the study variable Y and the indicator of the subset membership Z that remains after conditioning their joint distribution on the auxiliary information, symbolically represented with $(Y \angle Z) \mid \mathbf{X}$, violates one of the assumptions of SITA. In such a situation, in words, there exists information in the subset membership Z about Y that is not available for adjustment. This factor, named SITAVIO1, was now added to the 4 previously analysed ones: COVSTR, SSIZE, METHOD, and OBSERVED (excluding OBSERVED= X_1).

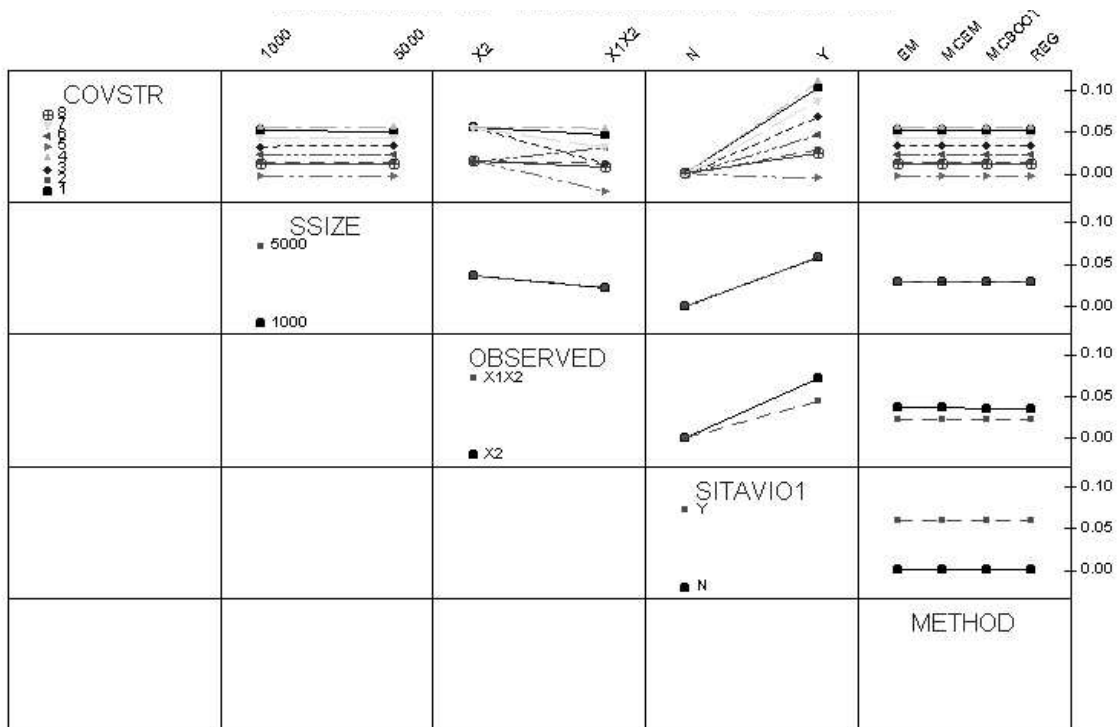


Figure 9: Case 2, interaction plot for *MeanBias*. (Identification of COVSTR levels given in the text.)

A change in the covariance structure from a $\rho_{.4} = 0$ to a negative $\rho_{.4}$ moved the MI-adjusted point estimates in the positive direction, thereby increasing *MeanBias*. (Figure 9) and decreasing *Clevel* (Figure 10). The effect on both statistics was differential with respect to the covariance structures, that is, the original bias depending on ρ_{23} . For the point estimate, the largest effect in terms of both *MeanBias* and *prb* was on the structure number 4 (whose *prb* was practically annulled, being reduced to only 10%), and then, in the descending order of any of the two statistics, 1, 7, 3, 6, 2, 8, and 5. The structures with the high original bias, that is, those with high correlation between the participation variable and the study variable (2, 5, 6, and 8), were also more robust to the violation, compared with the other group (1, 3, 4, and 7): the average *prb* for the groups was 95% and 55% respectively. Again, an analogous effect was noticed in the related simulation study of the efficiency of the propensity score weighting (Lorenc, 2003b). All the included first order effects but SSIZE had a highly significant contribution in the ANOVA decomposition (Table X in the Appendix).

Introduction of a serious bias in the adjusted point estimates, achieved through introducing a nonzero $\rho_{.4}$, had a consequence even for the confidence levels based on these estimates. In general, they were not able to hold the nominal level, the highest empirical *Clevel* being .859 (related to structure 8), and the lowest (related to structure 7) being .286. Presumably because they were more robust to the violation in the case of the preceding point estimates, the structures with the high original bias (2, 5, 6, and 8)

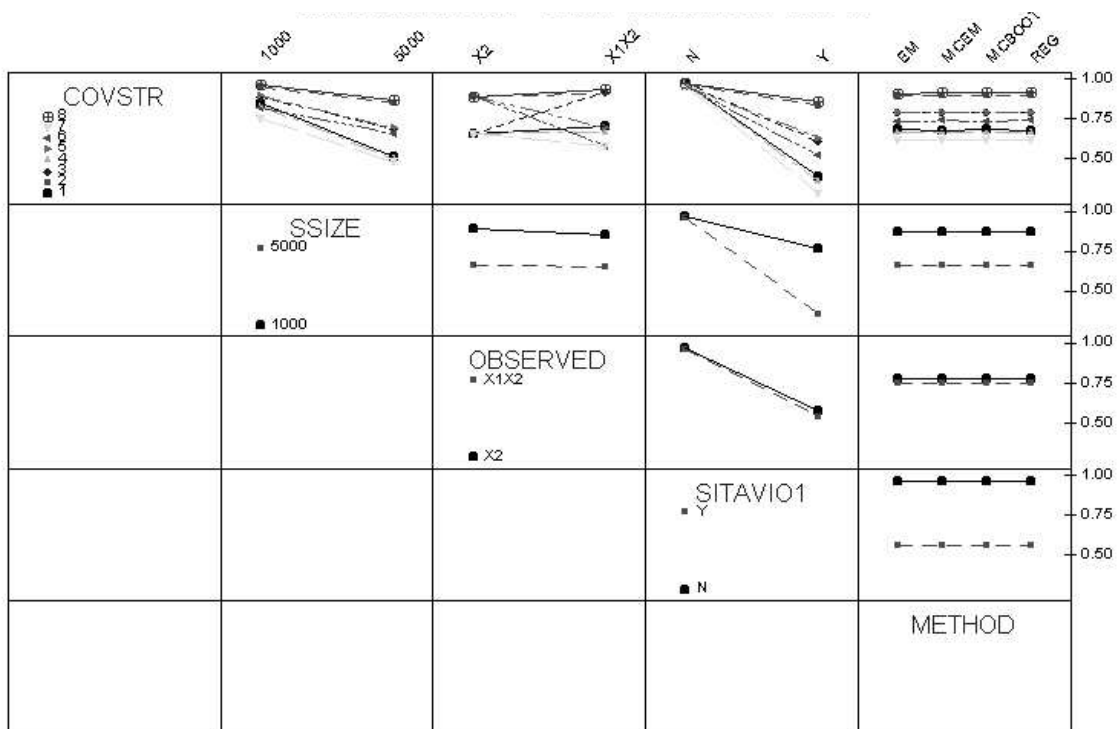


Figure 10: Case 2, interaction plot for *Clevel*. (Identification of COVSTR levels given in the text.)

were also more robust in the case of the confidence levels; the difference was though not as pronounced as previously: .715 for the high original bias group and .411 for the low original bias group.

3.4 Case 3: Not all units given a positive probability to appear in r (SITA violation #2)

When the determining property of the subset was $Z = I_{\max(V,0) < X_2}$, instead of $Z = I_{V < X_2}$ which was applied thus far, no unit in the subset may have taken on a negative value of X_2 . In other words, units with the negative X_2 had no chance of appearing in the sample from the subset, r . This level, termed “SITA violation #2”, was introduced in the analysis next.

Existence of the level SITAVIO2=YES in the simulations had little impact on the point estimates and on the confidence levels for these estimates. The overall *MeanBias* for this factor, across all the other levels, changed from .0007 to .0002 (it may be recalled that the level KNRATIO=1/2, on which the simulations after Case 0 were run, did have a residual bias that was estimated to be about .0009 across the other levels). The most important interaction between SITAVIO2 with COVSTR was through the covariance structure 4, whose mean bias was lowered by almost .002 (Figure 11 and Table XV in The Appendix).

There was still a considerable resemblance between the interaction plots for *MeanBias* in the present case and in Case 0 (Figure 11 compared to Figure 5); similar comparison

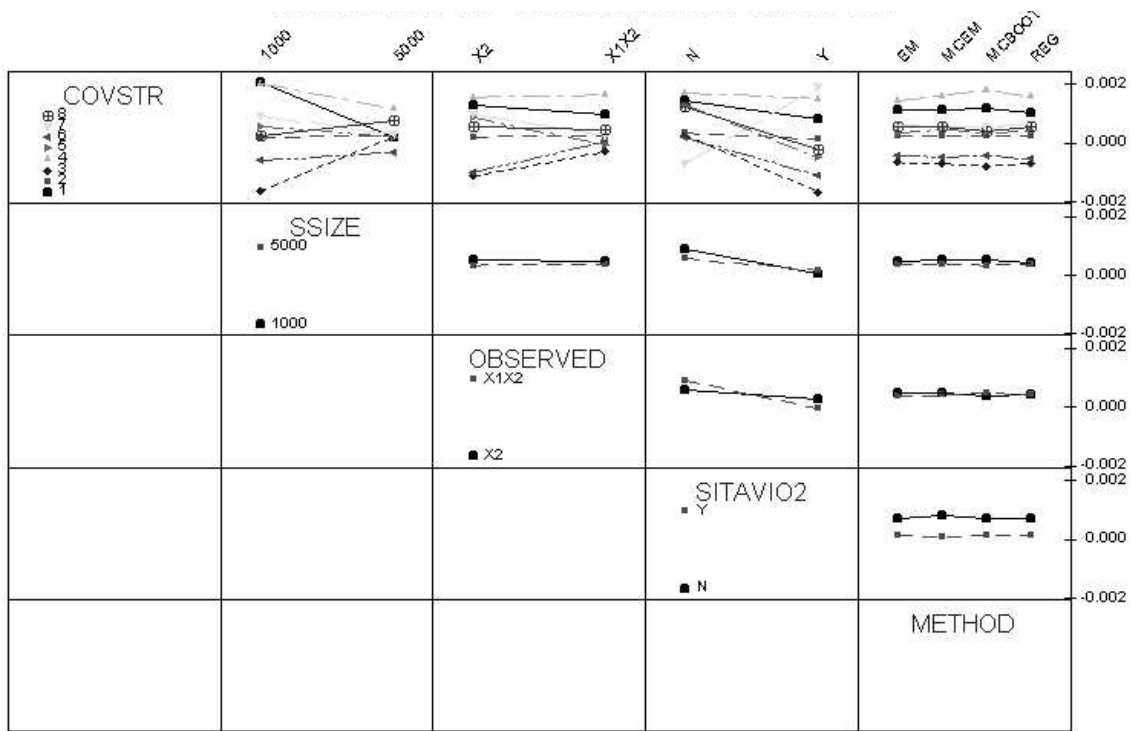


Figure 11: Case 3, interaction plot for *MeanBias*. (Identification of COVSTR levels given in the text.)

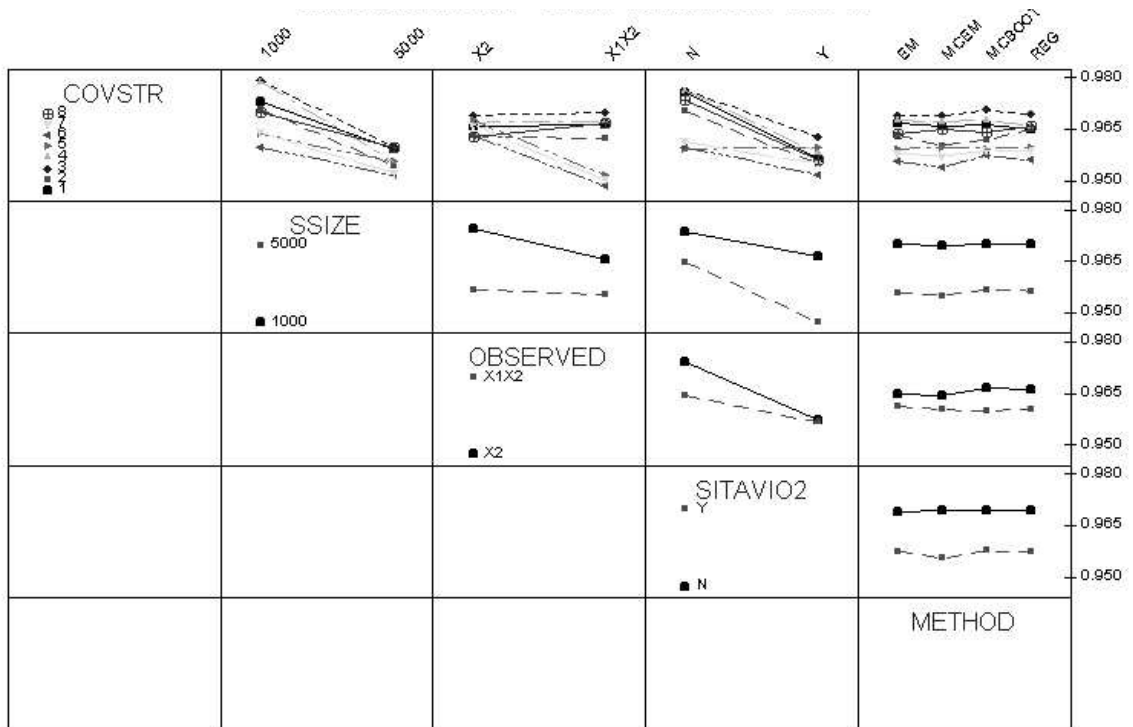


Figure 12: Case 3, interaction plot for *Clevel*. (Identification of COVSTR levels given in the text.)

for *Clevel*—Figure 12 with Figure 6—also showed little change. The percentage reduction in bias, too, has changed little by introducing SITA violation #2, as column *prb* corresponding to $\hat{Y}_{MI,3}$ in Table 4 illustrates.

3.5 Case 4: Participation variable not observed and SITA violation #1

The cases 4-6 present pairwise combinations of the SITA violations investigated under Case 1 – Case 3. Manner of the presentation is brief, but the details concerning all the simulation statistics can still be found in the Appendix.

It may be recalled that observing only X_1 —treated as Case 1 above—did have a selective impact on the *prb* of the MI-adjusted estimator, depending on the covariance structure: being still of some help for some of them (3, 4, and 8), of little help for the others (1, 2, 5, and 6), and devastating in one case (7), more than doubling the original bias. And, also that introduction of SITA violation #1—treated as Case 2—in general moved the point estimates in the positive direction, to which the structures with the high original bias were more robust (*prb* of about 95% after adjustment) than those with the low bias (*prb* of about 55% after adjustment). For no structure was the adjusted estimator, with SITA violation #1 present, more biased than the unadjusted estimator \tilde{Y}_r .

Introduction of both violations simultaneously did have much more detrimental consequence on the point estimator than their individual introduction have had (*MeanBias* given in Figure 13 and column *prb* corresponding to $\hat{Y}_{MI,4}$ in Table 4). It was here, too, conditional on the covariance structure, following a pattern similar to that discussed concerning Case 1. Even here, the effect of the strong violations was an unusable adjusted point estimator, for half of the structures actually increasing the original bias, with an unusable confidence level, not larger than 60% but for most of the structures actually zero.

In the cases where the point estimate is seriously biased there is little sense in building confidence intervals around these wrongly placed points, why the results concerning *Clevel* are not presented here (but can be found in the Appendix, Tables XVII-XVIII).

3.6 Case 5: Participation variable not observed and SITA violation #2

The effect of the sole assumption violation *SITAVIO2=YES*, presented as Case 3, was with the MI-adjusted estimator negligible with respect to both *MeanBias* and *Clevel*. This effect also dampened here the strong and differential (conditional on the covariance structures) influence of observing only X_1 , reducing the adjusted estimators bias somewhat (column *prb* corresponding to $\hat{Y}_{MI,5}$ in Table 4 compared to column *prb* corresponding to $\hat{Y}_{MI,1}$ in Table 4). But, even with this moderating effect of *SITAVIO2=YES*, the consequences were still damaging for *Clevel* with most of the covariance structures, the exceptions being 3 and 4 with about 77% confidence level.

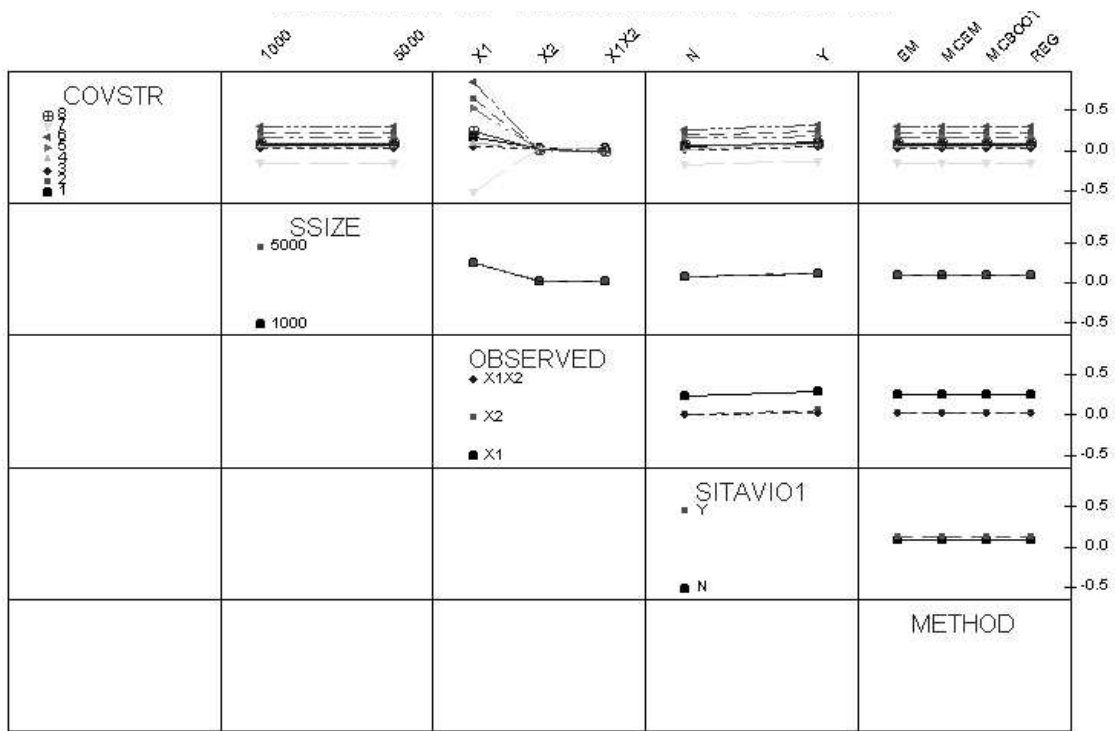


Figure 13: Case 4, interaction plot for *MeanBias*. (Identification of COVSTR levels given in the text.)

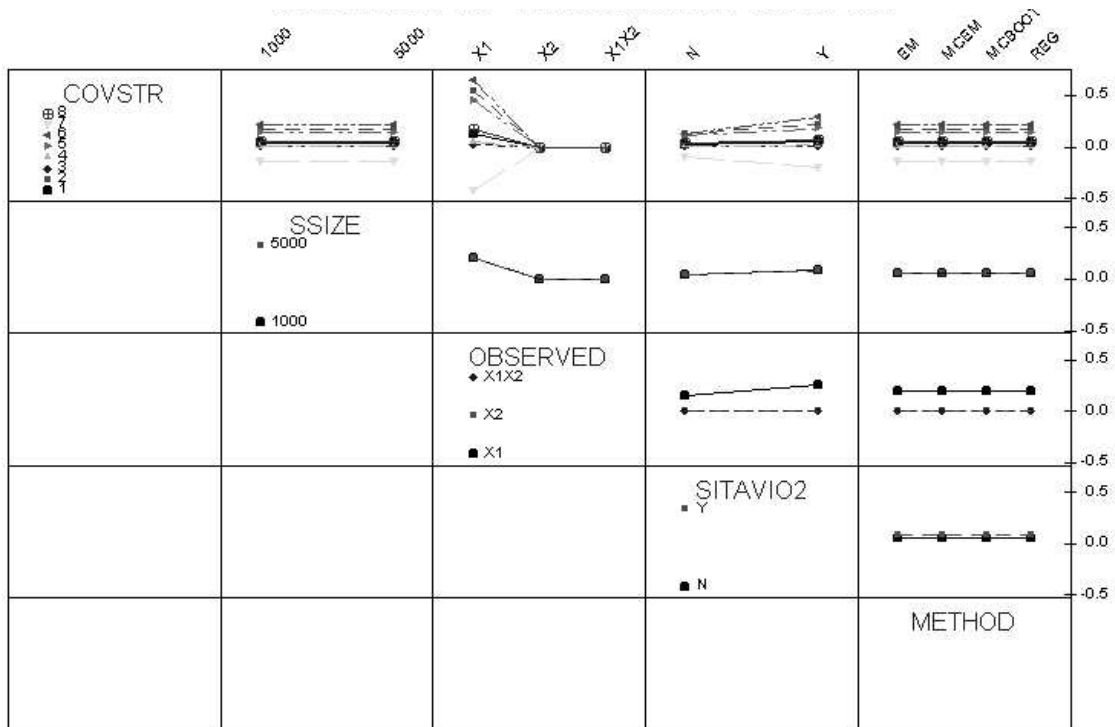


Figure 14: Case 5, interaction plot for *MeanBias*. (Identification of COVSTR levels given in the text.)

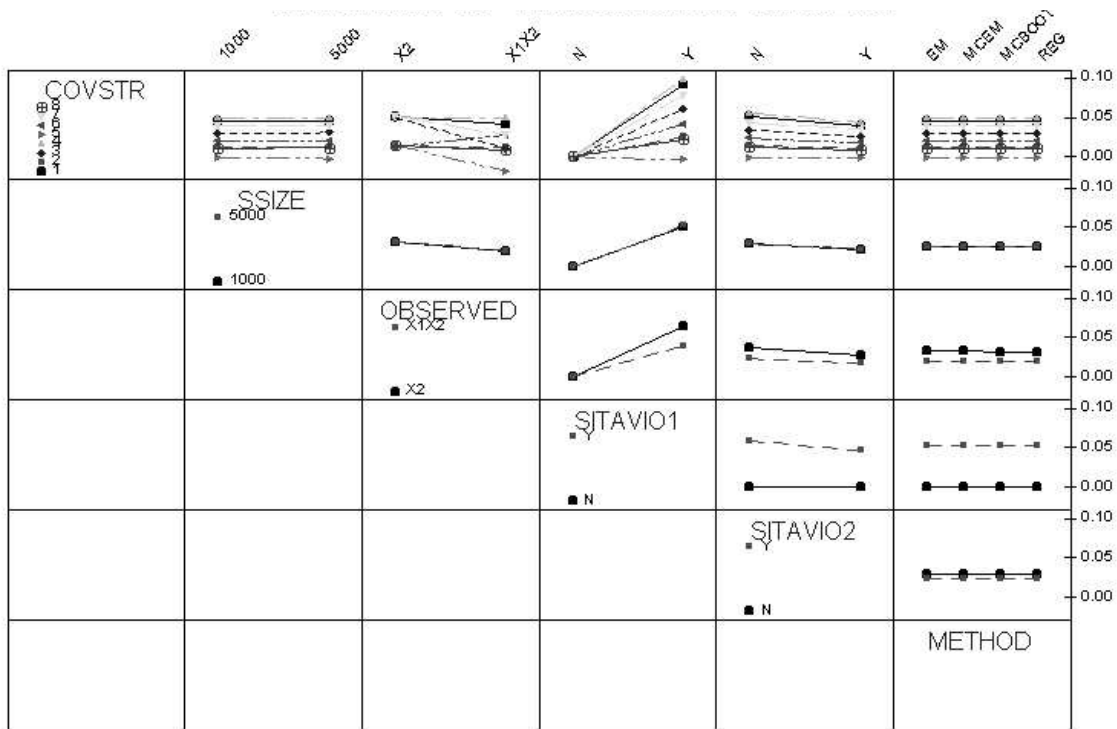


Figure 15: Case 6, interaction plot for *MeanBias*. (Identification of the COVSTR levels given in the text.)

3.7 Case 6: SITA violations #1 and #2

The situation with all the information observed, but with the two SITA violations present at the same time, did not differ much from the one with only the first of the SITA violations (Case 2). Actually, even here—as with the preceding Case 5—the violation SITA violation #2 dampened the influence of SITA violation #1, reducing the bias of the adjusted estimators caused by this latter violation somewhat (column $\hat{Y}_{MI,6}$ in Table 4 compared to column $\hat{Y}_{MI,2}$ in Table 4): the average *prb* for the more robust group with the high original bias (2, 5, 6, and 8) was 97%, while for the other group (1, 3, 4, and 7) it was 55%—for both somewhat higher than in Case 2. The confidence levels were thus also somewhat higher compared to Case 2: .796 for the high original bias group and .629 for the low original bias group.

3.8 Case 7: All the violations at the same time

Finally, simultaneous introduction of all the studied assumption violations into the analysis produced results that were in accord with the ones from the two preceding Cases. Presence of SITA violation #2 had a certain dampening effect on *MeanBias*, compared to the corresponding Case 4 (column $\hat{Y}_{MI,7}$ in Table 4 compared to column $\hat{Y}_{MI,4}$ in Table 4). This unfortunately did not suffice to produce usable confidence levels.

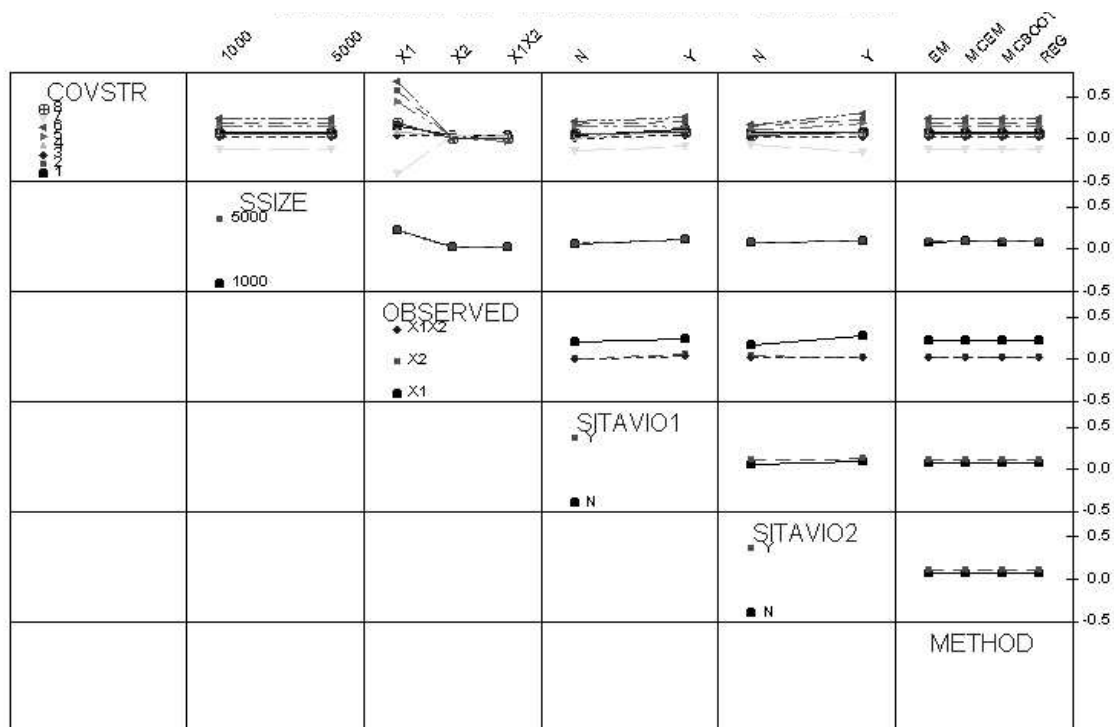


Figure 16: Case 7, interaction plot for *MeanBias*. (Identification of the COVSTR levels given in the text.)

3.9 The propensity score method

An analysis of the results of the multiple imputation adjustment when the technique for multiple imputation is the propensity score (SAS Institute Inc., 2001) was performed separately of that of the others due to the large differences that resulted from the use of this specific technique in comparison to the other four techniques. The analysis is presented in two parts, first the one concerning the situation where all the assumptions held (corresponding to Case 0 above), and the next the one where all the violations were present simultaneously (corresponding to Case 7 above).

3.9.1 Case 0: method=prop, all the assumptions held

Three factors were investigated when all the assumptions held, COVSTR, KNRATIO and SSIZE. All three had a significant effect in the ANOVA decomposition, in the order mentioned (Table XXVIII in the Appendix). There were two distinct groups of covariance structures with respect to both *MeanBias* and *Clevel*: those with a high original bias on the one hand (the structures 2, 5, 6, and 8) and those with a low original bias on the other (1, 3, 4, and 7). The residual bias was larger for the former group as well as the confidence level lower, compared with the latter group. Percentage reduction in bias was though the same in both groups—83%.

The other two factors had the following effects: raise in KNRATIO raised the bias of the MI-adjusted estimate (Figure 17) and—for the structures with high original bias—lowered the confidence level (Figure 18). Raise in SSIZE lowered (marginally) the bias

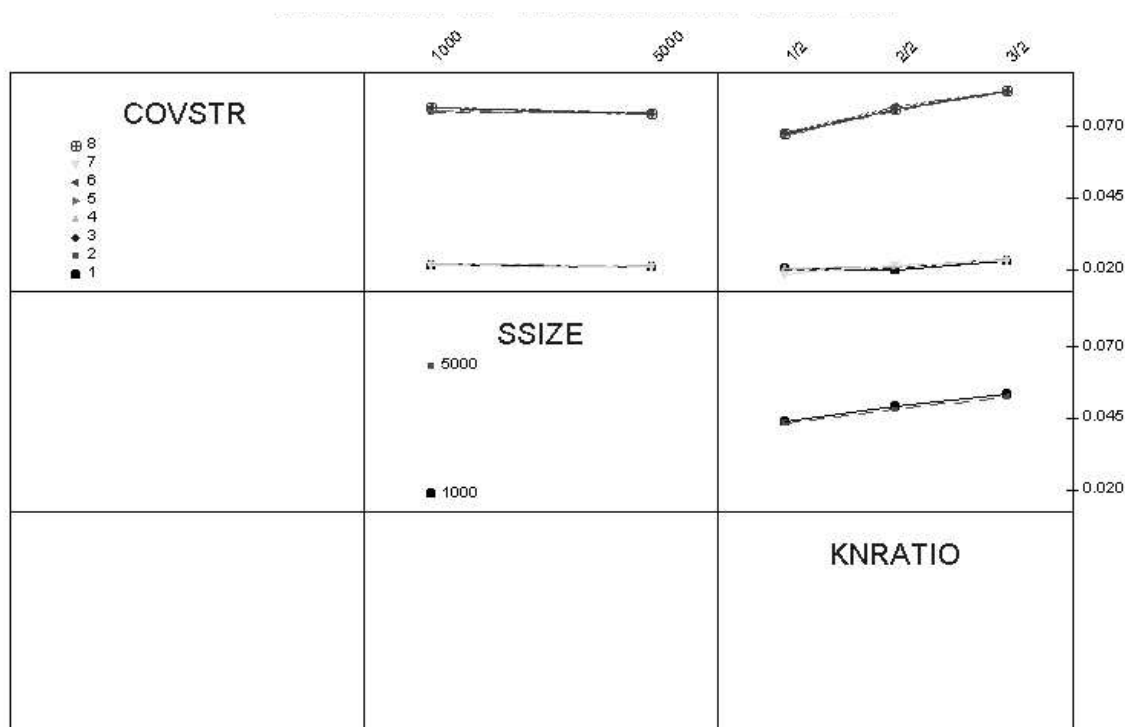


Figure 17: Case 0, only the level METHOD=PROP, interaction plot for *MeanBias*. (Identification of the COVSTR levels given in the text.)

but also lowered markedly the confidence level for the adjusted estimate, in relation to the nominal confidence level. Even here was the impact on *Clevel* dependent on the covariance structure: structures with a high original bias suffered a larger confidence level loss.

As results from a simulation study of the propensity score weighting (PS-weighting) based on the same population model as the present study were available (Lorenc, 2003b), it was also possible to compare the efficiency of that weighting technique, discussed under heading “The propensity score approach” in Section 1 above, with MI-adjustment when the propensity score was used for multiple imputations (i.e. METHOD=PROP in the present study). The corresponding results for *MeanBias* and *Clevel* from that study are given in Figures 19 and 20, respectively.

With the PS-weighting, change in KNRATIO did not have an effect on *MeanBias*, in contrast to MI-adjustment, where increase in KNRATIO increased the bias of the adjusted point estimator. The pattern of influence of KNRATIO on *Clevel* was similar between the adjustment approaches, the difference being that, for MI-adjustment, structures with low original bias did not interact with KNRATIO, but that for the PS-weighting there was an interaction: increase in KNRATIO decreased the confidence level.

The pattern of the interaction plots was in general similar between the approaches.

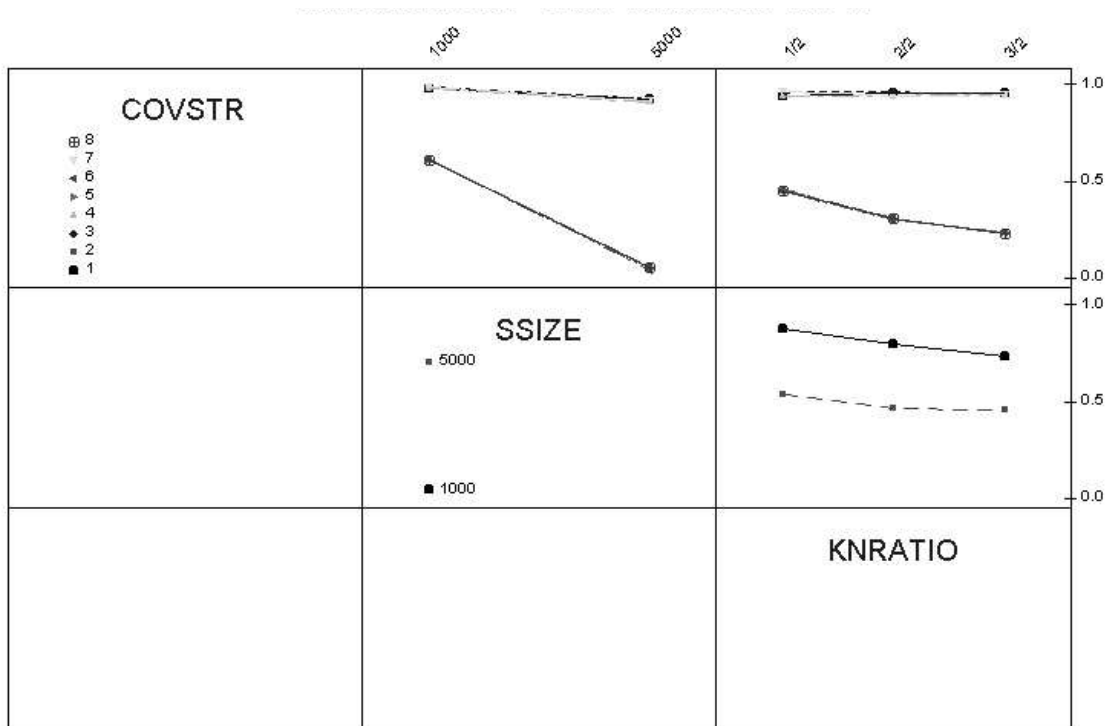


Figure 18: Case 0, only the level METHOD=PROP, interaction plot for *Cleveland*. (Identification of the COVSTR levels given in the text.)

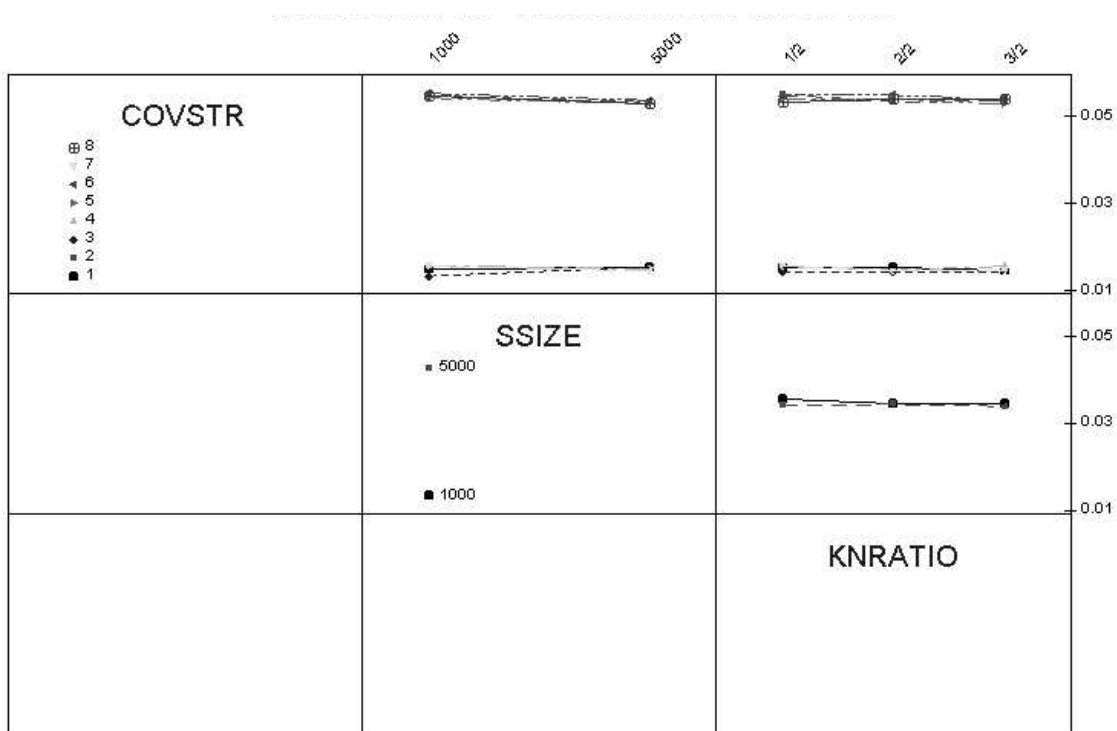


Figure 19: The propensity score weighting, Case 0, interaction plot for *MeanBias*. (Identification of the COVSTR levels given in the text.)

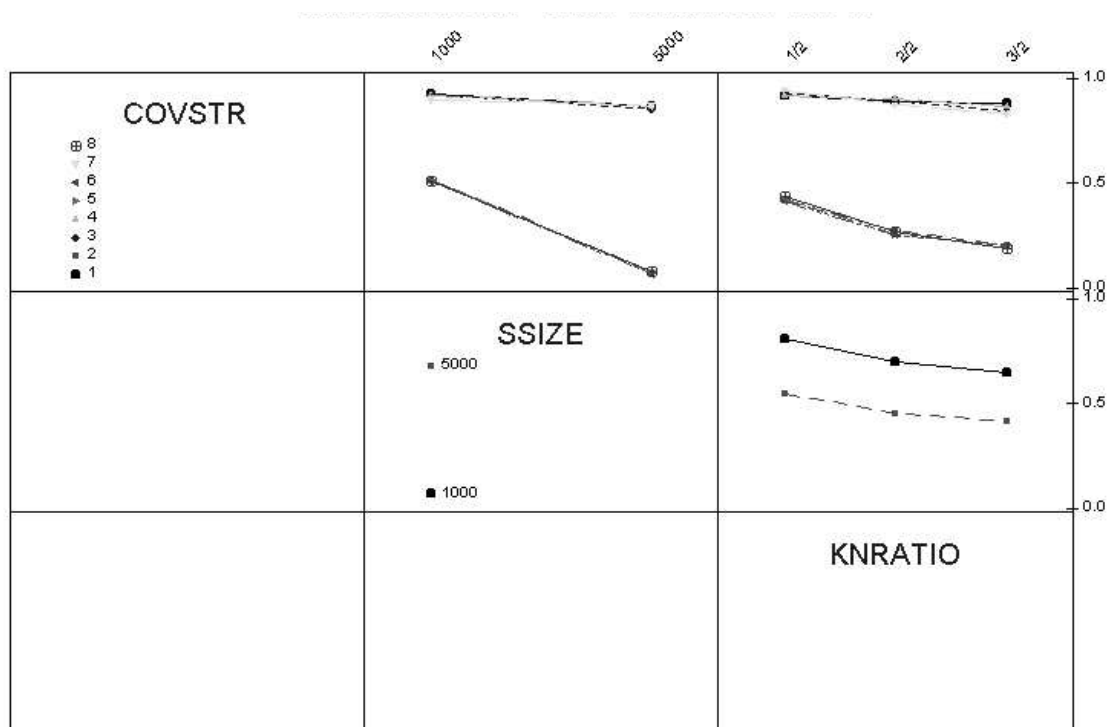


Figure 20: The propensity score weighting, Case 0, interaction plot for *Clevel*. (Identification of the COVSTR levels given in the text.)

It is though interesting to note that the PS-weighting produced point estimates that were approximately as biased as theoretically expected, while MI-adjustment with the propensity score technique produced estimates that were more biased “than necessary”. The means for the estimates adjusted with the PS-weighting were .015 and .054 for the low and high original bias structures respectively, agreeing with the analytic derivation in (Lorenc, 2003a), while the corresponding numbers for the MI-adjustment using the propensity score technique were higher, .021 and .075 respectively.

3.9.2 Case 7: method=prop, all the violations simultaneously present

Before presenting the effects of introduction of the assumption violations, a comment regarding a “nonviolating” difference between observing full and partial auxiliary information. With the other multiple imputation techniques, using all available information (i.e. both X_1 and X_2) was more effective that using only X_2 (Figure 5), primarily because of the bias reduction for the structures 5, 6, and 7. While from Figure 21 (pane COVSTR×OBSERVED) it appears that the opposite was the case with the propensity score as the MI-technique (i.e. that the bias had increased by a change from observing X_2 to observing both X_1 and X_2), the apparent large increase is the result of confounding in this two-way graphical representation of the differential effect of the SITA violations on the covariance structures. In an analysis that removed all the violations (not graphically presented here), the increase in bias due to observing both X_1 and X_2 was negligible.

All three assumption violations, SITAVIO2, OBSERVED=X1, and SITAVIO1, had a

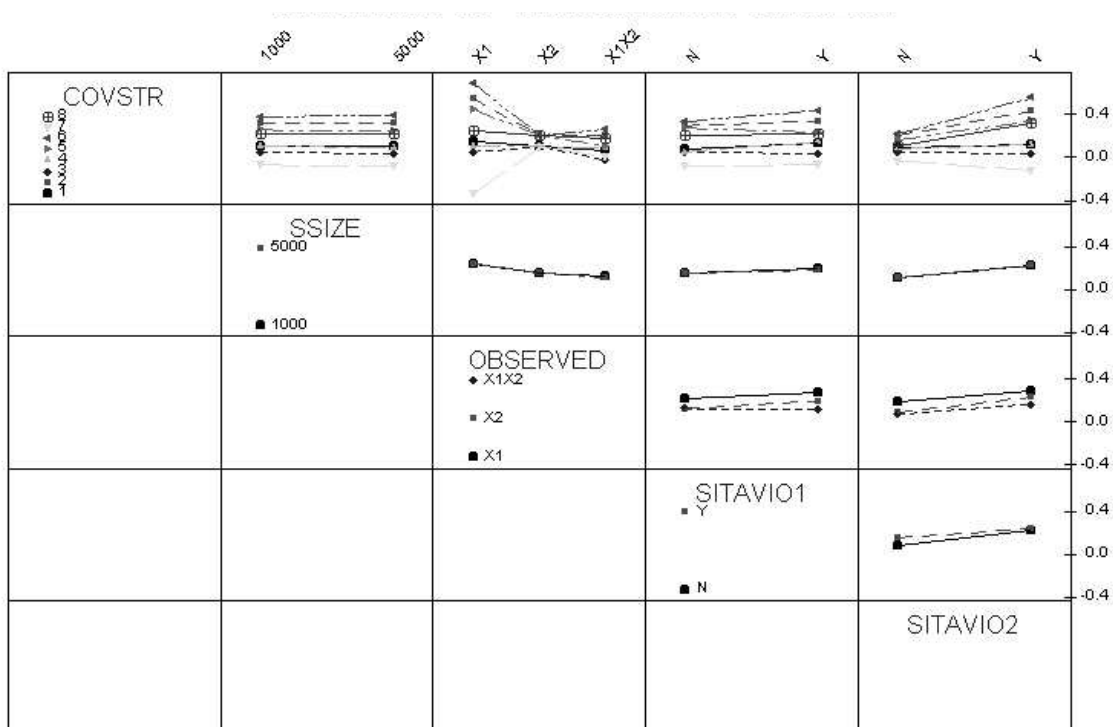


Figure 21: Case 7, only the level METHOD=PROP, interaction plot for *MeanBias*. (Identification of the COVSTR levels given in the text.)

strong impact on the bias of the adjusted point estimator, in the order mentioned (Figure 21 and Table XXXIII in the Appendix). Observing only X_1 created a considerable bias that was differential with respect to covariance structures (pane COVSTR \times OBSERVED = X1 in the figure), SITA violation #2 in general amplified that bias (pane COVSTR \times SITAVIO2), while SITA violation #1 moved some of the adjusted point estimates in the positive direction (pane COVSTR \times SITAVIO1). The only factor not significant in the ANOVA decomposition was SSIZE (Table XXXI in the Appendix).

Finally, as the same data exist for the simulation study of the PS-weighting mentioned previously, they are given here in Figure 22 for the sake of comparison. The same broad description of the effects of the violations is in effect here too: the level OBSERVED = X1 introduced a large bias that was differential with respect to covariance structures, SITAVIO2 amplified the bias of some of the structures, and while SITAVIO1 moved in general the adjusted point estimates in the positive direction. The significant difference was that the absolute level of bias introduced by the violations was much higher with the PS-weighting: the technique proved to be less robust to violation assumptions than the MI-adjustment.

4 Conclusions

The aim of this simulation study was to demonstrate the efficiency of multiple imputation as a bias reducing technique in situations with double samples. After a summary of the main results, that comes first, a discussion of the robustness of the technique is given

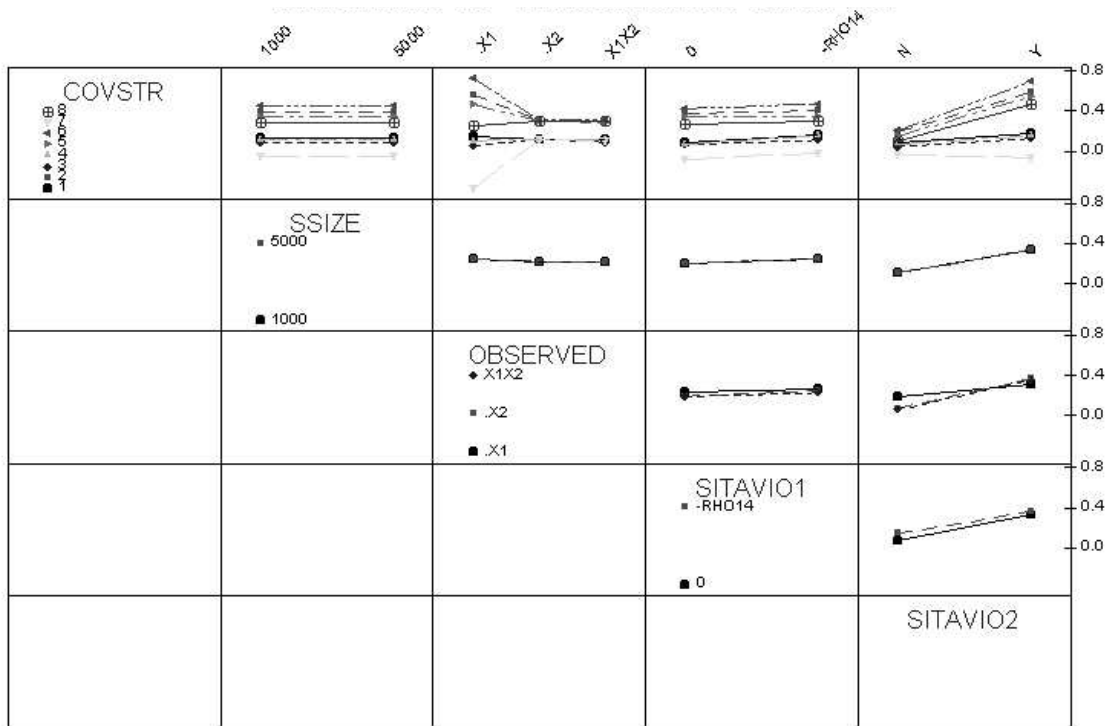


Figure 22: The propensity score weighting, Case 7, interaction plot for *MeanBias*. (Identification of the COVSTR levels given in the text.)

with a view at its practical application.

4.1 Effects of the factors studied

Among the factors conforming to the assumptions, increase in sample size had the expected general effects of producing point estimates with higher accuracy and empirical confidence levels closer to the nominally declared ones.

Ratio of the sample sizes did also have a positive correlation with accuracy and correct interval prediction, which at least partially must be accounted for by the mere increase in the number of observations available for the analysis associated with the higher levels of this factor. In other words, had a change in sample ratios not have changed the combined samples' size, a more proper view on the contribution of this factor would have been gained.

The technique for multiple imputation was the sole factor that consistently proved to be insignificant with respect to both point estimation and confidence level. Presumably, the model used for generating the populations did not let the advantages of the more robust, and more demanding in terms of computer power, MCMC techniques to show up.

Of the covariance structures studied, one turned up to be particularly dangerous for a statistician wishing to correct for the bias of the unadjusted estimate. The structure turns against the statistician when a strong participation variable is disregarded, and consists of a strong correlation between the auxiliary variable and the study variable as well as a strong correlation between the auxiliary variable and the participation variable,

but a weak correlation between the participation variable and the study variable. This latter fact assures that the bias due to the sample's origin in the subset is little, but two forceful adjustments are nevertheless performed, overadjusting the mild bias. The problem with this particular covariance structure disappears when the participation variable is observed. Of course, in practical applications an exact knowledge of the correlation between the participation variable and the study variable does not exist in advance, but reasoning and consulting existing results might give some prior information about these relations.

The factors violating the assumption did all have a significant effect on the simulation statistics. The primary purpose of introducing them was a demonstrative one: particular choices of the levels of the assumption violating factors were too many. Nevertheless, some wider statements may be made, as follows.

4.2 Robustness of the technique

The present study addressed even the issue of robustness to assumption violations of the estimates produced by multiple imputation adjustment. With the population model as used in the present study, to give to some units in the population a zero chance to appear in the restricted sample (e.g. a web sample) did not have any deteriorating effect on the quality of the point and variance estimates—they were on target practically just as with no assumption violations. The reason for this, technically, was that relation between the participation variable and the study variable was linear in both the population and the subset, with these two lines parallel. Whether multiple imputation adjustment would be as robust to this particular violation in practical applications as in the present simulation study will depend on whether the linear and parallel relation holds also for the variables in the real study.

Robustness of the technique to a residual correlation between the study variable and the subset indicator was dependent on the correlation between the study variable and the participation variable. With high correlations, there was a slight downgrading effect, about 5% on the bias reduction and somewhat more, by about .25 on average, on the confidence level produced from the variance estimator. With low correlations, the effect was much harsher, bringing down the bias reduction to a half, and producing confidence intervals of little use. This stresses the importance in real studies of collecting auxiliary information that strongly predicts participation of the units in the subset from which the restricted sample originates. From the present study, it is by far more important to collect this information than other auxiliary information, not related to participation. This position was apparently taken by Terhanian, whose procedure included collecting attitudinal and behavioural auxiliary information in addition to classical demographic variables (Terhanian et al., 2001).

Frail rather than robust was the technique to failure to observe the participation variable. (By assumption, all information relevant for sample assignment needs to be collected). While a pattern was observed regarding differences between the covariance structures in bias reduction caused by this factor, such that in some particular cases the effect was less damaging than in others, the results are not of other than academic interest as the correlation between the study variable and the other variables in the population is not known in advance.

4.3 Further work

Among the possible future work, three topics are brought up here.

In the situation when the assumptions held (Case 0 above), for five of the eight covariance structures the variance estimates were too conservative—the three nonconforming ones being those with almost singular covariance matrices, why say Y could have almost ideally be predicted from the other variables. In practical applications, where singularity need not be the case, the variance estimates using the multiple imputation adjustment would thus be too large and the resulting confidence intervals too wide. Ways of improving the variance estimates may thus be a topic of interest for future work.

Further, while it was the intention of the present study to address the difference between the multiple imputation techniques, this failed due to an inappropriate choice of the model for that end. Such an investigation, considering the pros and cons of the different techniques with respect to the underlying population remains for an eventual future study.

Related to this, the advantage of multiple imputation adjustment over the propensity score weighting, demonstrated in this study, may be accounted for by the simple, practically linear model used for the population. The propensity score is effective (i.e., it reduces most of the bias) even in much more complex variables structures; some of the used techniques for multiple imputation (e.g. the MCMC techniques) are that, too. It might be thus of interest to investigate the circumstances in which the propensity score weighting eventually would perform better than multiple imputation adjustment.

4.4 Summary

The present simulation study showed that multiple imputation adjustment in a double samples setting came close to being perfect both for the point estimates and the estimates of their variance. In order for this to be so, some assumptions needed to be fulfilled, but these were not stronger than for any weighting technique in the double samples setup (e.g., for the propensity score weighting). The study also demonstrated the impact of a number of factors on the efficiency of the technique, some of the factors related to the performance of the technique in general, and some related to violations of the assumptions. Of great importance turned out to be collection of information predictive of units' participation in the special subset from which the non-random sample comes, much greater than "usual" auxiliary information. With this information carefully collected, multiple imputation may in many cases give point estimates with most of the bias removed and with confidence levels not too far from their nominal levels.

Acknowledgment

Support for this study from the Bank of Sweden Tercentenary Foundation, Grant no. 2000-5063, is gratefully acknowledged.

References

- [1] Cochran, W.G. (1968). “The effectiveness of adjustment by subclassification in removing bias in observational studies”. *Biometrics*, 24:205-13.
- [2] Cochran, W.G. and Rubin, D.B. (1973). “Controlling bias in observational studies: a review”. *Sankya*, ser. A, 35:417-46.
- [3] Lorenc, B. (2003a). “Effectiveness of weighting by stratification on the propensity score using double samples”. Research report 2003:10. Department of statistics, Stockholm university.
- [4] Lorenc, B. (2003b). “Propensity score weighting with double samples: a simulation study”. Research report 2003:11. Department of statistics, Stockholm university.
- [5] Rubin, D.B. (1987). *Multiple Imputation for Nonresponse in Surveys*. New York, NY: Wiley.
- [6] Rosenbaum, P.R. and Rubin, D.B. (1983). “The central role of the propensity score in observational studies for causal effects”. *Biometrika*, 70:41-55.
- [7] Rosenbaum, P.R. and Rubin, D.B. (1984). “Reducing bias in observational studies using subclassification on the propensity score.” *Journal of the American Statistical Association*, 79:516-24.
- [8] SAS Institute Inc. (2001). *SAS/STAT Software: Changes and Enhancements, Release 8.2*. Cary, NC: SAS Institute Inc.
- [9] Terhanian, G., Marcus, S., Bremer, J., and Smith, R. (2001). “Reducing error associated with non-probability sampling through propensity scores: evidence from election 2000”. *Joint Statistical Meeting 2001*, August 5-9, 2001, Atlanta, GA, USA.