 Anna-Simone Josefine Frank*, David S. Matteson, Hiroko K. Solvang, Angela Lupattelli and Hedvig Nordeng

Extending balance assessment for the generalized propensity score under multiple imputation

<https://doi.org/10.1515/em-2019-0003>

Received January 30, 2019; accepted March 31, 2020

Abstract: This manuscript extends the definition of the Absolute Standardized Mean Difference (ASMD) for binary exposure ($M = 2$) to cases for $M > 2$ on multiple imputed data sets. The Maximal Maximized Standardized Difference (MMSD) and the Maximal Averaged Standardized Difference (MASD) were proposed. For different percentages, missing data were introduced in covariates in the simulated data based on the missing at random (MAR) assumption. We then investigate the performance of these two metric definition using simulated data of full and imputed data sets. The performance of the MASD and the MMSD were validated by relating the balance metrics to estimation bias. The results show that there is an association between the balance metrics and bias. The proposed balance diagnostics seem therefore appropriate to assess balance for the generalized propensity score (GPS) under multiple imputation.

Keywords: balance diagnostics; generalized propensity score; missing data; Monte Carlo simulations; multiple treatment exposure.

Introduction

It is impossible in observational studies to control how subjects are assigned into treatment groups, and this may potentially result in biased effect estimates caused by confounding (Hernán et al. 2008; Rubin 2004). The application of weights derived from the propensity score (PS) (the conditional probability of receiving treatment given observed covariates) aims to balance characteristics between treatment groups (Rubin 2004). When this is achieved, the result is reduced bias effect in the estimates (Austin 2011). For binary exposures, a common approach to checking whether balance between treatment and control group has been achieved is to calculate the Absolute Standardized Mean Difference (ASMD, d). For a continuous covariate x , the ASMD is defined as

$$d = \frac{|\bar{x}_{treatment} - \bar{x}_{control}|}{\sqrt{\frac{s_{treatment}^2 + s_{control}^2}{2}}},$$

where $\bar{x}_{treatment}$ and $\bar{x}_{control}$ denote respectively, the sample mean of the covariate in the treated and control groups, while $s_{treatment}^2$ and $s_{control}^2$ represent the sample variance in the respective treatment groups (Austin

*Corresponding author: Anna-Simone Josefine Frank, Department of Informatics, Computational Biology Unit (CBU), University of Bergen, Bergen, Hordaland, Norway, E-mail: Anna-Simone.Frank@uib.no; asfrank88@gmail.com. <https://orcid.org/0000-0002-3728-3476>

David S. Matteson: Department of Statistical Science, Cornell University, Ithaca, New York, USA; Department of Biological Statistics and Computational Biology, Cornell University, Ithaca, New York, USA

Hiroko K. Solvang: Havforskningstutttet, Bergen, Norway

Angela Lupattelli: Department of Pharmacy, Universitetet i Oslo, Oslo, Norway

Hedvig Nordeng: Department of Pharmacy, Universitetet i Oslo, Oslo, Norway; Nasjonalt folkehelseinstitutt, Department of Child Health and Development, Oslo, Norway

2018; McCaffrey et al. 2013). When defining d for categorical variables, mean values (\bar{x}) are replaced with proportions and variances (s^2) become functions of proportions (Austin 2018).

If balance has been achieved, the ASMD value should be below a pre-defined threshold after PS analysis has been performed (Austin 2018; McCaffrey et al. 2013).

There is however inconsistency in balance thresholds across the literature, where values vary between 0.1 and 0.25 for d (Austin 2011; Nguyen et al. 2017; Stuart, Lee and Leacy 2013). Given these differences, a recent study investigated the effectiveness of pre-defined thresholds for PS matching with binary exposure (Nguyen et al. 2017). The authors concluded that if balance was below 0.1 on most variables, and if variables with balance above 0.1 are included in the outcome model as adjustment variables, bias was small. Yet, a specific threshold for the generalized propensity score (GPS) has not been tested.

The GPS is the generalization of the PS for binary exposure. The GPS makes it possible to estimate the effect of multiple treatment exposure on the outcome (Imai and Van Dyk 2004; Lechner 2001). The GPS has been applied in various studies on observational data, for example (Feng et al. 2012; Jiang and Foster 2013; Spreeuwenberg et al. 2010; Sugihara 2010).


A common problem of observational data is missing information due to non-response, especially where data collection is performed via questionnaires (Pandis 2014).

Multiple imputation techniques are often applied to fill the information gaps (Miri et al. 2016). For binary exposure, several studies have combined imputation techniques with PS analysis (Doidge 2018; Eulenburg et al. 2016; Hayes and Groner 2008; Hill 2004; Hsu and Yu 2018; Kupzyk and Beal 2017; Lavori, Dawson, and Shera 1995; Mitra and Reiter 2016; Rosenbaum and Rubin 1984; Qu and Lipkovich 2009). A systematic review by Malla et al. (2018) summarizes how missing data are combined with PS analysis on actual patient data. The authors showed that the majority of reviewed articles performed PS analysis in combination with complete case data (Malla et al. 2018).

Unlike for PS analysis, only one recent study, by De Vries, Van Smeden, and Groenwold (2018), considered the combination of missing data and the GPS (De Vries, Van Smeden, and Groenwold 2018). The authors found that multiple imputations of data, followed by PS estimation using Classification and Regression Trees (CART) resulted in least biased estimates (De Vries, Van Smeden, and Groenwold 2018). However, no previous study has assessed balance for the GPS under multiple imputations, partly due to the computational burden involved (De Vries, Van Smeden, and Groenwold 2018).

This paper proposes an approach for balance assessment of GPS under multiple imputation. It uses simulated data to evaluate the proposed diagnostics. The article is organized in the following way. Section 2, the methodology section, defines the quantitative concepts applied in this study, such as the GPS, its estimation, a multiple treatment definition and multiple imputation. Existing balance diagnostics are reviewed for binary PS models and finally extended to the multiple treatment case under multiple imputations. Section 3, the simulation study section, describes the data-generation process and details about the implementation. The results of simulated data example are presented in Section 4, followed by the discussion of the results in Section 5 and finally the conclusions in Section 6.

Methodology

 In this section we describe, the methodological concepts and ideas that are later applied to simulated data examples.

Generalized PS

The aim of PS analysis is to reduce the dimension of observed pre-treatment variables X and bias, due to confounding, by re-weighting them.

For multiple treatments, this can be achieved by using the GPS (Imbens 2000) In the present study, we will estimate GPS using generalized boosted models (GBM), and combine it with inverse probability of treatment weights (IPTW) to obtain the average treatment effect (ATE) (McCaffrey, Ridgeway, and Morral 2004; McCaffrey et al. 2013).

We applied IPTW based on the results by Nian et al. 2019, who showed that this approach had preferred performance compared to other approaches, such as matching, stratification or GPS adjustment. In addition PS methods may differ in the population where an overall treatment estimate shall be calculated Kurth et al. 2005. The results in Kurth et al. (2005) have shown that IPTW is well suited to calculate the treatment effect of the total population.

Following, Imbens (2000) and Feng et al. (2012), for the subject index $i = 1, \dots, N$, let $\mathcal{T} = \{1, 2, \dots, M\}$ denote the set of M multiple treatments, and $Y_i(m)$ denote the potential outcome of subject i , if subject i has been assigned to treatment $m \in \mathcal{T}$ (Imbens 2000, Feng et al. 2012). Let T_i be the treatment that subject i received and the indicator that subject i receives treatment m is

$$I_m(T_i) = I(T_i = m). \quad (1)$$

Definition 2.1. (Generalized Propensity Scores (GPS) (Imbens 2000; McCaffrey et al. 2013)) Let X_i be the set of observed pre-treatment variables of subject i . The GPS $r(m, X_i)$ is the conditional probability of receiving a particular level of the treatment m given X_i :

$$r(m, X_i) = \Pr(T_i = m|X_i) = \mathbb{E}[I_m(T_i)|X_i]. \quad (2)$$

Hence, for the set \mathcal{T} of M different treatment groups, we obtain M GPSs, for all subjects i .

Given $r(m, X_i)$, the empirical expected outcome $\widehat{\mathbb{E}}(Y(m))$ is estimated by the weighted mean (McCaffrey et al. 2013), i. e.

$$\widehat{\mathbb{E}}(Y(m)) = \frac{\sum_{i=1}^N I_m(T_i)Y_i(m)w_i(m)}{\sum_{i=1}^N I_m(T_i)w_i(m)}, \quad (3)$$

Where

$$T_i = m, \quad w_i(m) = \frac{1}{r(m, X_i)}. \quad (4)$$

Then, the ATE is estimated, comparing treatment $t \in \mathcal{T}$ vs. treatment $\ell \in \mathcal{T}$ ($t \neq \ell$) across the subjects

$$\widehat{\text{ATE}}_{t\ell} = \widehat{\mathbb{E}}[Y(t)] - \widehat{\mathbb{E}}[Y(\ell)], \quad (5)$$

under the condition that treatment (T_i) is independent of the outcome $Y_i(T_i)$, and that the compared groups are representatives of the population (Feng et al. 2012). The latter condition does not hold in general, but can be achieved by the key assumption that treatment assignment is *weakly unconfounded* given observed covariates X . For multi-valued treatments, without missing data, the *weak unconfoundedness* assumption was defined by Imbens (2000), see Definition S1.1 in the Supplementary material section S1. This definition means that treatment and potential outcome are independent given the observed covariate (Imbens 2000). It was shown in Leyrat et al. 2019a, that this assumption holds for a binary exposure under multiple imputation. However, to be able to define the ATE for the present situation, the following assumptions have to hold for the GPS after multiple imputation:

For a proof of these assumptions, see the Supplementary material sub-section S1.1.

Previous studies estimated GPS with multinomial and ordinal logistic regression models (Bray et al. 2018; Feng et al. 2012). However, these parametric approaches have been shown to lead to less robust ATE estimates than non-parametric approaches, such as GBM, which we describe briefly below (McCaffrey, Ridgeway, and Morral 2004; McCaffrey et al. 2013).

Generalized boosted models

Firstly, we describe how GBM are applied to estimate the PS for binary treatment, i. e., $\mathcal{T} \in \{0, 1\}$ (McCaffrey, Ridgeway, and Morral 2004; McCaffrey et al. 2013). Then the GBM algorithm is generalized to the multiple treatment case analogous to McCaffrey, Ridgeway, and Morral (2004); McCaffrey et al. (2013).

Let X be the set of observed pre-treatment variables and let X_i be the set of observed pre-treatment variables for subject i . For the binary treatment case, the treatment indicators for subject i in Eq. (1) is given by $T_i = 1$, simplified with $I_1(T_i) = 1 - I_0(T_i)$ and $r(1, X_i) = 1 - r(0, X_i)$. Instead of directly estimating the PSs, the algorithm finds the maximum-likelihood estimate of the function $g(X)$, which is the log-odds of treatment assignment, i. e., $g(X) = \log(r(1, X)/(1 - r(1, X)))$. Therefore, the GBM algorithm iteratively adds regression trees together to fit a non-linear logistic regression model to treatment indicator $I_1(T_i)$ (McCaffrey, Ridgeway, and Morral 2004). Initially, the algorithm sets $g_0(X) = \log(\bar{I}/(1 - \bar{I}))$, where \bar{I} is the average treatment assignment indicator for the whole sample. Then to improve the PS fit to the data, at each new iteration j ($j = 1, \dots, J$), a new regression tree $h(X)$ is added to the current model $g_{j-1}(X)$, if it is the best fit to the residuals $I_1(T_i) - g_{j-1}(X_i)$ and provides the greatest increase in the log-likelihood for the data. When the regression trees are combined a shrinkage coefficient α is introduced to improve smoothness of the resulting piecewise constant model. In order to avoid overfitting of the data, GBM selects a number of trees in order to minimize imbalance on pre-treatment covariates across “treatment” and “control” groups. More details about the boosting algorithm can be found in McCaffrey, Ridgeway, and Morral (2004).

GBM extension to more than two treatment groups ($M > 2$)

The approach in Section 2.1.1 can be extended to more than two treatment groups in the following way (McCaffrey, Ridgeway, and Morral 2004; McCaffrey et al. 2013): Firstly, we define indicator functions for each of the M treatment groups as in Eq. (1). The GBM algorithm is then applied respectively to all M treatment indicator functions $I_m(T_i)$, resulting in M PSs $r(m, X_i) = Pr(T_i = m|X_i)$. The respective IPTW $w_i(m) = \frac{1}{r(m, X_i)}$ are then applied in Eq. (3).

As mentioned previously, missing data are common in observational data. Therefore, techniques, which fill information gaps, have been invented, such as Multiple imputation by chained equations (MICE), which we describe below (Van Buuren and Groothuis-Oudshoorn 2011).

Multiple imputation by chained equations

Let \mathcal{S} be a data set and let \mathcal{Z} be the subset for the elements that are missing in \mathcal{S} . MICE consists of three main steps (Van Buuren and Groothuis-Oudshoorn 2011): Firstly, there is imputation of the missing data component \mathcal{Z} , using chained equations. After an initial random filling of all missing data, an MCMC algorithm draws iteratively from the conditional distributions based on a collection of observed and missing variables (Chen and Ip 2015; Van Buuren and Groothuis-Oudshoorn 2011) This results in Q multiple imputed data sets $\mathcal{S}^{(q)}$, where $q = 1, \dots, Q$, that are identical on each observed element, but differ on the initially missing entries. On each of the q imputed data sets, PS analysis is performed, and leads to Q effect estimates of the PS analysis.

Finally, all Q effect estimates and their variance are pooled into one estimate. Following Murray (2018), as many variables as possible, including the outcome variable, should be incorporated into the imputation model (Leyrat et al. 2019a; Murray 2018).

Initially, the literature recommended that $5 \leq Q \leq 10$ imputed sets are sufficient (Azur et al. 2011). However more recent developments showed that in order to also detect small effect sizes a minimum of $Q = 40$ imputed sets are needed (this however depends on percentage of missing data), see Graham, Olchowski, and Gilreath (2007).

There are three approaches to combine multiple imputation and PS analysis (Leyrat et al. 2019a): After multiple imputation, one approach averages the estimated PSs to generate one outcome analysis (Leyrat et al. 2019a; Mitra and Reiter 2016). Another approach, calculates the PS based on the pooled covariates over all imputed data sets. A third approach averages the Q treatment effects, which were estimated based on the PS on each imputed data set. Leyrat et al. 2019a compared all three approaches and concluded that the third one leads to least biased results. This approach was therefore applied in the below illustrative example (see Section 3).

Next we review and extend the definition of balance.

Balance assessment

In order to assess if the PS analysis was able to balance treatment groups across observed pre-treatment variables X , balance is calculated. Previous studies extended balance diagnostics to GPS for continuous and multiple treatment for complete data (Austin 2018; Bray et al. 2018; Fong, Hazlett, and Imai 2018; McCaffrey et al. 2013; Zhu, Coffman, and Ghosh 2015). We first review the definition of balance for binary and multiple treatment groups and finally extend this definition for the case under multiple imputation. In addition to previous notations, let $k = 1, \dots, K$ denote the covariate index and let $C = \binom{M}{2}$ denote the total number of pairwise comparisons of M treatments and let c be an index over comparison pairs, i. e., $c = 1, \dots, C$.

Balance assessment without multiple imputation

Definition 2.2. (Absolute Standardized Mean Difference (ASMD) (McCaffrey et al. 2013)) For each covariate k and binary treatment m , the ASMD equals the absolute value of the difference between the weighted mean of the covariate in the treatment group ($m = 1$) minus the weighted mean of the covariate in the control group ($m = 0$), divided by the unweighted standard deviation of the pooled population for the ATE (see Eq. (5)). Given covariate k ,

$$ASMD_k = \frac{|\bar{X}_{k,1} - \bar{X}_{k,0}|}{\hat{\sigma}_k},$$

where $\bar{X}_{k,m}$ is the weighted mean for covariate k and treatment ($m = 1$) or control ($m = 0$) group, and $\hat{\sigma}_k$ is the unweighted averaged (pooled) within standard deviation for all treatment groups.

For GPS, Burgette, Griffin, and McCaffrey (2017) suggest that, balance shall be assessed, via the ASMD, for all pairwise comparisons when determining the ATE.

Definition 2.3 extends Definition 2.2 to M treatment groups $m = 1, \dots, M$, which will result in C -pairwise comparisons (Burgette, Griffin, and McCaffrey 2017; McCaffrey, Ridgeway, and Morral 2004).

Definition 2.3. (Multiple ASMD (McCaffrey et al. 2013)) For M multiple treatment groups, and given the set of K covariates and treatment pairs (t, ℓ) , $t \neq \ell$,

$$ASMD_{k,c} = \frac{|\bar{X}_{k,t} - \bar{X}_{k,\ell}|}{\hat{\sigma}_k}$$

$\bar{X}_{k,t}$ denotes the weighted mean for covariate k for treatment group t , while $\bar{X}_{k,\ell}$ is the weighted mean for covariate k for treatment group ℓ . The denominator $\hat{\sigma}_k$ is the unweighted averaged (pooled) standard deviation of all treatment groups (same as in Definition 2.2).

A similar definition holds for categorical variables, where the mean is replaced by weighted proportions.

Except for McCaffrey et al. (2013), two recent articles (Li and Li 2019; Yang et al. 2016) included discussions on balance assessments for more than two treatment arms: The Multiple ASMD is a special case of the pairwise absolute standardized differences (ASD) defined by Li and Li (2019), with $w_i(m) = \frac{1}{r(m, X_i)}$ and tilting function $h(X) = 1$. In addition, from McCaffrey et al. (2013), Li and Li (2019) extended the population standardized difference (PDS) for varying weight functions ($w_i(m)$). Instead of comparing weighted covariate means pairwise between treatment groups, the PDS compares weighted covariate means between the treatment group and the target population. The balance metric in Yang et al. (2016) are defined for matching and stratification, therefore $w_i(m) = 1$. Another major difference is that covariate means for each treatment group m are not compared pairwise, but with the covariate mean of all other treatment groups combined, except m (i. e., $m^c = T$). Finally, Yang et al. (2016) also proposed a metric that allows the assessment of balance in covariate distributions, where the GPS are considered instead of the covariates directly.

In the next section, Definition 2.3 is extended to multi-treatment exposure under multiple imputation.

Balance assessment for GPS under multiple imputation

Definition 2.4. (Maximal Maximized Standardized Difference (MMSD)) For each covariate k , the MMSD is defined by,

$$MMSD_k = \max_c \max_{q=1, \dots, Q} ASMD_{q,k,c}, \quad (6)$$

where $ASMD_{q,k,c}$ refers to the ASMD for covariate k and pairwise comparison c on the multiple imputed set q .

Definition 2.5. (Maximal Averaged Standardized Difference (MASD)) Similarly, for each covariate k , the MASD is defined by,

$$MASD_k = \max_c \frac{1}{Q} \sum_{q=1}^Q ASMD_{q,k,c}, \quad (7)$$

where $ASMD_{q,k,c}$ is defined as above in Definition 2.4.

It is assumed that balance is obtained when, after weighting, the Maximized Standardized Difference (MMSD) or the Maximal Averaged Standardized Difference (MASD) are below a pre-defined threshold. If all covariates are balanced under MMSD then they are automatically balanced under MASD. The reverse argument does not hold.

Pre-defined thresholds have been reported to vary between 0.1 and 0.25 for balance assessment in the literature (Stuart, Lee, and Leacy 2013). For validation of the balance diagnostic, we follow a similar approach as Franklin et al. (2014), which creates a summary score of balance diagnostics over all covariates. Therefore, for the balance summary score, we decided to consider a threshold of 0.1 appropriate.

It is also important to assess, if the *positivity assumption* is fulfilled. The *positivity assumption* states that each treatment group has a positive probability of receiving the treatment (Austin and Stuart 2015; McCaffrey, Ridgeway, and Morral 2004; McCaffrey et al. 2013). Across all imputed data sets, we therefore select the minimum and maximal weights, in order to make sure that there are no extreme values. Occurrence of extreme

values could be indicative of a violation of the *positivity assumption*. Hence in cases they occur, the weights are truncated or stabilized weights are used as alternatives (Austin and Stuart 2015).

Besides checking that the positivity assumption is fulfilled, we also checked that the selected GPS model was well specified by quantifying the standard deviation of the weights (Austin and Stuart 2015; McCaffrey, Ridgeway, and Morral 2004; McCaffrey et al. 2013).

Simulation study

The aim of this simulation study is to evaluate balance metrics, MASD and MMSD, in their ability to detect imbalances in confounders after IPTW based on the GPS on multiple imputed data sets. We will compare the imbalance after weighting with the biased estimates of a continuous outcome.

Motivation of simulated data: clinical problem

The simulation example for this study is motivated by observational studies investigating the associations between prenatal exposure to medications on pregnancy outcomes (Lupattelli et al 2017; Nezvalová-Henriksen et al. 2016; Nordeng et al. 2012). Such studies are necessary as clinical studies are rarely ethical among pregnant women (Blehar et al. 2013). As many diseases have three or more therapeutic options, we decided to focus on a multi-group comparison of three medication alternatives. For example, hyperthyroidism during pregnancy can be treated with methimazole/carbimazole (MMI/CMZ), propylthouracil (PTU) or left untreated (Alexander et al. 2017; Moleti et al. 2019). Although, we motivated the generation of simulation data on the above mentioned clinical example, the main focus of this manuscript is of primary methodological nature and allows no clinical implications and conclusions. The manuscript uses simulated data and has therefore no ethical issues.

Data generation

Data sets of sample size $n = 1000$ were generated, in order to mirror an observational study comparing three treatment options ($M = 3$), $\mathbf{T} = (T_1, T_2, T_3)$, with T_2 as reference treatment, on a fully observed continuous outcome Y (birth weight in gram) with five measured covariates $\mathbf{X} = (X_1, X_2, X_3, X_4, X_5)$. The covariates represent, respectively Body Mass Index (BMI), maternal educational level and age, marital status and parity. The covariates X_1, X_2 and X_3 were categorical, and X_4 and X_5 binary. While the variables X_1 and X_2 were partially observed, the other covariates were fully observed. Simulation details on covariates are outlined in Table 1. The data generation process is summarized in a directed acyclic graph (DAG) (see Figure 1).

Figure 1 summarized the data generation in a DAG. T represents the treatment and Y the outcome, while the blue arrow from T to Y represents the treatment effect. X_{obs} and X_{miss} represent, respectively, the completely and partially observed covariates, and R the missing data indicator.

Treatment assignment

For each subject i , the probability of treatment assignment was determined from a multinomial logit function for $M > 2$ based on a linear combination of all five covariates \mathbf{X} and with treatment reference $T_2 = 2$ via (Menard 2002, Borooah 2002):

Table 1: Simulation details for covariates.^a

Covariates	Categories	Properties, %
X_1	BMI ^b ≤ 18 , $19 - 29$, ≥ 30	2.0, 82.0, 16%
X_2	Education in years <9 , $9-12$, $13-16$, and >16	1.4, 25.2, 42.5 and 30.9%
X_3	Age in years ≤ 24 , $25-29$, $30-34$ and ≥ 35	45.0, 17.0, 30.0 and 8.0%
X_4	Married/Cohabiting vs. Other	94.6 vs. 5.4%
X_5	Primiparity vs. Multiparity	42.3 vs. 57.7%

^a Numbers were inspired by real world data in Frank (2019), Frank et al. (2018) ^b BMI categories (kg/m^2), $19 \leq \text{BMI} \leq 24$ (normal weight) and $25 \leq \text{BMI} < 29$ (overweight) are merged into one category in this dataset Abbreviations: μ , mean, σ , standard deviation.

$$\ln\left(\frac{\mathbb{P}(T = m)}{\mathbb{P}(T = 2)}\right) = \alpha_m + \sum_{c=1}^C \beta_{m,c} X_{i,c} = Z_{m,i}$$

With $M > 2$ treatment classes, there are $M - 1$ predicted log-odds, one for each treatment category, with respect to the reference category. For our simulation example with $M = 3$, $Z_{1,i}$ and $Z_{3,i}$ were defined as follows:

$$Z_{1,i} = \log(0.2) + \log(1.4)X_1 + \log(2.0)X_2 + \log(1.25)X_3 + \log(1.3)X_4 + \log(0.2)X_5,$$

$$Z_{3,i} = \log(0.5) + \log(1.4)X_1 + \log(1.2)X_2 + \log(1.25)X_3 + \log(1.3)X_4 + \log(0.8)X_5$$

The coefficients in models $Z_{1,i}$ and $Z_{3,i}$ were chosen empirically, such that $p_m > 0$, $m = 2, \dots, M$ and $p_1 > 0$. Then, for $m \in \{1, 3\}$, treatment assignment probabilities were calculated the following way:

$$p_m := \mathbb{P}(T = m) = \frac{\exp(Z_{m,i})}{1 + \sum_{m \in \{1,3\}} \exp(Z_{m,i})}$$

The reference group probability was calculated, as

$$p_2 := \mathbb{P}(T = 2) = \frac{1}{1 + \sum_{m \in \{1,3\}} \exp(Z_{m,i})} = 1 - \sum_{m \in \{1,3\}} p_m,$$

$$\text{since } p_2 + \sum_{m \in \{1,3\}} p_m = 1.$$

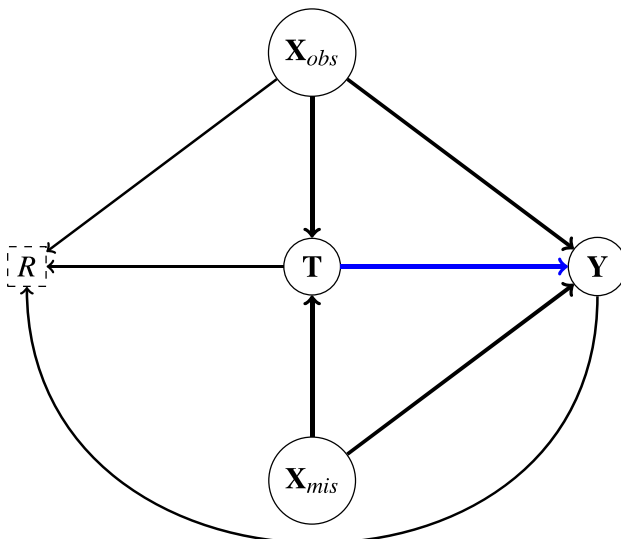


Figure 1: Diagram of data generation.

Treatment was assigned to all subjects by simulating random categorical treatment groups using the probabilities p_m , $m \in \{1, 2, 3\}$.

Continuous outcome

The continuous outcome was simulated via a linear model based on the five covariates and the treatment,

$$Y_i = \beta_T T_i + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \varepsilon \quad (8)$$

with $\varepsilon \sim \mathcal{N}(\mu = 3000, \sigma = 250)$, coefficients $\beta = (\log(2.4), \log(4), \log(1.25), \log(1.3), \log(1.1))$ and with treatment effect of $\beta_T = 200$. This means that we modeled a linear mean effect of treatment. The above models define the full data situation. Next we will generate missing data with difference missing percentages in covariates the X_1 and X_2 .

Missing data mechanism

In this simulation study, we assume that data are Missing At Random (MAR). Hence, the missing information in the variables, X_1 and X_2 , depends on the fully observed covariates X_3 and X_4 (but not X_5 , the treatment assignment variable T , as well as the outcome Y). To introduce missing information, we defined missing indicators R , respectively, for X_1 and X_2 via logistic models:

$$\log it(p(R_1 = 0|T, X_3, Y)) = \gamma_0 + \gamma_1 T + \gamma_2 X_3 + \gamma_4 Y \quad (9)$$

$$\log it(p(R_2 = 0|T, X_3, X_4, Y)) = \eta_0 + \eta_1 T + \eta_2 X_3 + \eta_3 X_4 + \eta_4 Y \quad (10)$$

The values were set to NA when the missing indicators were equal to zero. Parameters describing the missing indicators are presented in Table 2. We varied several factors *at-a-time*, in order to generate several missing percentages.

Missing data percentages of 14.9, 31.8 and 57.5% are well within empirically observed values, which vary between 2 and 65% (Dong and Peng 2013; Karahalios et al. 2012; Marston et al. 2010). Though the 83.1% level of missing data falls outside the empirical range, it was chosen to test ability of the methodology to accommodate the extreme values, because the literature documents evaluation of imputation techniques for 80% missing information (Lee and Huber 2011). For missing percentages larger than 10%, it is assumed that the analysis is potentially biased and imputation analysis is therefore recommended (Dong and Peng 2013).

Estimates

We estimated the mean treatment effects $\hat{\beta}$, intercepts $\hat{\beta}_0$ and approximate 95% confidence intervals (CIs) on the full data set, and the imputed data sets. Over the Q imputed data sets, the effect estimates and CIs were averaged applying Rubin's rule (Leyrat et al. 2019a; Murray 2018).

Methods

For all missing scenarios, we imputed $Q = 10$ data sets via chained equations (using R package "MICE" from Van Buuren and Groothuis-Oudshoorn (2011)), due to the computational burden involved when calculating the GPS (Azur et al. 2011). The imputation model included the outcome, all covariates and the treatment (Frank 2019; Murray 2018,). For each scenario, we then estimated the

Table 2: Missing percentages and parameter values for Eqs. (9) and (10).

Missing, %	Y_0, η_0	Y_1, η_1	Y_2, η_2	Y_3, η_3	Y_4, η_4
14.9	$\log(0.2), \log(0.3)$	$\log(1.75), \log(1.5)$	$\log(4.0), \log(2.0)$	–, $\log(0.8)$	$\log(1.0), \log(1.001)$
31.8	$\log(0.2), \log(0.3)$	$\log(1.25), \log(1.0)$	$\log(3.0), \log(2.0)$	–, $\log(0.8)$	$\log(1.0), \log(1.001)$
57.5	$\log(0.2), \log(0.2)$	$\log(1.12), \log(1.0)$	$\log(2.0), \log(1.1)$	–, $\log(0.5)$	$\log(1.0), \log(1.001)$
83.1	$\log(0.2), \log(0.2)$	$\log(1.12), \log(1.0)$	$\log(1.1), \log(0.9)$	–, $\log(0.5)$	$\log(1.0), \log(1.001)$

GPS based on all five covariates, and calculated the MMSD and MASD for each covariate separately, before and after IPTW. Then similar to the approach in Franklin et al. (2014), we created balance summary scores that average the balance metrics over all covariates. To evaluate whether the balance metrics describe appropriately imbalance in covariates after IPTW, we plotted the estimated bias vs. the balance summary score, after IPTW. Balance assessment was based on the function “bal.table” in R package “TWANG” (Burgette, Griffin, and McCaffrey 2017; Ridgeway et al. 2017). The GPS and IPTW were calculated based on GBM using the R function “mnps” in the R package “TWANG” (Burgette, Griffin, and McCaffrey 2017; Ridgeway et al. 2017). With help of the “SURVEY” R package, IPTW was performed on each imputation set by using the “imputationList”-tool from the package “MITOOLS” (Lumley 2018, 2015). The pooled effect estimates and CIs across the multiple imputed data sets were obtained by applying Rubin’s rule via the function “MIcombine” from the “MITOOLS” R package (Lumley 2015).

Performance measures

For each scenario, and the treatments, $T_{1,3}$, we estimated bias, coverage, and relative percent error in the model standard error (RelError), together with the Monte Carlo standard error (MCSE) for $\hat{\beta}$. These performance measures were calculated with help of the R package “rsimsum” (Gasparini and Lang 2018). Based on conservative estimates from an initial small simulation run, we assumed that $SD(\hat{\theta}) \leq 1676$, and that the MCSE of the bias should be lower than 52 g (Morris, White, and Crowther 2019). Using equations for the MCSE of the bias and coverage, as well as the maximized MCSE of the coverage, as provided by Morris, White, and Crowther (2019), we are required to simulate $n_{sim} = 1053$ repetitive runs (see Figure S1 in Supplementary Material Section S2). The MCSE of the bias estimate of 52 g is very conservative, as this was based on initial runs. It is therefore likely that the required n_{sim} value may be higher in reality. Figure S1 was created in MATLAB (Version 9.4.0.813654 (R2018a)).

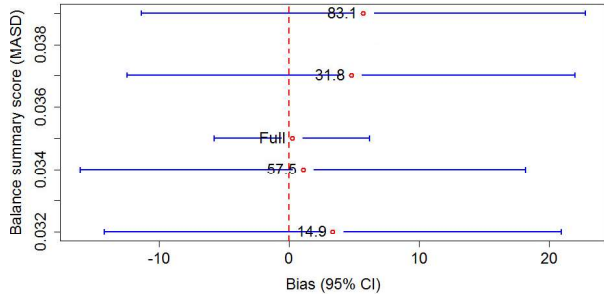
Given the computational burden involved in calculating the GPS (De Vries, Van Smeden, and Groenwold 2018), we were unable to perform analysis over all required Monte Carlo simulations ($n_{sim} = 1053$). We therefore iteratively reduced the number of n_{sim} , and used the number of simulated sets for which we had enough memory capacity (that was the case for $n_{sim} = 10$). These 10 sets were randomly selected out of the 1,053 simulated data sets and statistical analysis was performed on them. Since, for each of the 10 selected data sets, we imputed 10 data sets for the missing data cases, we analyzed the full data set over $n_{sim} = 100$ Monte Carlo runs.

The results of the analysis and balance assessment are presented next.

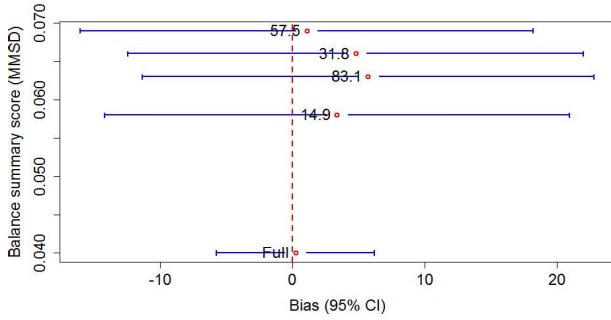
Results of simulation study

In the panels a–d, each line represents a missing data level setting. The line indicated with label “Full” represents the results obtained from the full data set, before missing values were introduced. Similarly, the lines notated with labels “14.9%”, “31.8%”, “57.5%”, and “83.1%”, represent results of the imputed data sets with missing values of respectively, 14.9, 31.8, 57.5, and 83.1%.

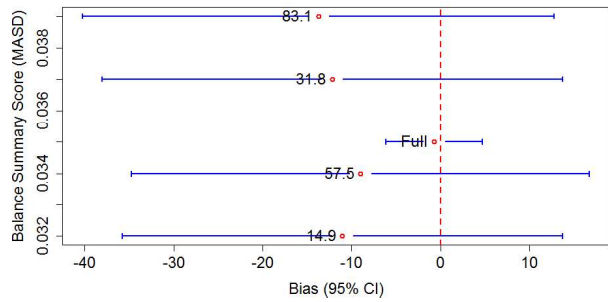
For all simulation scenarios, Figure 2 shows the bias on the x -axis vs. the balance summary score on the y -axis, respectively for the MASD and MMSD. For high imbalance, we would expect a high (i. e., ≥ 0.1) balance summary score, and for low imbalance, a low (i. e., < 0.1) balance summary score. In Figure 2a–d, the full data set had consistently low bias and balance summary score, for each treatment case (T_1 and T_3), when estimated over both, the MASD and MMSD. Compared to the imputed data scenarios, the CIs of the full data set analysis were narrow. For the imputed data sets, we see higher bias together with larger balance summary scores (Figure 2a–d), and larger variation in bias, as compared to the full data set. In total, we find a non-linear, but consistent association between the balance measures, and covariate imbalance, due to multiple imputation.



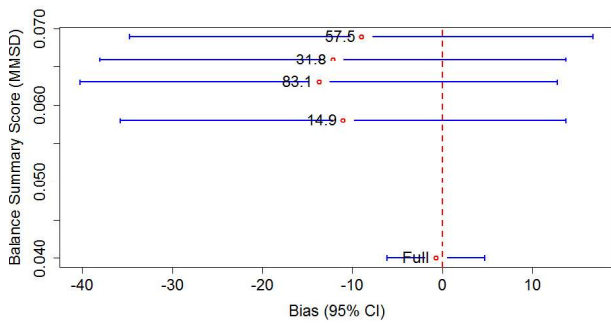
a. Bias vs Balance summary score (MASD) for T_1



b. Bias vs Balance summary score (MMSD) for T_1



c. Bias vs Balance summary score (MASD) for T_3



d. Bias vs Balance summary score (MMSD) for T_3

Figure 2: Bias vs. Balance summary scores after IPTW.

For this illustrative example presented, the low balance summary scores (Table 3) reflected low imbalance in covariates. Table 3 also shows a clear reduction of balance summary scores from unweighted to the weighted sets.

Table 3: Balance summary scores.

Balance metric	Full data set	Imputed data sets ^a with missing values of			
		14.9%	31.8%	57.5%	83.1%
Unweighted [MASD]	0.096	0.108	0.107	0.113	0.115
Weighted [MASD]	0.035	0.032	0.037	0.034	0.039
Unweighted [MMSD]	0.122	0.138	0.141	0.165	0.149
Weighted [MMSD]	0.040	0.058	0.066	0.069	0.063

^a The balance summary scores are the averaged MASD and MMSD over all covariate values.

Table 4: Performance measures per treatment group.

Treatment	Full data set	Imputed data sets with missing values of			
		14.9%	31.8%	57.5%	83.1%
T_1					
Bias ^a (MCSE)	0.22 (3.05)	3.35 (8.96)	4.76 (8.79)	1.07 (8.74)	5.7 (8.71)
Coverage (MCSE)	0.96 (0.02)	0.90 (0.09)	0.90 (0.09)	0.90 (0.09)	0.90 (0.09)
RelError (MCSE)	-2.32 (7.60)	6.19 (25.69)	8.96 (26.2)	11.69 (27.04)	8.71 (25.94)
T_3					
Bias ^a (MCSE)	-0.70 (2.75)	-11.01 (12.64)	-12.17 (13.21)	-8.97 (13.14)	-13.71 (13.52)
Coverage (MCSE)	0.96 (0.02)	0.90 (0.09)	0.90 (0.09)	0.90 (0.09)	0.90 (0.09)
RelError (MCSE)	9.68 (8.28)	-26.47 (17.66)	-29.45 (16.87)	-27.60 (17.40)	-31.60 (16.25)

^a Bias is measured in grams (g). Abbreviations: MCSE, Monte Carlo standard error, ModSE, Model standard error, RELError, Relative % error in ModeSE.

The estimates of the full and imputed data sets are presented in Table S1 in Supplementary Material Section S2. Table 4 presents performance measures for all simulated scenarios. The mean coverage rate was above the 95% rate for the full data set, however slightly lower in the imputed case scenarios. For all scenarios, we see low MCSE of the coverage rate. The relative percentage error in the model standard error (Table 4) shows that the results are conditioned on adequate number of simulated Monte Carlo runs. However, even for 10 Monte Carlo simulations (due to computational burden of the GPS algorithm), we can see in Table 4 that the relative percentage error is still relatively low.

For all imputed sets, diagnostic tools, considering the distribution of the observed and imputed data, as well as convergence of the MICE algorithm for multiple imputations, indicated that the imputed values were plausible. Given that the multiple imputations seemed plausible, the GPS for each scenario was calculated and balance assessed. Furthermore, on each imputed data set, and for each missing case, the maximal standard deviation of the weights were small (i. e. ≤ 1.5), indicating that the PS model was well specified (Xu et al. 2010).

Discussion

This study extended the definition of balance diagnostic to include multiple treatment exposure under multiple imputation. It has used simulations to investigate the performance of the proposed approach. Results from the analysis show that the mean differences after multiple imputation were unbiased for 14.9 to 83.1% missing data. The results were in accordance with the MASD and the MMSD balance summary scores after weighting, which indicate that covariates were balanced between groups.

The significance of this study can be appreciated when one considers the preponderance of data with missing covariates, especially in pharmacoepidemiological studies involving multiple treatment exposure

(Bandoli et al. 2018; Franklin et al. 2013). Specifically, the current approach allows incorporating external validation data, as suggested by Webb-Vargas et al. (2017), and for multiple treatment exposure. A recent study by Frank et al. (2019), combined the incorporation of validation data and multiple treatment exposure after multiple imputation. For balance assessment, the authors applied the MASD balance metric. In contrast to the present study, where the balance summary score of the metrics was calculated for validation purposes, the authors in Frank et al. (2019) assessed the balance of each covariate separately. This approach seems appropriate when the balance metrics are applied to real world data. When analyzing and interpreting balance metrics for each covariate separately, it should be taken into account that PS methods cannot balance covariates with small sample sizes (Rubin 1997). Hence, when there are small sample sizes within treatment groups among a specific covariate, the balance metrics might fail to reduce balance for that covariate.

Each PS method (i. e., weighting, stratification, matching or adjustment), has its own merits and limitations. Our choice of IPTW in combination of GPS has been confirmed in the recent published article by Nian et al. (2019). However, the simulation results by Yang et al. (2016), identified IPTW with multiple exposures as the worst performing approach. Although, more research is needed to identify under which condition, a GPS approach performed best, such a methodological comparison exceeds the scope of this current study.

With the illustrative example presented, we could validate the proposed balance metric for multiple treatment groups under multiple imputation. However, we did not see a linear relationship (as proportional to the missing data percentage) between the summary scores and bias of imputed data sets. One explanation could be that we were not able to perform analysis on the required number of Monte Carlo runs, and therefore the result might be more prone to noise due to multiple imputations. Although, we cannot rule out some influence of chance, it should be noted that low imbalance and low balance summary scores were consistent for all presented data cases. The results were also in accordance with the data-generation process. Nevertheless, the conclusions made, based on the simulation data example presented, should be considered taking into account the drawbacks of the data generation, as well as the low number of Monte-Carlo runs. The latter limitation should however be attributed to the method applied. In the present study the MASD and MMSD diagnostics resulted in the same conclusion. It seems however that the MMSD diagnostic maintains the balance characteristics better than the MASD for small percentages of missing data (i. e., <30%). Therefore, when MASD should be applied vs. MMSD needs further investigation. A recently published article considered the use of GPS under multiple imputation of missing covariates (De Vries, Van Smeden, and Groenwold 2018). The authors found that bias was introduced when the Classification and Regression Trees (CART) algorithm was applied directly for multiple imputation. However the authors concluded that the CART approach worked fine, when multiple imputation was applied before the estimation of the PS. Results presented in this paper confirm the finding.

The results reported come with some limitations. The determination of GPS and balance assessment are based on algorithms by Burgette, Griffin, and McCaffrey (2017); McCaffrey et al. (2013); McCaffrey, Ridgeway, and Morral (2004). Therefore, limitations and strength of their approach will also hold for the present study. One such limitation is the improper CIs, while strength of the determination of GPS via boosted regression models is the relatively small weights. Another limitation, is that our proposed balance diagnostic does not assess interaction and higher order terms (De Vries, Van Smeden, and Groenwold 2018). Based on our results, this however did not seem to be a major problem, as long as the imputation and PS models are well-specified.

Furthermore, in this illustrative simulation example, we assumed that for missing data the MAR assumption holds. However, future studies need to investigate how balance diagnostics perform under the missing completely at random and missing not at random assumptions, as well as when the missing data process is complex, i. e. linear and non-linear dependents on more than two observed variables. In addition, our example assumed that there are no unmeasured confounders. For real data, this assumption is unrealistic. Therefore, one could evaluate the effect of bias due to unmeasured confounders with trimming methods for the IPTW, as presented in Yoshida et al. (2018). There was not time varying confounding present in the illustrative example presented. However, if this should be the case, methods, such as proposed by Imai and Ratkovic (2015); Jackson (2016) should be implemented in the analysis.

In this study, we choose to propose and evaluate the diagnostic based on the standardized difference, as this is most commonly used in the pharmacoepidemiology literature, although other diagnostics exist, such as

the Kolmogorov–Smirnov test statistic or the C-statistic (Austin and Stuart 2015; Franklin et al. 2014). As other balance measures (e. g., C-statistic), showed promising results in assessing imbalance in covariates (Franklin et al. 2014), the MASD and MMSD could be alternated based on these balance measures. Such alternations however require separate validation.

The different scenarios of missing data percentages investigated can be considered as strength.

Different methods exist, which combine PSs and multiple imputation for a binary exposure as presented in Leyrat et al. (2019a). Such approaches should also be applied to the GPS in future studies, in order to better understand the effect of balance and missing data on causal effect estimates.

Given that, to our knowledge, no previous study has proposed or analyzed balance diagnostics for the GPS under multiple imputed data, this study is the first of its kind and will hopefully influence future research.

Conclusion

Based on simulation data, the proposed balance diagnostics seemed appropriate for balance assessment of the GPS after multiple imputation. Further research is needed to validate the results of the present study under different assumptions and conditions, and to apply the proposed diagnostics to real world data.

References

- Alexander, E. K., E. N. Pearce, G. A. Brent, R. S. Brown, H. Chen, C. Dosiou, W. A. Grobman, P. Laurberg, J. H. Lazarus, S. J. Mandel, and R. P. Peeters. 2017. “2017 Guidelines of the American Thyroid Association for the Diagnosis and Management of Thyroid Disease during Pregnancy and the Postpartum.” *Thyroid* 27: 315–89, <https://doi.org/10.1089/thy.2016.0457>.
- Austin, P. C. 2011. “An Introduction to Propensity Score Methods for Reducing the Effects of Confounding in Observational Studies.” *Multivariate Behavioral Research*, 46: 399–424, <https://doi.org/10.1080/00273171.2011.568786>.
- Austin, P. C. 2018. “Assessing Covariate Balance When Using the Generalized Propensity Score with Quantitative or Continuous Exposures.” *Statistical Methods in Medical Research*, 0962280218756159.
- Austin, P. C., and E. A. Stuart. 2015. “Moving Towards Best Practice When Using Inverse Probability of Treatment Weighting.(Iptw) Using the Propensity Score to Estimate Causal Treatment Effects in Observational Studies.” *Statistics in Medicine* 34: 3661–79, <https://doi.org/10.1002/sim.6607>.
- Azur, M. J., E. A. Stuart, C. Frangakis, and P. J. Leaf. 2011. “Multiple Imputation by Chained Equations: What is it and How Does it Work?” *International Journal of Methods in Psychiatric Research* 20: 40–9, <https://doi.org/10.1002/mpr.329>.
- Bandoli, G., G. M. Kuo, R. Sugathan, C. D. Chambers, M. Rolland, and K. Palmsten. 2018. “Longitudinal trajectories of Antidepressant Use in Pregnancy and the Postnatal Period.” *Archives of Women's Mental Health* 21: 411–19, <https://doi.org/10.1007/s00737-018-0809-2>.
- Bia, M. 2007. “The Propensity Score Method in Public Policy Evaluation: A Survey.” POLIS Working Paper Series.
- Billingsley, P. 2008. *Probability and Measure*. Hoboken: John Wiley & Sons.
- Blehar, M. C., C. Spong, C. Grady, S. F. Goldkind, L. Sahin, and J. A. Clayton. 2013. “Enrolling Pregnant Women: Issues in Clinical Research.” *Women's Health Issues* 23: e39–45, <https://doi.org/10.1016/j.whi.2012.10.003>.
- Borooh, V. K. 2002. *Logit and Probit: Ordered and Multinomial Models*, vol. 138. New York: Sage.
- Bray, B. C., J. J. Dziak, M. E. Patrick, and S. T. Lanza. 2018. “Inverse Propensity Score Weighting With a Latent Class Exposure: Estimating the Causal Effect of Reported Reasons for Alcohol Use on Problem Alcohol Use 16 Years Later.” *Prevention Science* 1–13, <https://doi.org/10.1007/s11121-018-0883-8>.
- Burgette, L., B. A. Griffin, and D. McCaffrey. 2017. “Propensity Scores for Multiple Treatments: A Tutorial for the Mnps Function in the Twang Package.” *R package*. Santa Monica: Rand Corporation, (Accessed July 2018).
- Chen, S. H. and E. H. Ip. 2015. “Behaviour of the Gibbs Sampler When Conditional Distributions are Potentially Incompatible.” *Journal of Statistical Computation and Simulation* 85: 3266–75, <https://doi.org/10.1080/00949655.2014.968159>.
- De Vries, B. B. P., M. Van Smeden, and R. H. Groenwold. 2018. “Propensity Score Estimation Using Classification and Regression Trees in the Presence of Missing Covariate Data.” *Epidemiologic Methods*.
- Doidge, J. C. 2018. “Responsiveness-Informed Multiple Imputation and Inverse Probability-Weighting in Cohort Studies with Missing Data that are Non-Monotone or not Missing at Random.” *Statistical Methods in Medical Research* 27: 352–63, <https://doi.org/10.1177/0962280216628902>.
- Dong, Y., and C. Y. J. Peng. 2013. “Principled Missing Data Methods for Researchers.” *SpringerPlus* 2: 222, <https://doi.org/10.1186/2193-1801-2-222>.

- Eulenburt, C., A. Suling, P. Neuser, A. Reuss, U. Canzler, T. Fehm, A. Luyten, M. Hellriegel, L. Woelber, and S. Mahner. 2016. "Propensity Scoring After Multiple Imputation in a Retrospective Study on Adjuvant Radiation Therapy in Lymph-Node Positive Vulvar Cancer." *PLoS One* 11: e0165705, <https://doi.org/10.1371/journal.pone.0165705>.
- Feng, P., X. H. Zhou, Q. M. Zou, M. Y. Fan, and X. S. Li. 2012. "Generalized Propensity Score for Estimating the Average Treatment Effect of Multiple Treatments." *Statistics in Medicine* 31: 681–97, <https://doi.org/10.1002/sim.4168>.
- Fong, C., C. Hazlett, and K. Imai. 2018. "Covariate Balancing Propensity Score for a Continuous Treatment: Application to the Efficacy of Political Advertisements." *The Annals of Applied Statistics* 12: 156–77, <https://doi.org/10.1214/17-aoas1101>.
- Frank, A. S., A. Lupattelli, D. S. Matteson, and H. Nordeng. 2018. "Maternal Use of Thyroid Hormone Replacement Therapy Before, During, and After Pregnancy: Agreement Between Self-Report and Prescription Records and Group-Based Trajectory Modeling of Prescription Patterns." *Clinical Epidemiology* 10: 1801–16, <https://doi.org/10.2147/clep.s175616>.
- Frank, A. S., A. Lupattelli, D. S. Matteson, H. M. Meltzer, and H. Nordeng. 2019. "Thyroid Hormone Replacement Therapy Patterns in Pregnant Women and Perinatal Outcomes in the Offspring." *Pharmacoepidemiology and Drug Safety*.
- Frank, A. S. J. 2019. "Thyroid Hormone Replacement Therapy During Pregnancy—Quantifying Medication Patterns and Associated Outcomes in the Offspring." In *Series of dissertations submitted to the Faculty of Mathematics and Natural Sciences*, vol. 1–251. Oslo: University of Oslo, URL <http://urn.nb.no/URN:NBN:no-73653>.
- Franklin, J. M., W. H. Shrank, J. Pakes, G. Sanf elix-Gimeno, O. S. Matlin, T. A. Brennan, and N. K. Choudhry. 2013. "Group-Based Trajectory Models: A New Approach to Classifying and Predicting Long-Term Medication Adherence." *Medical Care* 51: 789–96, <https://doi.org/10.1097/mlr.0000000000000002>.
- Franklin, J. M., J. A. Rassen, D. Ackermann, D. B. Bartels, and S. Schneeweiss. 2014. "Metrics for Covariate Balance in Cohort Studies of Causal Effects." *Statistics in Medicine* 33: 1685–99, <https://doi.org/10.1002/sim.6058>.
- Gasparini, A. and M. Lang. 2018. "rsimsum: Summarise Results from Monte Carlo Simulation Studies." *Journal of Open Source Software* 3: 739, <https://doi.org/10.21105/joss.00739>.
- Graham, J. W., A. E. Olchowski, and T. D. Gilreath. 2007. "How Many Imputations are Really Needed? Some Practical Clarifications of Multiple Imputation Theory." *Prevention Science* 8: 206–13, <https://doi.org/10.1007/s1121-007-0070-9>.
- Hayes, J. R. and J. I. Groner. 2008. "Using Multiple Imputation and Propensity Scores to Test the Effect of Car Seats and Seat Belt Usage on Injury Severity from Trauma Registry Data." *Journal of Pediatric Surgery* 43: 924–7, <https://doi.org/10.1016/j.jpedsurg.2007.12.043>.
- Hern n, M. A., A. Alonso, R. Logan, F. Grodstein, K. B. Michels, M. J. Stampfer, W. C. Willett, J. E. Manson, and J. M. Robins. 2008. "Observational Studies Analyzed Like Randomized Experiments: An Application to Postmenopausal Hormone Therapy and Coronary Heart Disease." *Epidemiology. Cambridge, Mass* 19: 766–79, <https://doi.org/10.1097/EDE.0b013e3181875e61>.
- Hill, J. 2004. "Reducing Bias in Treatment Effect Estimation in Observational Studies Suffering from Missing Data." Report no.04-01. US: Columbia University, January 2004.
- Hsu, C. H., and M. Yu. 2018. "Cox Regression Analysis with Missing Covariates via Nonparametric Multiple Imputation." *Statistical Methods in Medical Research*, 0962280218772592.
- Imai, K., and M. Ratkovic. 2015. "Robust Estimation of Inverse Probability Weights for Marginal Structural Models." *Journal of the American Statistical Association* 110: 1013–23, <https://doi.org/10.1080/01621459.2014.956872>.
- Imai, K., and D. A. Van Dyk. 2004. "Causal Inference with General Treatment Regimes: Generalizing the Propensity Score." *Journal of the American Statistical Association* 99: 854–66, <https://doi.org/10.1198/016214504000001187>.
- Imbens, G. W. 2000. "The Role of the Propensity Score in Estimating Dose-Response Functions." *Biometrika* 87: 706–10, <https://doi.org/10.1093/biomet/87.3.706>.
- Jackson, J. W. 2016. "Diagnostics for Confounding of Time-Varying and Other Joint Exposures." *Epidemiology. Cambridge, Mass* 27: 859, <https://doi.org/10.1097/ede.0000000000000547>.
- Jiang, M., and E. M. Foster. 2013. "Duration of Breastfeeding and Childhood Obesity: A Generalized Propensity Score Approach." *Health Services Research* 48: 628–51, <https://doi.org/10.1111/j.1475-6773.2012.01456.x>.
- Karahalios, A., L. Baglietto, J. B. Carlin, D. R. English, and J. A. Simpson. 2012. "A Review of the Reporting and Handling of Missing Data in Cohort Studies with Repeated Assessment of Exposure Measures." *BMC Medical Research Methodology* 12: 96, <https://doi.org/10.1186/1471-2288-12-96>.
- Kupzyk, K. A., and S. J. Beal. 2017. "Advanced Issues in Propensity Scores: Longitudinal and Missing Data." *The Journal of Early Adolescence* 37: 59–84, <https://doi.org/10.1177/0272431616636229>.
- Kurth, T., A. M. Walker, R. J. Glynn, K. A. Chan, J. M. Gaziano, K. Berger, and J. M. Robins. 2005. "Results of Multivariable Logistic Regression, Propensity Matching, Propensity Adjustment, and Propensity-Based Weighting Under Conditions of Nonuniform Effect." *American Journal Of Epidemiology* 163: 262–70, <https://doi.org/10.1093/aje/kwj047>.
- Lavori, P. W., R. Dawson, and D. Shera. 1995. "A Multiple Imputation Strategy for Clinical Trials with Truncation of Patient Data." *Statistics in Medicine* 14: 1913–25, <https://doi.org/10.1002/sim.4780141707>.
- Lechner, M. 2001. "Identification and Estimation of Causal Effects of Multiple Treatments Under the Conditional Independence Assumption." In *Econometric Evaluation of Labour Market Policies*, 43–58. Berlin: Springer.
- Lee, J. H., and J. Huber, Jr. 2011. "Multiple Imputation with Large Proportions of Missing Data: How Much is too Much?" In *United Kingdom Stata Users' Group Meetings 2011*, No. 23, Stata Users Group.

- Leyrat, C., S. R. Seaman, I. R. White, I. Douglas, L. Smeeth, J. Kim, M. Resche-Rigon, J. R. Carpenter, and E. J. Williamson. 2019a. "Propensity Score Analysis with Partially Observed Covariates: How Should Multiple Imputation be Used?" *Statistical Methods in Medical Research* 28: 3–19, <https://doi.org/10.1177/0962280217713032>.
- Li, F., and F. Li. 2019. "Propensity Score Weighting for Causal Inference with Multiple Treatments." *The Annals of Applied Statistics* 13: 2389–415, <https://doi.org/10.1214/19-aos1282>.
- Lumley, T. 2015. *Mitools: Tools for Multiple Imputation of Missing Data*. <https://cran.r-project.org/web/packages/mitools/mitools.pdf> (accessed July 2018).
- Lumley, T. 2018. *Survey: Analysis of Complex Survey Samples*. <http://r-survey.r-forge.r-project.org/survey/> (accessed July 2018).
- Lupattelli, A., M. Wood, K. Lapane, O. Spigset, and H. Nordeng. 2017. "Risk of Preeclampsia After Gestational Exposure to Selective Serotonin Reuptake Inhibitors and Other Antidepressants: A Study from the Norwegian Mother and Child Cohort Study." *Pharmacoepidemiology and Drug Safety* 26: 1266–76, <https://doi.org/10.1002/pds.4286>.
- Malla, L., R. Perera-Salazar, E. McFadden, M. Ogero, K. Stepniewska, and M. English. 2018. "Handling Missing Data in Propensity Score Estimation in Comparative Effectiveness Evaluations: A Systematic Review." *Journal of Comparative Effectiveness Research* 7: 271–9, <https://doi.org/10.2217/ce-2017-0071>.
- Marston, L., J. R. Carpenter, K. R. Walters, R. W. Morris, I. Nazareth, and I. Petersen. 2010. "Issues in Multiple Imputation of Missing Data for Large General Practice Clinical Databases." *Pharmacoepidemiology and Drug Safety* 19: 618–26, <https://doi.org/10.1002/pds.1934>.
- McCaffrey, D. F., B. A. Griffin, D. Almirall, M. E. Slaughter, R. Ramchand, and L. F. Burgette. 2013. "A Tutorial on Propensity Score Estimation for Multiple Treatments Using Generalized Boosted Models." *Statistics in Medicine*, 32, 3388–414, <https://doi.org/10.1002/sim.5753>.
- McCaffrey, D. F., G. Ridgeway, and A. R. Morral. 2004. "Propensity Score Estimation with Boosted Regression for Evaluating Causal Effects in Observational Studies." *Psychological Methods* 9: 403–25, <https://doi.org/10.1037/1082-989x.9.4.403>.
- Menard, S. 2002. *Applied Logistic Regression Analysis*, Vol. 106, New York: Sage.
- Miri, H. H., J. Hassanzadeh, A. Rajaeefard, M. Mirmohammadhani, and K. A. Angali. 2016. "Multiple Imputation to Correct for Nonresponse Bias: Application in Non-Communicable Disease Risk Factors Survey." *Global Journal Health Science* 8: 133–58, <https://doi.org/10.5539/gjhs.v8n1p133>.
- Mitra, R., and J. P. Reiter. 2016. "A Comparison of Two Methods of Estimating Propensity Scores After Multiple Imputation." *Statistical Methods in Medical Research* 25: 188–204, <https://doi.org/10.1177/0962280212445945>.
- Moleti, M., M. Di Mauro, G. Sturniolo, M. Russo, and F. Vermiglio. 2019. "Hyperthyroidism in the Pregnant Woman: Maternal and Fetal Aspects." *Journal of Clinical & Translational Endocrinology* 100190.
- Morris, T. P., I. R. White, and M. J. Crowther. 2019. "Using Simulation Studies to Evaluate Statistical Methods." *Statistics in Medicine* 38: 2074–102, <https://doi.org/10.1002/sim.8086>.
- Murray, J. S. 2018. "Multiple Imputation: A Review of Practical and Theoretical Findings." *Statistical Science* 33: 142–59, <https://doi.org/10.1214/18-sts644>.
- Nezvalová-Henriksen, K., O. Spigset, R. E. Brandlistuen, E. Ystrom, G. Koren, and H. Nordeng. 2016. "Effect of Prenatal Selective Serotonin Reuptake Inhibitor (SSRI) Exposure on Birthweight and Gestational Age: A Sibling-Controlled Cohort Study." *International Journal of Epidemiology* 45: 2018–29, <https://doi.org/10.1093/ije/dyw049>.
- Nguyen, T. L., G. S. Collins, J. Spence, J. P. Daurès, P. Devereaux, P. Landais, and Y. Le Manach. 2017. "Double-Adjustment in Propensity Score Matching Analysis: Choosing a Threshold for Considering Residual Imbalance." *BMC Medical Research Methodology* 17: 78, <https://doi.org/10.1186/s12874-017-0338-0>.
- Nian, H., C. Yu, J. Ding, H. Wu, W. D. Dupont, S. Brunwasser, T. Gebretsadik, T. V. Hartert, and P. Wu. 2019. "Performance Evaluation of Propensity Score Methods for Estimating Average Treatment Effects with Multi-Level Treatments." *Journal of Applied Statistics* 46: 853–73, <https://doi.org/10.1080/02664763.2018.1523375>.
- Nordeng, H., M. M. Van Gelder, O. Spigset, G. Koren, A. Einarson, and M. Eberhard-Gran. 2012. "Pregnancy Outcome after Exposure to Antidepressants and the Role of Maternal Depression: Results from the Norwegian Mother and Child Cohort Study." *Journal of Clinical Psychopharmacology* 32: 186–94, <https://doi.org/10.1097/jcp.0b013e3182490eaf>.
- Pandis, N. 2014. "Bias in Observational Studies." *American Journal of Orthodontics and Dentofacial Orthopedics* 145: 542–3, <https://doi.org/10.1016/j.ajodo.2014.01.008>.
- Qu, Y., and I. Lipkovich. 2009. "Propensity Score Estimation with Missing Values Using a Multiple Imputation Missingness Pattern. (MIMP) Approach." *Statistics in Medicine* 28: 1402–14, <https://doi.org/10.1002/sim.3549>.
- Reardon, S. F., and S. W. Raudenbush. 2009. "Assumptions of Value-Added Models for Estimating School Effects." *Education Finance and Policy* 4: 492–519, <https://doi.org/10.1162/edfp.2009.4.4.492>.
- Ridgeway, G., D. McCaffrey, A. Morral, L. Burgette, and B. A. Griffin. 2017. *twang: Toolkit for Weighting and Analysis of Nonequivalent Groups*. <https://cran.r-project.org/web/packages/twang/twang.pdf> (accessed July 2018).
- Rosenbaum, P. R., and D. B. Rubin. 1983. "The Central Role of the Propensity Score in Observational Studies for Causal Effects." *Biometrika* 70: 41–55, <https://doi.org/10.1093/biomet/70.1.41>.
- Rosenbaum, P. R., and D. B. Rubin. 1984. "Reducing Bias in Observational Studies Using Subclassification on the Propensity Score." *Journal of the American Statistical Association* 79: 516–24, <https://doi.org/10.1080/01621459.1984.10478078>.

- Rubin, D. B. 1986. “Comment: Which Ifs Have Causal Answers.” *Journal of the American Statistical Association* 81: 961–2, <https://doi.org/10.2307/2289065>.
- Rubin, D. B. 1997. “Estimating Causal Effects from Large Data Sets Using Propensity Scores.” *American Journal of Epidemiology* 127: 757–63, https://doi.org/10.7326/0003-4819-127-8_part_2-199710151-00064.
- Rubin, D. B. 2004. “On Principles for Modeling Propensity Scores in Medical Research.” *Pharmacoepidemiology and Drug Safety* 13: 855–7, <https://doi.org/10.1002/pds.968>.
- Spreuuenberg, M. D., A. Bartak, M. A. Croon, J. A. Hagenaars, J. J. Busschbach, H. Andrea, J. Twisk, and T. Stijnen. 2010. “The Multiple Propensity Score as Control for Bias in the Comparison of More Than Two Treatment Arms: An Introduction from a Case Study in Mental Health.” *Medical Care* 48: 166–74, <https://doi.org/10.1097/mlr.0b013e3181c1328f>.
- Stuart, E. A., B. K. Lee, and F. P. Leacy. 2013. “Prognostic Score–Based Balance Measures can be a Useful Diagnostic for Propensity Score Methods in Comparative Effectiveness Research.” *Journal of Clinical Epidemiology* 66: S84–S90.e1, <https://doi.org/10.1016/j.jclinepi.2013.01.013>.
- Sugihara, M. 2010. “Survival analysis using inverse probability of treatment weighted methods based on the generalized propensity score.” *Pharmaceutical Statistics: The Journal of Applied Statistics in the Pharmaceutical Industry* 9: 21–34, <https://doi.org/10.1002/pst.365>.
- Van Buuren, S., and K. Groothuis-Oudshoorn. 2011. “mice: Multivariate Imputation by Chained Equations in r.” *Journal of Statistical Software* 1–68. Also available at <https://www.jstatsoft.org/v45/i03/> (accessed July 20, 2018), URL <https://www.jstatsoft.org/v45/i03/>.
- Webb-Vargas, Y., K. E. Rudolph, D. Lenis, P. Murakami, and E. A. Stuart. 2017. “An Imputation-Based Solution to Using Mismeasured Covariates in Propensity Score Analysis.” *Statistical Methods in Medical Research* 26: 1824–37, <https://doi.org/10.1177/0962280215588771>.
- Xu, S., C. Ross, M. A. Raebel, S. Shetterly, C. Blanchette, and D. Smith. 2010. “Use of Stabilized Inverse Propensity Scores as Weights to Directly Estimate Relative Risk and its Confidence Intervals.” *Value in Health* 13: 273–7, <https://doi.org/10.1111/j.1524-4733.2009.00671.x>.
- Yang, S., G. W. Imbens, Z. Cui, D. E. Faries, and Z. Kadziola. 2016. “Propensity Score Matching and Subclassification in Observational Studies with Multi-Level Treatments.” *Biometrics* 72: 1055–65, <https://doi.org/10.1111/biom.12505>.
- Yoshida, K., D. H. Solomon, S. Haneuse, S. C. Kim, E. Patorno, S. K. Tedeschi, H. Lyu, J. M. Franklin, T. Stürmer, S. Hernández-Díaz, and R. J. Glynn. 2018. “Multinomial Extension of Propensity Score Trimming Methods: A Simulation study.” *American Journal of Epidemiology* 188: 609–16, <https://doi.org/10.1093/aje/kwy263>.
- Zhu, Y., D. L. Coffman, and D. Ghosh. 2015. “A Boosting Algorithm for Estimating Generalized Propensity Scores with Continuous Treatments.” *Journal of Causal Inference* 3: 25–40, <https://doi.org/10.1515/jci-2014-0022>.

Supplementary Material: This article contains supplementary material <https://doi.org/10.1515/em-2019-0003>.