

# Supplementary Material

**smdi: An R package to perform structural missing data investigations on partially observed confounders in real-world evidence studies**

**Authors:** Janick Weberpals<sup>1</sup>, RPh, PhD, Sudha R. Raman<sup>2</sup>, PhD, Pamela A. Shaw<sup>3</sup>, PhD, MS, Hana Lee<sup>4</sup>, PhD, Bradley G. Hammill<sup>2</sup>, DrPH, Sengwee Toh<sup>5</sup>, ScD, John G. Connolly<sup>5</sup>, ScD, Kimberly J. Dandreo<sup>5</sup>, MS, Fang Tian<sup>6</sup>, PhD, Wei Liu<sup>6</sup>, PhD, Jie Li<sup>6</sup>, PhD, José J. Hernández-Muñoz<sup>6</sup>, PhD, Robert J. Glynn<sup>1</sup>, PhD, ScD, Rishi J. Desai<sup>1</sup>, PhD

Author affiliations:

<sup>1</sup> Division of Pharmacoepidemiology and Pharmacoeconomics, Department of Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, MA

<sup>2</sup> Department of Population Health Sciences, Duke University School of Medicine, Durham, NC

<sup>3</sup> Biostatistics Division, Kaiser Permanente Washington Health Research Institute, Seattle, WA

<sup>4</sup> Office of Biostatistics, Center for Drug Evaluation and Research, US Food and Drug Administration, Silver Spring, MD

<sup>5</sup> Department of Population Medicine, Harvard Medical School and Harvard Pilgrim Health Care Institute, Boston, MA

<sup>6</sup> Office of Surveillance and Epidemiology, Center for Drug Evaluation and Research, US Food and Drug Administration, Silver Spring, MD

# 1 Supplementary Methods

## 1.1 Missingness assumptions

For the design and validation of the `smdi` R package functions, we employed comprehensive simulations and a real-world example [1]. These simulations followed structural assumptions on realistic missingness generating mechanisms that could be expected in electronic health record (EHR) data. The below directed acyclic diagrams (DAG) [2,3], where a directed edge represents a causal effect of one variable on another variable, outline and illustrate these assumptions along with a brief explanations and exemplary real-world EHR scenarios. Since the `smdi` package was developed with focus on partially observed confounders (as opposed to missingness in exposure or outcome), we simplified the below DAGs by leaving out the nodes for exposure and outcome.

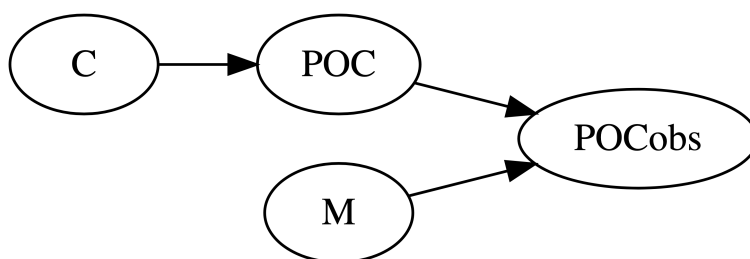
**Notation:** C = Fully observed confounder, POC = True values of the partially observed confounder, POCobs = Observed portion of the partially observed confounder, M = Missingness of POC with M=0 fully observed and M=1 fully missing, U = Unobserved confounder.

### 1.1.1 Missing completely at random (MCAR)

The missingness (M) of the partially observed confounder (POC) is independent of any observed and unobserved confounders. That is, no edge directs to M and the missingness is completely at random (MCAR).

Real-world example: A machine breaks that is used to measure and analyze a certain biomarker.

Supplementary Figure 1: Directed acyclic graph displaying a missing completely at random (MCAR) scenario.

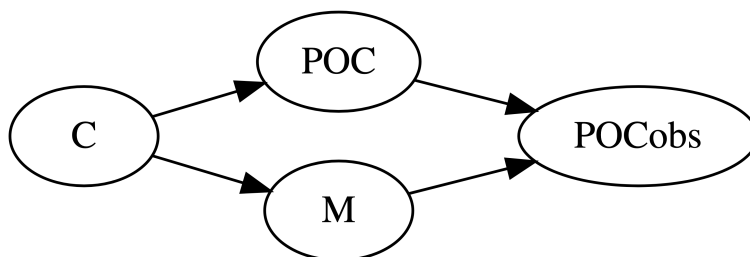


### 1.1.2 Missing at random (MAR)

The missingness (M) of the partially observed confounder (POC) can be explained by one or multiple fully observed confounders (C). That is, there is an edge between C to M and the missingness is at random (MAR).

Real-world example: Older patients are more likely to receive a certain biomarker test and age is a fully observed confounder which is measured for every patient in the dataset.

Supplementary Figure 2: Directed acyclic graph displaying a missing completely at random scenario.

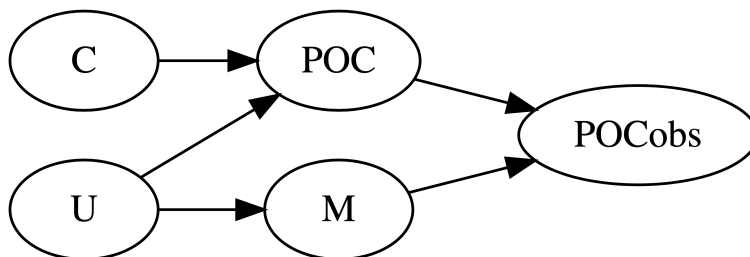


### 1.1.3 Missing not at random - unmeasured ( $\text{MNAR}_{\text{unmeasured}}$ )

The missingness (M) of the partially observed confounder (POC) can only be explained by one or multiple confounder(s) which are not observed in the dataset.

Real-world example: Patients with a certain performance status, which not observed in the dataset, are more likely to receive a lab test.

Supplementary Figure 3: Directed acyclic graph displaying a missing not at random (unmeasured) scenario.

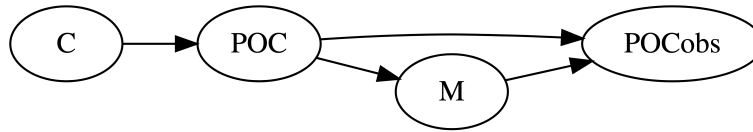


### 1.1.4 Missing not at random - value ( $MNAR_{\text{value}}$ )

The missingness (M) of the partially observed confounder (POC) depends on the true value of the partially observed confounder itself.

Real-world example: Patients with a history of normal biomarker values are systematically less likely to be tested again for the same biomarker.

Supplementary Figure 4: Directed acyclic graph displaying a missing not at random (value) scenario.



## 1.2 Missingness characterization

This section provides extensive information and details on key statistical principles and methods employed in the `smdi` package. Details on all package functions can also be found in the corresponding documentation section of the `smdi` package website

<https://janickweberpals.gitlab-pages.partners.org/smdi/reference/index.html>

or accessed in R by executing the the function name preceded with a question mark, e.g.:

```
?smdi_diagnose()
```

### 1.2.1 Group 1 diagnostics

The featured group 1 diagnostics in `smdi` focus on investigating potential differences in the distribution of characteristics between patients with and without an observed value for a partially observed confounder. To that end, different analytic approaches are leveraged which are described in detail for each function below.

- `smdi_hotelling()`: This function computes a two-sample Hotelling's T-squared test for the difference in two multivariate means based on the Hotelling's T-Square test statistic [4,5]. It assesses whether there is a statistically significant difference between the means of two groups in multivariate data. This test is an extension of the univariate t-test that examines for two groups (i.e., patients with and without a value observed for the confounder of interest) with  $p$  variables (e.g., patient or disease characteristics) whether

their mean vectors (p-dimensional vectors of means) are significantly different. To that end, the  $T^2$  statistic is derived by a [vector of covariate means](#) (i.e, one element for each covariate) [6], which are represented as  $\bar{x}_1$  (vector of covariate means in group 1) and  $\bar{x}_2$  (vector of covariate means in group 2) in Equation 1 below (with  $n_1$  and  $n_2$  representing the sample sizes of group 1 and 2, respectively, and  $S_p$  representing the pooled variance-covariance matrix) (Equation 1 adapted from [7]). To implement Hotelling’s test in the `smdi` package, we built a wrapper around the `Hotelling::hotelling.test()` function [5].

$$T^2 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \{ \mathbf{S}_p (\frac{1}{n_1} + \frac{1}{n_2}) \}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) \quad (1)$$

- `smdi_little()`: Little’s test, as opposed to Hotelling’s test, performs a single global test to evaluate the missing completely at random (MCAR) assumption [8]. It was developed since with an approach like Hotelling’s test, there can be concerns regarding multiplicity due to the multiple testing of each variable. Little’s test first [identifies a set of subgroups that share the same missing data patterns  \$j\$](#) . Across these missing data patterns ( $n_j$ ), it then tests for mean differences on every variable by comparing the observed means ( $\hat{\mu}_j$ ) versus expected population means ( $\hat{\mu}_j^{(ML)}$ ) which are [estimated using a maximum likelihood \(ML\) expectation-maximization \(EM\) algorithm](#) as implemented by the `norm::prelim.norm()` and `norm::em.norm()` functions in R (Equation 2 adapted from [9]). Under the assumption of MCAR, the test statistic ( $d^2$ ) approximates a chi-square distribution and is derived by computing the sum of squared standardized differences between the subgroup means and the expected population means weighted by both the estimated variance-covariance matrix (where  $\hat{\Sigma}_j$  is the maximum likelihood estimate of the covariance matrix) and the respective subgroup sizes [9,10]. To implement Little’s test in the `smdi` package, we built a wrapper around the `naniar::mcar_test()` function [11].

$$d^2 = \sum_{j=1}^J n_j (\hat{\mu}_j - \hat{\mu}_j^{(ML)})^T \hat{\Sigma}_j^{-1} (\hat{\mu}_j - \hat{\mu}_j^{(ML)}) \quad (2)$$

For both `smdi_hotelling()` and `smdi_little()`, a high test statistic and a low p-value would indicate differences in the groups compared. A limitation of both Hotelling’s and Little’s test is that they assume continuous variables following multivariate normality. Real-world data, however, are often of binary nature and to consider categorical data in the computation of these tests, [categorical data are one-hot encoded to binary dummy variables](#) before performing each test. In addition, these tests are sensitive to just small differences when study sizes are large which is why we recommend to perform these tests jointly with `smdi_asmd`.

- `smdi_asmd()`: The `smdi_asmd` functions computes the absolute standardized mean difference (ASMD) as a statistical measure for assessing dissimilarities in disease and patient

characteristics between individuals with and without observed value for a specific partially observed confounder (POC) of interest. If the median/average ASMD is high, this may indicate imbalance in patient covariate distributions which may be indicative of the POC following a missing at random (MAR) mechanism, i.e. the missingness is explainable by other observed covariates. Similarly, no imbalance between observed covariates may be indicative that missingness cannot be explained with observed covariates and the underlying missingness mechanism may be completely at random (MCAR) or not at random (e.g., missingness is only associated with unobserved factors or through the POC itself). This methodological approach follows the same theory that is often applied to measure covariate balance between two groups before and after matching or weighting on propensity scores, which are balancing scores often used to reduce confounding in real-world evidence studies [12]. While there isn't a universally established standard for the threshold of the ASMD indicating significant imbalance, an ASMD below 0.1 is often considered indicative of a negligible difference in the mean or prevalence of a covariate between two groups [12,13]. The advantage of ASMDs to assess the balance of patient characteristics between two groups with and without an observed value for the POC is that they are not influenced by sample size, do not come with many assumptions, are easy to interpret and applicable to various types of covariates. Equation 3 and Equation 4 illustrate how ASMDs are computed for continuous and binary covariates with  $\bar{x}$ ,  $s^2$  and  $\hat{p}$  indicating the mean, sample variance and prevalence of the binary variable of the covariate in patients with a missing ( $M$ ) and an observed ( $O$ ) value for the POC, respectively. In the `smdi` package, ASMDs are [extracted from tableone objects](#) using the `ExtractSmd` function [14].

$$ASMD_{continuous} = \left| \frac{(\bar{x}_M - \bar{x}_O)}{\sqrt{\frac{s_M^2 + s_O^2}{2}}} \right| \quad (3)$$

$$ASMD_{binary} = \left| \frac{(\hat{p}_M - \hat{p}_O)}{\sqrt{\frac{\hat{p}_M(1-\hat{p}_M) + \hat{p}_O(1-\hat{p}_O)}{2}}} \right| \quad (4)$$

### 1.2.2 Group 2 diagnostics

- `smdi_rf()`: The Group 2 diagnostics, implemented in the `smdi_rf()` function, assesses the ability to predict missingness based on observed covariates via the `smdi_rf()` function. This function trains and fits a random forest classification model [3,15] to predict the missing indicator of each POC given exposure, outcome, follow-up time, and covariates plus missingness indicator for other POC as the predictors. We chose a random forest classification model for the purpose of this package due its many beneficial features within the context of data structures that are frequently encountered in routine healthcare databases:

- Ability to implicitly model nonlinear and non-additive relationships between observed variables (i.e., higher order terms do not need to be explicitly specified).
- Recursive partitioning models like random forests have been found to work particularly well with sparse tabular data, which is the typical data type that is used for real-world evidence studies [16,17].
- Random Forests provide transparent feature importance measures, indicating the variables contributing significantly to predicting missingness. This aids in identifying key features driving the missingness mechanisms in the dataset.
- Multiple other studies have also reported good result of this random forest-based approach [3,18].

The default training parameters of the random forest model include a train-test split of the data (default if 70% training and 30% testing) which can be changed using the `train_test_ratio` parameter and the number of trees to grow using `ntree` (default is 1000 decision trees since a higher number of trees typically give more stable results). Users can optionally select `tune = TRUE` which will perform a 5-fold cross validation and a random search for the optimal number of variables randomly sampled as candidates at each split (`mtry`). However, users should be aware that this may lead to longer computation times with larger datasets which is why the default is set to `FALSE`. In summary, the following parameters are part of this function:

- `data`: dataframe or tibble object with partially observed/missing variables
- `covar`: character covariate or covariate vector with partially observed variable/column name(s) to investigate. If `NULL`, the function automatically includes all columns with at least one missing observation and all remaining covariates will be used as predictors
- `train_test_ratio`: numeric vector to indicate the test/train split ratio, e.g. `c(.7, .3)` which is the default
- `set_seed`: seed for reproducibility, defaults to 42
- `ntree`: integer, number of trees (defaults to 1000 trees)
- `n_cores`: integer, if `>1`, computations will be parallelized across the amount of cores specified in `n_cores`. This is important especially for larger datasets as the random forest training can be very time-consuming if done sequentially.

### 1.2.3 Group 3 diagnostics

To evaluate the association of the missing indicator variable of a POC ( $M_{POC}$ ) and the studied outcome in Group 3 Diagnostics, conventional univariate and multivariate regression models which are found used in the vast majority of real-world evidence studies are employed. The selection of regression type needs to be chosen using the `model` parameter and is determined

by the type of outcome that is studied. The possible outcome models include linear regression [19] for continuous outcomes (Equation 5), logistic regression [20] for binary outcomes (Equation 6) and a Cox proportional hazards regression [21] for time-to-event outcomes (Equation 7) with  $X_j$  indicating other covariates included in the multivariate models. The `form_lhs` specifies the left-hand side of the outcome formula, which, in case of `model = "linear"` or `model = "logistic"` just reflects the name of the column that contains the outcome and the form `Surv(time, status)` for `model = "cox"` indicating the time-to-event variable and the event status (0/1). The output of the functions automatically returns results for both the univariate and the multivariate models and a user has the choice to have the  $\beta_1 M_{POC}$  estimate exponentiated or not using the logical `exponentiated` parameter.

$$Y = \beta_0 + \beta_1 M_{POC} + \dots + \beta_2 X_j \quad (5)$$

$$\ln\left(\frac{p}{p-1}\right) = \beta_0 + \beta_1 M_{POC} + \dots + \beta_2 X_j \quad (6)$$

$$h(t) = h_0(t) e^{\sum \beta_1 M_{POC} + \dots + \beta_2 X_j} \quad (7)$$



## References

- 1 Weberpals J, Raman SR, Shaw PA, *et al.* A principled approach to characterize and analyze partially observed confounder data from electronic health records: A plasmode simulation study. *Submitted.* 2023.
- 2 Moreno-Betancur M, Lee KJ, Leacy FP, *et al.* [Canonical Causal Diagrams to Guide the Treatment of Missing Data in Epidemiologic Studies.](#) *American Journal of Epidemiology.* 2018;187:2705–15.
- 3 Sondhi A, Weberpals J, Yerram P, *et al.* A systematic approach towards missing lab data in electronic health records: A case study in non-small cell lung cancer and multiple myeloma. *CPT: Pharmacometrics & Systems Pharmacology.* Published Online First: 15 June 2023. doi: [10.1002/psp4.12998](#)
- 4 Hotelling H. [The Generalization of Student’s Ratio.](#) *The Annals of Mathematical Statistics.* 1931;2:360–78.
- 5 Curran J, Hersh T. [Hotelling: Hotelling’s  \$t^2\$  test and variants.](#) Published Online First: 2021.
- 6 Multivariate analysis of means for two groups. [https://rpubs.com/juanhklopper/multivariate\\_comparison\\_of\\_means\\_of\\_two\\_groups](https://rpubs.com/juanhklopper/multivariate_comparison_of_means_of_two_groups) (accessed 19 November 2023)
- 7 Applied multivariate statistical analysis (PennState). <https://online.stat.psu.edu/stat505/lesson/7/7.1/7.1.15> (accessed 19 November 2023)
- 8 Little RJA. [A Test of Missing Completely at Random for Multivariate Data with Missing Values.](#) *Journal of the American Statistical Association.* 1988;83:1198–202.
- 9 Enders CK. *Applied missing data analysis.* Guilford Publications 2022.
- 10 Little’s missing completely at random (MCAR) test. <https://search.r-project.org/CRAN/refmans/misty/html/na.test.html> (accessed 19 November 2023)
- 11 Tierney N, Cook D. Expanding tidy data principles to facilitate missing data exploration, visualization and assessment of imputations. 2023;105. doi: [10.18637/jss.v105.i07](#)
- 12 Austin PC. [An Introduction to Propensity Score Methods for Reducing the Effects of Confounding in Observational Studies.](#) *Multivariate Behavioral Research.* 2011;46:399–424.
- 13 Normand S-LT, Landrum MB, Guadagnoli E, *et al.* [Validating recommendations for coronary angiography following acute myocardial infarction in the elderly.](#) *Journal of Clinical Epidemiology.* 2001;54:387–98.
- 14 Yoshida K, Bartel A. [Tableone: Create ‘table 1’ to describe baseline characteristics with or without propensity score weights.](#) Published Online First: 2022.
- 15 Liaw A, Wiener M. [Classification and regression by randomForest.](#) 2002;2:18–22.

- 16 Shwartz-Ziv R, Armon A. [Tabular data: Deep learning is not all you need](#). *Information Fusion*. 2022;81:84–90.
- 17 Grinsztajn L, Oyallon E, Varoquaux G. Why do tree-based models still outperform deep learning on tabular data? 2022. <https://arxiv.org/abs/2207.08815>
- 18 Beaulieu-Jones BK, Lavage DR, Snyder JW, *et al.* [Characterizing and Managing Missing Structured Data in Electronic Health Records: Data Analysis](#). *JMIR Medical Informatics*. 2018;6:e11.
- 19 R Core Team. [R: A language and environment for statistical computing](#). Published Online First: 2022.
- 20 Friedman J, Tibshirani R, Hastie T. Regularization paths for generalized linear models via coordinate descent. 2010;33. doi: [10.18637/jss.v033.i01](https://doi.org/10.18637/jss.v033.i01)
- 21 Therneau TM. [A package for survival analysis in r](#). Published Online First: 2023.