

Overview:

The goal of this project is to outperform existing techniques in the literature related to nonmonotone missing data.

Completed:

- **Implemented simulation of monotone MAR data:** This is correspondingly easier than the subsequent nonmonotone MAR simulation. For this simulation we use the following approach:

1. Generate X , Y_1 , and Y_2 for elements $i = 1, \dots, n$.
2. Using the covariate X , determine the probability p_1 of Y_1 being observed for each element i .
3. Based on p_1 , determine if $R_1 = 1$.
4. If $R_1 = 0$, then $R_2 = 0$. Otherwise, using variables X and Y_1 , determine the probability p_{12} .
5. Based on p_{12} determine if $R_2 = 1$.

At the end of the algorithm, we have determined the values of binary variables R_1 and R_2 for each i and if either of them are equal to 1, the corresponding level of Y_k . As is common in this literature, the values of R_1 and R_2 determine if the corresponding variable Y_1 or Y_2 is missing or observed with $R = 1$ indicating Y being observed.

- **Implemented simulation of nonmonotone MAR data:** Following the approach of [robins1997non], I construct a nonmonotone MAR simulation with two response variables Y_1 and Y_2 and one covariate X . The algorithm to generate the data is the following:

1. Generate X , Y_1 , and Y_2 for elements $i = 1, \dots, n$.
2. Using the covariate X_i , generate probabilities for each element i p_0 , p_1 , and p_2 such that $p_0 + p_1 + p_2 = 1$.
3. Select one option based on the three probabilities for each element i . If 0 is selected: $R_1 = 0$ and $R_2 = 0$. If 1 is selected $R_1 = 1$. If 2 is selected, $R_2 = 1$.
4. We take the next step in multiple cases. If 0 was selected, we are done. If 1 was selected, we generate probabilities p_{12} based on X and Y_1 . Then based on this probability, we determine if $R_2 = 1$. In the same manner, if 2 was selected in the previous step, we generate probabilities p_{21} based on X and Y_2 . Then based on this probability, we determine if $R_2 = 1$.

Like the monotone MAR simulation this algorithm produces similar final results with the determination of binary variables R_1 and R_2 and variables X , Y_1 , and Y_2 . Unlike the monotone MAR case, the nonmonotone MAR includes observations with Y_2 observed and Y_1 missing.

- **Simulation 1 with Monotone MAR:** Following the algorithm described in the monotone MAR simulation bullet, we first generate data from the following distributions:

$$X_i \stackrel{iid}{\sim} N(0, 1)$$

$$Y_{1i} \stackrel{iid}{\sim} N(0, 1)$$

$$Y_{2i} \stackrel{iid}{\sim} N(\theta, 1)$$

Then, we create the probabilities $p_1 = \text{logistic}(x_i)$ and $p_{12} = \text{logistic}(y_{1i})$. Since, both x_i and y_1 are standard normal distributions, each of these probabilities is approximately 0.5 in expectation.

The goal of this simulation is to estimate θ . Alternatively, we can express this as solving the estimating equation:

$$g(\theta) \equiv Y_2 - \theta = 0.$$

We estimate θ using the following procedures:

- Oracle: This computes \bar{Y} using *both* the observed and missing data.
- IPW-Oracle: This is an IPW estimator using only the observed values of Y_2 . The weights (inverse probabilities) use the actual probabilities.
- IPW-Est: This is an IPW estimator using the probabilities that have been estimated by a logistic model.
- Semi: This is the monotone semiparametric efficient estimator from Slide 11 (Equation 2) of Dr. Kim’s Nonmonotone Missingness presentation.

We run this simulation with different values of θ , sample size of 2000, and 2000 Monte Carlo replications. Each algorithm for each replication generates $\hat{\theta}$. In the subsequent tables, we compute the bias, standard deviation (sd), t-statistic (where we test for a significant difference between the Monte Carlo mean $\hat{\theta}$ and the true θ) and the p-value of the t-statistic.

Table 1: True Value is -5

algorithm	bias	sd	tstat	pval
oracle	0.001	0.033	0.680	0.248
ipworacle	-0.012	0.392	-0.973	0.165
ipwest	0.007	0.186	1.178	0.120
semi	0.001	0.074	0.538	0.295

Table 2: True Value is 0

algorithm	bias	sd	tstat	pval
oracle	-0.001	0.031	-1.091	0.138
ipworacle	-0.001	0.085	-0.201	0.420
ipwest	0.000	0.085	-0.029	0.488
semi	0.000	0.079	0.112	0.455

Table 3: True Value is 5

algorithm	bias	sd	tstat	pval
oracle	0.000	0.033	-0.468	0.320
ipworacle	0.010	0.383	0.857	0.196
ipwest	-0.006	0.176	-1.020	0.154
semi	0.000	0.077	-0.049	0.481

Overall, these results are mostly what I would have expected. All of the algorithms estimate the true value of θ correctly in each case, with the oracle estimate having the smallest variance followed by the semiparametric algorithm. If there is anything surprising it is that the IPW estimator has better performance with the estimated weights compared to the true weights. However, I think that this is a known phenomenon.

- **Simulation 1 with Nonmonotone MAR:**

We generate variables (X, Y_1, Y_2) using the following setup:

$$\begin{bmatrix} X_i \\ \varepsilon_{1i} \\ \varepsilon_{2i} \end{bmatrix} \stackrel{iid}{\sim} N \left(\begin{bmatrix} 0 \\ 0 \\ \theta \end{bmatrix}, \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & \sigma_{yy} \\ 0 & \sigma_{yy} & 1 \end{bmatrix} \right).$$

Then,

$$y_{1i} = x_i + \varepsilon_{1i} \text{ and } y_{2i} = x_i + \varepsilon_{2i}.$$

Since we have nonmonotone data, our “Stage 1” probabilities are different. We compute the true Stage 1 probabilities being proportional to the following values:

$$p_0 = 0.2$$

$$p_1 = 0.4$$

$$p_2 = 0.4$$

However, we keep the same structure for the Stage 2 probabilities with: $p_{12} = \text{logistic}(y_1)$ and $p_{21} = \text{logistic}(y_2)$. The goal remains to estimate θ . We continue to use the Oracle algorithm and the IPW-Oracle algorithm. Since we have nonmonotone MAR data, we use the “Proposed” algorithm that is described on Slide 25 (Equation 12) of Dr. Kim’s presentation. The outcome models were estimated using logistic regression and OLS and correctly specified. The response model used the oracle estimates of the probabilities. This yields the following results:

Table 4: True Value is -5. $\text{Cor}(Y_1, Y_2) = 0$

algorithm	bias	sd	tstat	pval
oracle	0.000	0.032	0.285	0.388
ipworacle	-0.003	0.381	-0.318	0.375
proposed	0.000	0.038	0.492	0.311

Table 5: True Value is 5. $\text{Cor}(Y_1, Y_2) = 0$

algorithm	bias	sd	tstat	pval
oracle	0.000	0.032	0.285	0.388
ipworacle	-0.001	0.098	-0.479	0.316
proposed	0.000	0.037	0.505	0.307

- **Simulation 2 with Nonmonotone MAR:** We also want to simulate data that is correlated. For this simulation, we focus on $\text{Cov}(Y_1, Y_2)$. The data generating process now has $\sigma_{yy} \neq 0$. We are still interested in \bar{Y}_2 and we still run 2000 simulation with 2000 observations. In all of the next simulations the true value of $\theta = 0$. The results are the following:

Table 6: True Value is 0. $\text{Cor}(Y_1, Y_2) = 0.1$

algorithm	bias	sd	tstat	pval
oracle	0.001	0.031	1.623	0.052
ipworacle	0.001	0.077	0.762	0.223
proposed	0.001	0.037	1.366	0.086

Table 7: True Value is 0. $\text{Cor}(Y_1, Y_2) = 0.5$

algorithm	bias	sd	tstat	pval
oracle	0.001	0.032	1.486	0.069
ipworacle	0.004	0.086	1.890	0.029
proposed	0.000	0.041	0.172	0.432

Table 8: True Value is 0. $\text{Cor}(Y_1, Y_2) = 0.9$

algorithm	bias	sd	tstat	pval
oracle	0.001	0.032	0.706	0.240
ipworacle	0.003	0.098	1.395	0.082
proposed	-0.002	0.062	-1.339	0.090

- **Simulation 3 with Nonmonotone MAR:** This simulation aims to see if the proposed algorithm is doubly robust. First, we check with a misspecified outcome model. In this case the data generating procedure is the following:

$$\begin{bmatrix} X_i \\ \varepsilon_{1i} \\ \varepsilon_{2i} \end{bmatrix} \stackrel{iid}{\sim} N \left(\begin{bmatrix} 0 \\ 0 \\ \theta \end{bmatrix}, \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & \sigma_{yy} \\ 0 & \sigma_{yy} & 1 \end{bmatrix} \right).$$

Then,

$$y_{1i} = x_i + x_i^2 \varepsilon_{1i} \text{ and } y_{2i} = -x_i + x_i^3 + \varepsilon_{2i}.$$

This procedure causes X to influence both Y_1 and Y_2 and we still have correlation in the error terms of Y_1 and Y_2 . However, since neither Y_1 nor Y_2 are linear in X , the model will be misspecified. The response mechanisms are first generated MCAR with a probability of either Y_1 or Y_2 being the first variable observed to be 0.4. (There is a 0.2 probability neither is observed.) Then the probability of the other variable being observed is proportional to $\text{logistic}(y_k)$ where y_k is the y that has been observed. To ensure that the proposed method has the correct propensity score we use the oracle probabilities instead of estimating them. This yields the following:

Table 9: True Value is 0. $\text{Cor}(Y_1, Y_2) = 0$

algorithm	bias	sd	tstat	pval
oracle	0.000	0.075	0.014	0.494
ipworacle	0.002	0.107	0.876	0.191
proposed	-0.002	0.084	-1.063	0.144

Table 10: True Value is 0. $\text{Cor}(Y_1, Y_2) = 0.1$

algorithm	bias	sd	tstat	pval
oracle	-0.002	0.074	-1.479	0.070
ipworacle	0.000	0.106	-0.196	0.422
proposed	-0.003	0.083	-1.464	0.072

Table 11: True Value is 0. $\text{Cor}(Y_1, Y_2) = 0.5$

algorithm	bias	sd	tstat	pval
oracle	-0.003	0.074	-1.567	0.059
ipworacle	-0.002	0.108	-0.818	0.207
proposed	-0.003	0.083	-1.633	0.051

Thus, the proposed method is unbiased with a misspecified outcome model. We now show a simulation where the outcome model is correctly specified but the response model is not.

- **Simulation 4 with Nonmonotone MAR:** Continuing to test if the proposed algorithm is doubly robust, this simulation checks a misspecified response model. Instead of using oracle weights as in Simulation 3, we estimate the weights for the proposed method. The algorithms to which we compare still use the oracle weights.

Table 12: True Value is 0. $\text{Cor}(Y1, Y2) = 0$

algorithm	bias	sd	tstat	pval
oracle	0.000	0.032	-0.318	0.375
ipworacle	-0.001	0.079	-0.475	0.317
proposed	0.000	0.038	0.174	0.431

Table 13: True Value is 0. $\text{Cor}(Y1, Y2) = 0.1$

algorithm	bias	sd	tstat	pval
oracle	0.000	0.031	0.394	0.347
ipworacle	0.001	0.082	0.560	0.288
proposed	0.000	0.037	0.529	0.299

Table 14: True Value is 0. $\text{Cor}(Y1, Y2) = 0.5$

algorithm	bias	sd	tstat	pval
oracle	0.000	0.031	0.318	0.375
ipworacle	0.001	0.093	0.683	0.247
proposed	0.000	0.042	0.182	0.428

Comments from JK

- Three variables (X, Y_1, Y_2) should be correlated.