

Multivariate Analysis

Estimation with Missing Not at Random Data

Guillaume Blanc, Phd. Candidate,
University of Geneva

February 2022

Multivariate Data with Missing Values

	Y1	Y2	Y3
1	-0.34	-0.64	0.78
2	-1.16	-2.95	1.94
3	NA	-0.59	0.19
4	NA	-0.20	0.24
5	-1.75	-2.38	2.39
6	-0.99	-2.35	0.94
7	-1.30	-1.28	2.10
8	-1.32	-2.74	1.86
9	-0.55	0.06	0.59
10	-0.70	-1.54	2.41
11	NA	-0.42	1.08

	M1	M2	M3
1	0	0	0
2	0	0	0
3	1	0	0
4	1	0	0
5	0	0	0
6	0	0	0
7	0	0	0
8	0	0	0
9	0	0	0
10	0	0	0
11	1	0	0

Goal: estimate μ

Problem: missing values

Solution:

1. Impute by the mean (?!)
2. Correct the bias.

Requires:

- Generative model for \mathbf{Y}
- Missingness mechanism

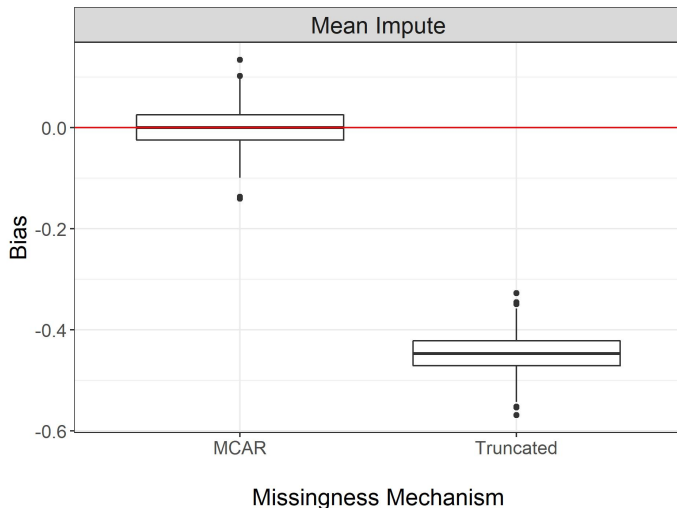
Figure: Data with missing values and the mask.

Experiment: Imputation by the Mean

Experiment: Simulate 1000 datasets with missing data. For each, get estimates of μ_1 .

Two settings:

- MCAR : Missing Completely at Random.
- Truncated : Only smallest 50% of the data are observed.



Generative Model for \mathbf{Y} : Factor Analysis / PPCA

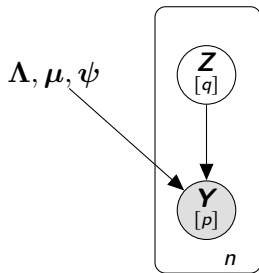


Figure: Generative model of the data.

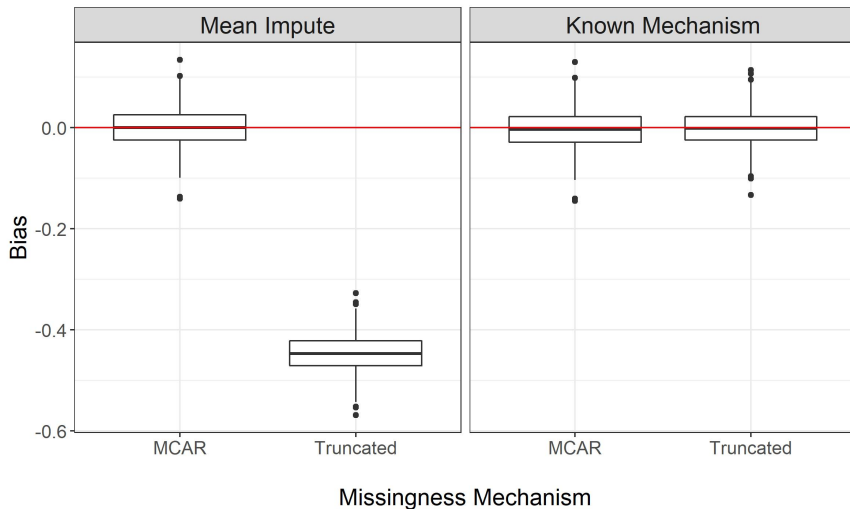
Assumptions: for $j, j' = 1, \dots, p$:

- $Y_j^{obs} \perp Y_{j'}^{obs} | \mathbf{Z}$ whenever $j \neq j'$,
- $Y_j^{obs} | \mathbf{Z} = \mathbf{z} \sim N(\mathbf{z}^\top \Lambda_j + \mu_j, \psi_j)$.

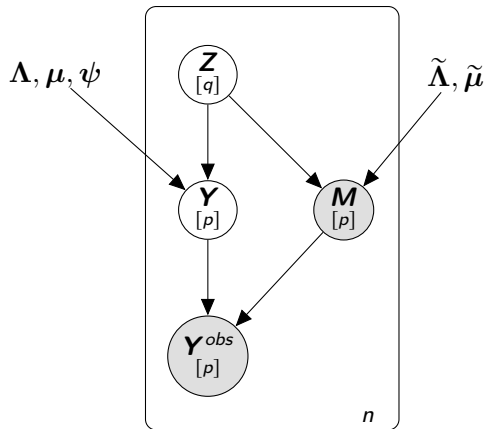
Consequence: the marginal distribution is

$$\mathbf{Y}^{obs} \sim MN(\mu, \Lambda \Lambda^\top + \psi).$$

Bias Correction with Known Missigness Mechanism



Modeling the Missingness Mechanism



Assumptions: for $j, j' = 1, \dots, p$,

- $Y_j^{obs} \perp Y_{j'}^{obs} | \mathbf{Z}$ whenever $j \neq j'$,
- $Y_j^{obs} | \mathbf{Z} = \mathbf{z} \sim N(\mathbf{z}^\top \Lambda_j. + \mu_j, \psi_j)$,
- $M_j \perp M_{j'} | \mathbf{Z}$ whenever $j \neq j'$,
- $M_j | \mathbf{Z} = \mathbf{z} \sim \text{Bernoulli}(\text{sigmoid}(\mathbf{z}^\top \tilde{\Lambda}_j. + \tilde{\mu}_j))$.

Figure: Generative model of the data and missingness mechanism.

Bias Correction with Estimated Missigness Mechanism

