

Missing Data Basics

Utrecht University Winter School: Missing Data in R



**Utrecht
University**

Kyle M. Lang

Department of Methodology & Statistics
Utrecht University

2022-02-03

Introduction

1. What's your name?
2. Where are you from/where do you work?
3. What type of research do you do?
4. What type of missing data problems do you encounter in your research?
5. What statistical software do you use/do you have programming experience?
6. What's your math background?



Outline

Missing Data Descriptives

Missing Data Mechanisms

Missing Data Treatments



A Little Notation

$Y :=$ An $N \times P$ Matrix of Arbitrary Data

$Y_{mis} :=$ The *missing* part of Y

$Y_{obs} :=$ The *observed* part of Y

$R :=$ An $N \times P$ pattern matrix encoding nonresponse



What are Missing Data?

Missing data are empty cells in a dataset where there should be observed values.

- The missing cells correspond to true population values, but we haven't observed those values.



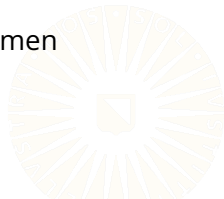
What are Missing Data?

Missing data are empty cells in a dataset where there should be observed values.

- The missing cells correspond to true population values, but we haven't observed those values.

Not every empty cell is a missing datum.

- Quality-of-life ratings for dead patients in a mortality study
- Firm profitability after the company goes out of business
- Self-reported severity of menstrual cramping for men
- Empty blocks of data following “gateway” items



Missing Data Descriptives



Missing Data Pattern

Missing data (or response) patterns represent unique combinations of observed and missing items.

- P items $\Rightarrow 2^P$ possible patterns.

	X	Y
1	x	y
2	x	.
3	.	y
4	.	.

Patterns for $P = 2$

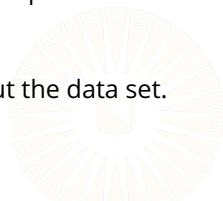
	X	Y	Z
1	x	y	z
2	x	y	.
3	x	.	z
4	.	y	z
5	x	.	.
6	.	.	z
7	.	y	.
8	.	.	.

Patterns for $P = 3$

Missing Data Pattern

The concept of a “missing data pattern” can also be used to classify the spatial arrangement of missing cells on a data set.

- Univariate
 - Missing data occur on only one variable
- Monotone
 - The proportion of complete elements, in both rows and columns, decreases when traversing the data set.
 - The observed cells can be arranged into a “staircase” pattern.
- Arbitrary
 - Missing values are “randomly” scattered throughout the data set.



Example Missing Data Patterns

	X	Y	Z
1	x	y	z
2	x	y	z
3	x	y	z
4	x	y	z
5	x	y	z
6	x	.	z
7	x	.	z
8	x	.	z
9	x	.	z
10	x	.	z

Univariate Pattern

	X	Y	Z
1	x	y	z
2	x	y	z
3	x	y	z
4	x	y	.
5	x	y	.
6	x	y	.
7	x	.	.
8	x	.	.
9	x	.	.
10	.	.	.

Monotone Pattern

	X	Y	Z
1	x	.	z
2	x	y	z
3	x	y	z
4	x	.	z
5	x	y	z
6	x	.	z
7	.	y	z
8	x	y	z
9	x	.	.
10	x	y	.

Arbitrary Pattern

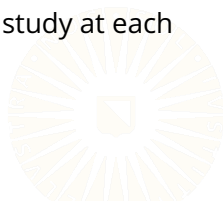
Nonresponse Rates

Proportion Missing

- The proportion of cells containing missing data
- Good early screening measure
- Should be computed for each variable, not for the entire dataset

Attrition Rate

- The proportion of participants that drop-out of a study at each measurement occasion



Nonresponse Rates

Proportion of Complete Cases

- The proportion of observations with no missing data
- Often reported but nearly useless quantity

Fraction of Missing Information

- Associated with an estimated parameter, not with an incomplete variable
- Like an R^2 for the missing data
- Most important diagnostic value for missing data problems
- Can only be computed after treating the missing data

Coverage Measures

Covariance Coverage

$$CC_{jk} = N^{-1} \sum_{n=1}^N r_{nj} r_{nk}$$

- The proportion of cases available to estimate a given pairwise relationship (e.g., a covariance between two variables)
- Very important to have adequate coverage of the parameters you want to estimate

Coverage Measures

Inbound Statistic

$$I_{jk} = \frac{\sum_{n=1}^N (1 - r_{nj}) r_{nk}}{\sum_{n=1}^N (1 - r_{nj})}$$

- The proportion of missing cases in Y_j for which Y_k is observed

Outbound Statistic

$$O_{jk} = \frac{\sum_{n=1}^N r_{nj} (1 - r_{nk})}{\sum_{n=1}^N r_{nj}}$$

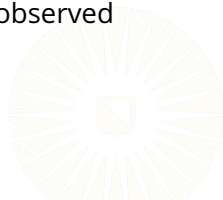
- The proportion of observed cases in Y_j for which Y_k is missing

Coverage Measures

Influx Coefficient

$$I_j = \frac{\sum_{k=1}^P \sum_{n=1}^N (1 - r_{nj}) r_{nk}}{\sum_{k=1}^P \sum_{n=1}^N r_{nk}}$$

- The proportion of observed cells in Y that exists in cases for which Y_j is missing
- How well the missing values in Y_j connect to the observed values in Y_{-j}



Coverage Measures

Outflux Coefficient

$$O_j = \frac{\sum_{k=1}^P \sum_{n=1}^N r_{nj}(1 - r_{nk})}{\sum_{k=1}^P \sum_{n=1}^N (1 - r_{nk})}$$

- The proportion of missing cells in Y that exists in cases for which Y_j is observed
- How well the observed values in Y_j connect to the missing values in Y_{-j}



Covariance Coverage Examples

- What is the coverage for $\text{cov}(X, Y)$?
- What is the coverage for $\text{cov}(W, Y)$?
- What about $\text{cov}(X, Z)$?

	W	X	Y	Z
1	w	x	y	.
2	w	x	y	.
3	w	x	y	.
4	w	x	y	.
5	w	x	y	.
6	w	.	y	z
7	w	.	y	z
8	w	.	y	z
9	w	.	y	z
10	w	.	y	z

Nonresponse Rate Examples

- What is the percent missing at Time 2?
- What is the attrition rate at Time 3?

	T1	T2	T3	T4
1	x1	x2	x3	x4
2	x1	x2	x3	x4
3	x1	x2	x3	x4
4	x1	x2	x3	.
5	x1	x2	x3	.
6	x1	x2	.	.
7	x1	x2	.	.
8	x1	.	.	.
9	x1	.	.	.
10	x1	.	.	.

Missing Data Mechanisms



Missing Data Mechanisms

MCAR:

$$P(R|Y_{mis}, Y_{obs}) = P(R)$$

MAR:

$$P(R|Y_{mis}, Y_{obs}) = P(R|Y_{obs})$$

MNAR:

$$P(R|Y_{mis}, Y_{obs}) \neq P(R|Y_{obs})$$



Simulate Some Toy Data

```
nObs <- 5000 # Sample Size
pm   <- 0.3  # Proportion Missing

sigma <- matrix(c(1.0, 0.5, 0.0,
                  0.5, 1.0, 0.3,
                  0.0, 0.3, 1.0),
               ncol = 3)

simDat <- as.data.frame(rmvnorm(nObs, c(0, 0, 0), sigma))
colnames(simDat) <- c("y", "x", "z")

x <- simDat$x
y <- simDat$y
z <- simDat$z

cor(y, x) # Check correlation between X and Y

[1] 0.5031885
```

MCAR Example

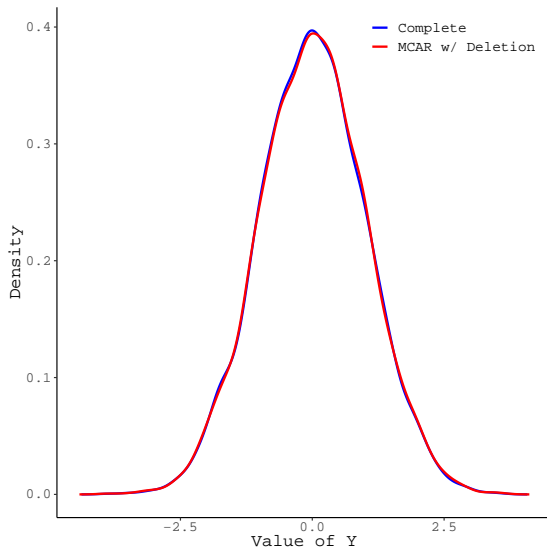
```
## Simulate MCAR Missingness:
rVec1 <- sample(1 : length(y), size = pm * length(y))

y2 <- y
y2[rVec1] <- NA

cor(y2, x, use = "pairwise") # Look at correlation

[1] 0.5229792
```

MCAR Example



MAR Example

```
## Simulate MAR Missingness:
rVec2 <- x < quantile(x, probs = pm)
mean(rVec2)

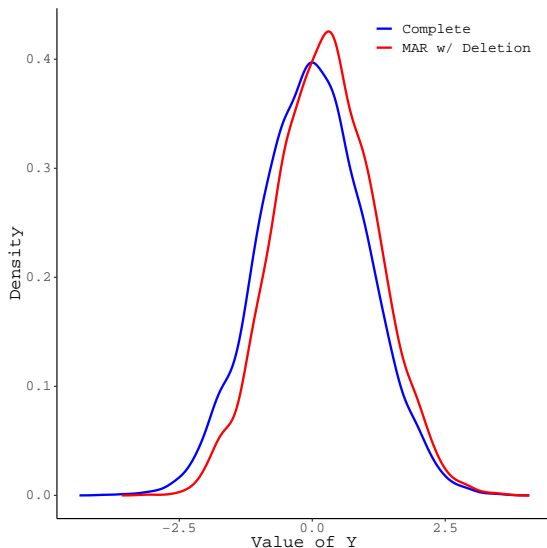
[1] 0.3

y3 <- y
y3[rVec2] <- NA

cor(y3, x, use = "pairwise") # Not looking so good :(

[1] 0.3870092
```


MAR Example



MNAR Example

```
## Simulate MNAR Missingness:
rVec3 <- y < quantile(y, probs = pm)
mean(rVec3)

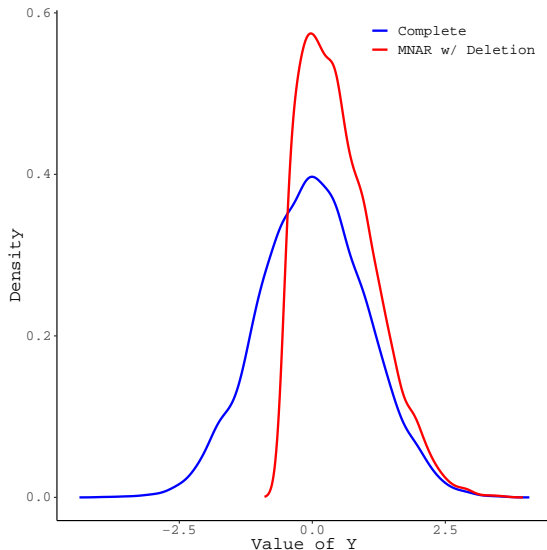
[1] 0.3

y4 <- y
y4[rVec3] <- NA

cor(y4, x, use = "pairwise") # Hmm...looks pretty bad.

[1] 0.3873756
```

MNAR Example



Crucial Nuance

In our previous MAR example, ignoring the predictor of missingness actually produces *Indirect MNAR*.

Crucial Nuance

In our previous MAR example, ignoring the predictor of missingness actually produces *Indirect MNAR*.

Question: What happens if we ignore the predictor of missingness, but that predictor is independent of our study variables?

Crucial Nuance

In our previous MAR example, ignoring the predictor of missingness actually produces *Indirect MNAR*.

Question: What happens if we ignore the predictor of missingness, but that predictor is independent of our study variables?

```
rVec3      <- z < quantile(z, probs = pm)
y5          <- y
y5[rVec3]   <- NA

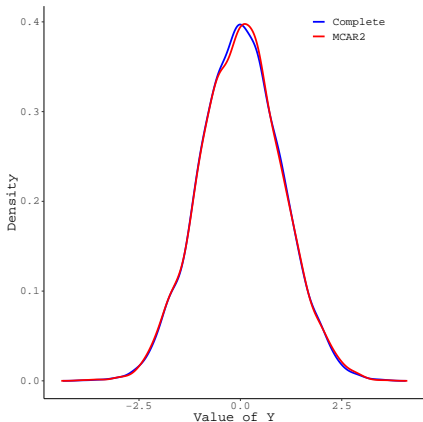
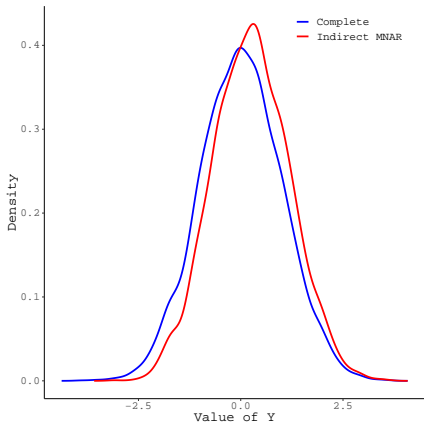
cor(y5, x, use = "pairwise")

[1] 0.5161666
```

Answer: We get back to MCAR :)

Crucial Nuance

The missing data mechanisms are not simply characteristics of an incomplete dataset; we also need to account for the analysis.



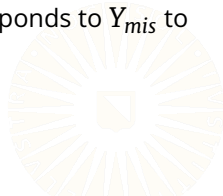
Testing the Missing Data Mechanism

We cannot test for MAR or MNAR.

- To do so would require knowing the values of the missing data.

We can test for MCAR (sort of).

- With MCAR, the missing data and the observed data should have the same distribution.
- We can test for MCAR by testing the distributions of *auxiliary variables*, Z .
 - Use a t -test to compare the subset of Z that corresponds to Y_{mis} to the subset corresponding to Y_{obs} .



Little's MCAR Test



Missing Data Treatments



Bad Methods (These almost never work)

Listwise Deletion (Complete Case Analysis)

- Use only complete observations for the analysis
 - Very wasteful (can throw out lots of useful data)
 - Loss of statistical power

Pairwise Deletion (Available Case Analysis)

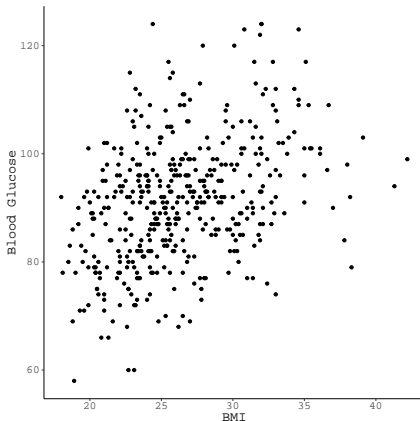
- Use only complete pairs of observations for analysis
 - Different samples sizes for different parameter estimates
 - Can cause computational issues



Bad Methods (These almost never work)

(Unconditional) Mean Substitution

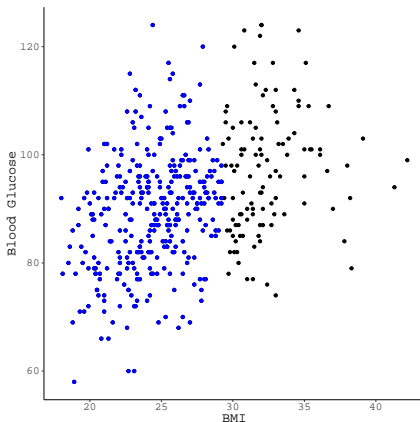
- Replace Y_{mis} with \bar{Y}_{obs}
 - Negatively biases regression slopes and correlations
 - Attenuates measures of linear association



Bad Methods (These almost never work)

(Unconditional) Mean Substitution

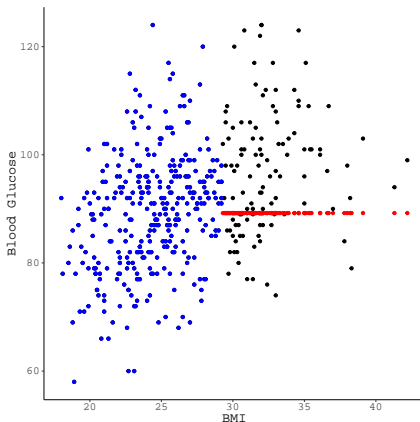
- Replace Y_{mis} with \bar{Y}_{obs}
 - Negatively biases regression slopes and correlations
 - Attenuates measures of linear association



Bad Methods (These almost never work)

(Unconditional) Mean Substitution

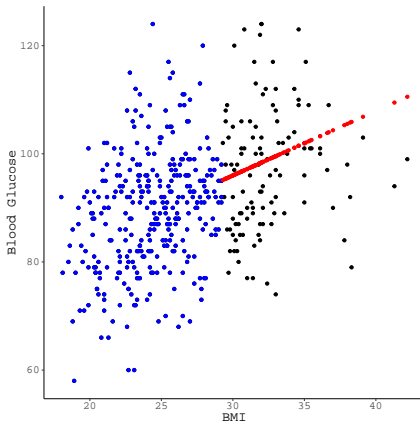
- Replace Y_{mis} with \bar{Y}_{obs}
 - Negatively biases regression slopes and correlations
 - Attenuates measures of linear association



Bad Methods (These almost never work)

Deterministic Regression Imputation (Conditional Mean Substitution)

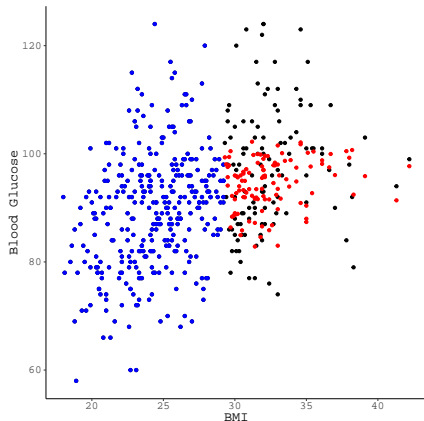
- Replace Y_{mis} with \hat{Y}_{mis} from some regression equation
 - Positively biases regression slopes and correlations
 - Inflates measures of linear association



Bad Methods (These almost never work)

Deterministic Regression Imputation (Conditional Mean Substitution)

- Replace Y_{mis} with \hat{Y}_{mis} from some regression equation
 - Positively biases regression slopes and correlations
 - Inflates measures of linear association



Bad Methods (These almost never work)

General Issues with Deletion-Based Methods

- Biased parameter estimates unless data are MCAR
- Generalizability issues

General Issues with Simple Single Imputation Methods

- Biased parameter estimates even when data are MCAR
- Attenuates variability in any treated variables

Bad Methods (These almost never work)

Averaging Available Items (Person-Mean Imputation)

- Compute aggregate scores using only available values
 - Missing data must be MCAR
 - Each item must contribute equally to the aggregate score

Last Observation Carried Forward (LOCF)

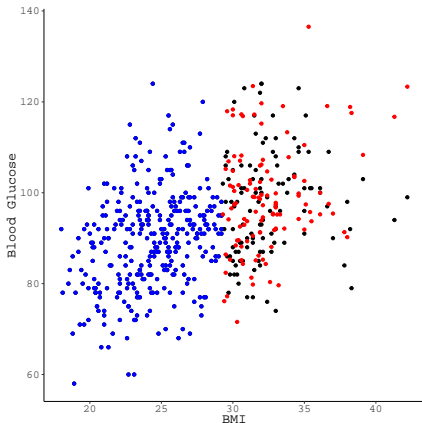
- Replace post-dropout values with the most recent observed value
 - Assume that dropouts would maintain their last known values
 - Attenuates estimates of growth/development

LOCF Viz

OK Methods (These work in some situations)

Stochastic Regression Imputation

- Fill Y_{mis} with \hat{Y}_{mis} plus some random noise.
 - Produces unbiased parameter estimates and predictions
 - Computationally efficient
 - Attenuates standard errors
 - Makes CIs and prediction intervals too narrow



OK Methods (These work in some situations)

Nonresponse Weighting

- Weight the observed cases to correct for nonresponse bias
 - Popular in survey research and official statistics
 - Only worth considering with *Unit Nonresponse*
 - Doesn't make any sense with *Item Nonresponse*



Expectation Maximization



Good Methods (These almost always work)

Multiple Imputation (MI)

- Replace the missing values with M plausible estimates
 - Essentially, a repeated application of stochastic regression imputation (with a particular type of regression model)
 - Produces unbiased parameter estimates and predictions
 - Produces “correct” standard errors, CIs, and prediction intervals
 - Very, very flexible
 - Computationally expensive



Good Methods (These almost always work)

What happens when we apply MI to our previous MAR example?

```
## Estimate imputation model:
miceOut1 <- mice(data      = data.frame(y3, x),
                 m         = 100,
                 maxit      = 1,
                 method     = c("norm", ""),
                 printFlag  = FALSE)

## Replace missing values with imputations:
impList1 <- list()
for(m in 1 : miceOut1$m)
  impList1[[m]] <- complete(miceOut1, m)
```


Good Methods (These almost always work)

```
## Estimate M correlations:
corList <-lapply(impList1,
                 FUN = function(impDat)
                   cor(impDat$x, impDat$y3)
                 )

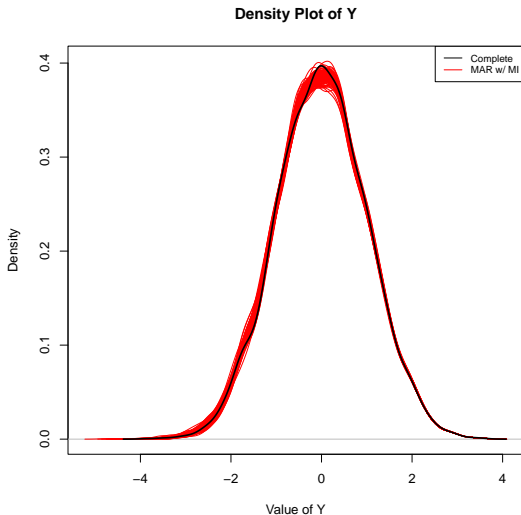
## Pool estimates:
mean(unlist(corList))

[1] 0.5108325
```

The MI-based parameter estimate looks good.

- MI produces unbiased estimates of the parameter when data are MAR.

Good Methods (These almost always work)



Good Methods (These almost always work)

What about applying MI to our MNAR example?

```
## Estimate imputation model:
miceOut2 <- mice(data      = data.frame(y4, x),
                 m         = 100,
                 maxit      = 1,
                 method     = c("norm", ""),
                 printFlag = FALSE)

## Replace missing values with imputations:
impList2 <- list()
for(m in 1 : miceOut2$m)
  impList2[[m]] <- complete(miceOut2, m)
```

Good Methods (These *almost* always work)

```
## Estimate M correlations:
corList2 <-lapply(impList2,
                  FUN = function(impDat)
                    cor(impDat$x, impDat$y4)
                  )

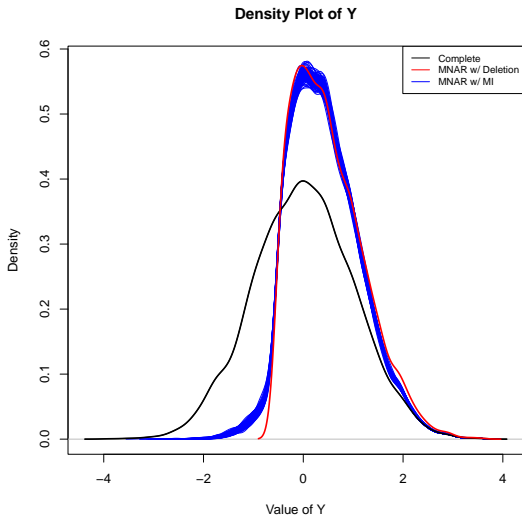
## Pool estimates:
mean(unlist(corList2))

[1] 0.4089334
```

The MI-based parameter estimate is still biased.

- MI cannot correct bias in parameter estimates when data are MNAR.

Good Methods (These *almost* always work)



Good Methods (These almost always work)

Bayesian Modeling

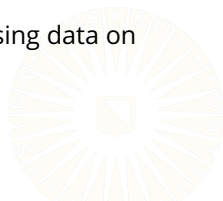
- Treat missing values as just another parameter to be estimated
 - Models can be directly estimated in the presence of missing data
 - Essentially, runs MI behind-the-scenes during model estimation
 - The predictors of nonresponse must be included in the model, somehow
 - Computationally expensive



Good Methods (These almost always work)

Full Information Maximum Likelihood (FIML)

- Adjust the objective function to only consider the observed parts of the data
 - Models are directly estimated in the presence of missing data
 - The predictors of nonresponse must be included in the model, somehow
 - Unless you write your own optimization program, FIML is only available for certain types of models
 - In linear regression models, FIML cannot treat missing data on predictors (if the predictors are taken as fixed)



References

