# Supervised learning with missing values

Julie Josse Head of the Inria-Inserm team PreMeDICaL: "Precision Medicine by Data Integration & Causal Learning"

February 9, 2023

#### Academic background:

- ▷ Assistant Professor at Institut Agro Rennes-Angers (2011-15)
- ▷ Visiting Scholar at Stanford University (2013-2015, 18 months)
- Professor at École Polytechnique (IP Paris) (2016 20)
- ▷ Visiting Researcher at Google Brain Paris (2019 2020). 2 days/week
- ▷ Senior Researcher at Inria Montpellier (Sept. 2020 -)
- $\triangleright~$  Visiting Researcher at Apple Paris (2023 ). 1 day/week

#### **Research topics:**

- ▷ Dimensionality reduction to visualize high dimensional heterogeneous data
- ▷ Missing values: EM alg., matrix completion, MNAR, supervised learning
- ▷ Causal inference: combining RCT & observational data, optimal policy
- ▷ Collaborations: medical (hospitals, SANOFI, etc.), energy (EDF), ecology

#### Software:

- R community: book R for Statistics, R foundation, R Forwards (widen the participation of minorities), R packages, R Task Views (missing, causal inf.)
- Website on missing values (R-miss-tastic), mobile application (ICUBAM)

## Traumabase project: decision support for trauma patients

- ▷ 30 000 French trauma patients<sup>1</sup>
- $\triangleright~250$  features from the accident site to the hospital discharge
- ▷ 30 hospitals
- ▷ 4000 new patients/ year

Center	Accident	Age	Sex	Lactactes	BP	Shock	Platelet	
Beaujon	fall	54	m	NM	180	yes	292 000	
Pitie	gun	26	m	NA	131	no	323 000	
Beaujon	moto	63	m	3.9	NR	yes	318 000	
Pitie	moto	30	W	Imp	107	no	211 000	

<sup>&</sup>lt;sup>1</sup>www.traumabase.eu - https://www.traumatrix.fr/

## Traumabase project: decision support for trauma patients

- ▷ 30 000 French trauma patients<sup>1</sup>
- $\triangleright$  250 features from the accident site to the hospital discharge
- ▷ 30 hospitals
- ▷ 4000 new patients/ year

Center	Accident	Age	Sex	Lactactes	BP	Shock	Platelet	
Beaujon	fall	54	m	NM	180	yes	292 000	
Pitie	gun	26	m	NA	131	no	323 000	
Beaujon	moto	63	m	3.9	NR	yes	318 000	
Pitie	moto	30	W	Imp	107	no	211 000	
:								•.

 $\Rightarrow$  Explain and predict hemorrhagic shock given pre-hospital features.

Ex: logistic regression/ random forests with missing values in covariates

Prospective study: real-time testing of models in the ambulance via a mobile data collection application (ShockMatrix playstore)

<sup>&</sup>lt;sup>1</sup>www.traumabase.eu - https://www.traumatrix.fr/

#### Missing data: important bottleneck in statistical practice



"One of the ironies of Big Data is that missing data play an ever more significant role"<sup>2</sup>

<sup>&</sup>lt;sup>2</sup>Zhu, Wang, Samworth. High-dimensional PCA with heterogeneous missingness. *JRSSB*. 2022.

#### Missing data: important bottleneck in statistical practice



"One of the ironies of Big Data is that missing data play an ever more significant role"<sup>2</sup>

Complete case analysis: delete incomplete samples

- Bias: Resulting sample not representative of the target population
- Information loss: Take a matrix with d features where each entry is missing with probability 1/100, remove a row (of length d) when one entry is missing

$$d = 5 \implies \approx 95\%$$
 of rows kept  
 $d = 300 \implies \approx 5\%$  of rows kept

<sup>&</sup>lt;sup>2</sup>Zhu, Wang, Samworth. High-dimensional PCA with heterogeneous missingness. *JRSSB*. 2022.

## Inference with missing values

First analysis to perform with missing data (and any data): descriptive study Visualize their patterns for clues as to how & why they occur FactoMineR<sup>3</sup>

Anomaly	Osthmot.	Improv.	SBP	DBP
No	NA	NA	150	100
Yes	Mannitol	Yes	99	41
No	NA	NA	110	76
Yes	SSH	NA	114	50
No	NA	NA	116	NA

<sup>&</sup>lt;sup>3</sup>Husson, J., Le. FactoMineR: An R Package for Multivariate Analysis. JSS. (2008)

First analysis to perform with missing data (and any data): descriptive study Visualize their patterns for clues as to how & why they occur  $FactoMineR^3$ 

Anomaly	Osthmot.	Improv.	SBP	DBP
No	NA	NA	Obs	Obs
Yes	Mannitol	Yes	Obs	Obs
No	NA	NA	Obs	Obs
Yes	SSH	NA	Obs	Obs
No	NA	NA	Obs	NA

Multiple Correspondence Analysis with numeric values coded as **Obs** & missing as **NA** 

<sup>&</sup>lt;sup>3</sup>Husson, J., Le. FactoMineR: An R Package for Multivariate Analysis. JSS. (2008)

First analysis to perform with missing data (and any data): descriptive study Visualize their patterns for clues as to how & why they occur FactoMineR<sup>3</sup>

Anomaly	Osthmot.	Improv.	SBP	DBP
No	NA	NA	Obs	Obs
Yes	Mannitol	Yes	Obs	Obs
No	NA	NA	Obs	Obs
Yes	SSH	NA	Obs	Obs
No	NA	NA	Obs	NA

Multiple Correspondence Analysis with numeric values coded as **Obs** & missing as **NA** 



MCA factor map

<sup>&</sup>lt;sup>3</sup>Husson, J., Le. FactoMineR: An R Package for Multivariate Analysis. JSS. (2008)

First analysis to perform with missing data (and any data): descriptive study Visualize their patterns for clues as to how & why they occur FactoMineR<sup>3</sup>

Anomaly	Osthmot.	Improv.	SBP	DBP
No	NA	NA	Obs	Obs
Yes	Mannitol	Yes	Obs	Obs
No	NA	NA	Obs	Obs
Yes	SSH	NA	Obs	Obs
No	NA	NA	Obs	NA

Multiple Correspondence Analysis with numeric values coded as **Obs** & missing as **NA** 



• Detect <u>nested variables</u>: (Anomaly



 $\Rightarrow$  Not a 'true' missing value, does not mask an underlying value

<sup>&</sup>lt;sup>3</sup>Husson, J., Le. FactoMineR: An R Package for Multivariate Analysis. JSS. (2008)

First analysis to perform with missing data (and any data): descriptive study Visualize their patterns for clues as to how & why they occur FactoMineR<sup>3</sup>

Anomaly	Osthmot.	Improv.	SBP	DBP	Anomaly-Osthmot.
No	NA	NA	Obs	Obs	No
Yes	Mannitol	Yes	Obs	Obs	Yes Mannitol
No	NA	NA	Obs	Obs	No
Yes	SSH	NA	Obs	Obs	Yes SSH
No	NA	NA	Obs	NA	No

Multiple Correspondence Analysis with numeric values coded as **Obs** & missing as **NA** 





- $\Rightarrow$  Not a 'true' missing value, does not mask an underlying value
- $\Rightarrow$  Solution: <u>recode</u> with a 3-level variable 'Yes Mannitol', 'Yes SSH', 'no'
- $\Rightarrow$  Feedback on data collection/encoding process

<sup>&</sup>lt;sup>3</sup>Husson, J., Le. FactoMineR: An R Package for Multivariate Analysis. JSS. (2008)

## Missing values mechanism: Rubin's taxonomy<sup>4,5</sup>

- <u>Random Variables</u>:
  - $\triangleright \ X \in \mathbb{R}^d$  : complete unavailable data

 $\triangleright \ M \in \{0,1\}^d$ : missing pattern, or mask,  $M_j = 1$  if and only if  $X_j$  is missing

For a pattern m, obs(m) indices of observed entries,  $X_{obs(m)}$  the vector of observed components

<sup>&</sup>lt;sup>4</sup>Rubin. Inference and missing data. *Biometrika*. 1976.

<sup>&</sup>lt;sup>5</sup>What Is Meant by "Missing at Random"? Seaman, et al. Statistical Science. 2013.

## Missing values mechanism: Rubin's taxonomy<sup>4,5</sup>

- Random Variables:
  - $\triangleright \ X \in \mathbb{R}^d$  : complete unavailable data

▷  $M \in \{0, 1\}^d$ : missing pattern, or mask,  $M_j = 1$  if and only if  $X_j$  is missing For a pattern m, obs(m) indices of observed entries,  $X_{obs(m)}$  the vector of observed components

Ex: Simulated missing values according to the 3 mechanisms (Orange points will be missing) in Systolic Blood Pressure - GCS is always observed



<sup>4</sup>Rubin. Inference and missing data. *Biometrika*. 1976.

<sup>b</sup>What Is Meant by "Missing at Random"? Seaman, et al. Statistical Science. 2013.

## Solutions to handle M(C)AR values (in the covariates)

Abundant literature: Creation of Rmistatic platform<sup>6</sup> (> 150 packages) Inferential aim: Estimate parameters & their variance, i.e.  $\hat{\beta}$ ,  $\hat{V}(\hat{\beta})$  to get confidence intervals with the appropriate coverage

 <sup>&</sup>lt;sup>6</sup>Mayer, J. et al. A unified platform for missing values methods and workflows. *R journal*. 2022.
 <sup>7</sup>Jiang, J. et al. Logistic Regression with Missing Covariates, Parameter Estimation, Model Selection and Prediction. *CSDA*. 2019. - Implementation in the misaem package

## Solutions to handle M(C)AR values (in the covariates)

Abundant literature: Creation of Rmistatic platform<sup>6</sup> (> 150 packages)

Inferential aim: Estimate parameters & their variance, i.e.  $\hat{\beta}$ ,  $\hat{V}(\hat{\beta})$  to get confidence intervals with the appropriate coverage

Modify the estimation process to deal with missing values Maximum likelihood inference: Expectation Maximization algorithms

Pros: Tailored toward a specific problem Cons: Few softwares even for simple models. Ex: logistic regression<sup>7</sup> Need to design one specific algorithm for each statistical method

 <sup>&</sup>lt;sup>6</sup>Mayer, J. et al. A unified platform for missing values methods and workflows. *R journal*. 2022.
 <sup>7</sup>Jiang, J. et al. Logistic Regression with Missing Covariates, Parameter Estimation, Model Selection and Prediction. *CSDA*. 2019. - Implementation in the misaem package

## Solutions to handle M(C)AR values (in the covariates)

Abundant literature: Creation of Rmistatic platform<sup>6</sup> (> 150 packages)

Inferential aim: Estimate parameters & their variance, i.e.  $\hat{\beta}$ ,  $\hat{V}(\hat{\beta})$  to get confidence intervals with the appropriate coverage

Modify the estimation process to deal with missing values Maximum likelihood inference: Expectation Maximization algorithms

Pros: Tailored toward a specific problem Cons: Few softwares even for simple models. Ex: logistic regression<sup>7</sup> Need to design one specific algorithm for each statistical method

#### (Multiple) imputation to get a complete data set

Pros: Any analysis can be performed, mice R package Cons: Generic

 <sup>&</sup>lt;sup>6</sup>Mayer, J. et al. A unified platform for missing values methods and workflows. *R journal*. 2022.
 <sup>7</sup>Jiang, J. et al. Logistic Regression with Missing Covariates, Parameter Estimation, Model Selection and Prediction. *CSDA*. 2019. - Implementation in the misaem package

## Single imputation by the mean

$$\triangleright (x_{i1}, x_{i2}) \underset{\text{i.i.d.}}{\sim} \mathcal{N}_2((\mu_{x_1}, \mu_{x_2}), \Sigma_{x_1x_2})$$



$$\hat{\mu}_{x_2} = -0.01 \hat{\sigma}_{x_2} = 1.01 \hat{\sigma}_{x_2} = 1.01 \hat{\rho} = 0.66$$

#### Single imputation by the mean

- $\triangleright (x_{i1}, x_{i2}) \underset{i i d}{\sim} \mathcal{N}_2((\boldsymbol{\mu}_{x_1}, \boldsymbol{\mu}_{x_2}), \boldsymbol{\Sigma}_{x_1 x_2})$
- $\triangleright$  70 % of missing entries completely at random on  $X_2$



$$\hat{\mu}_{x_2} = 0.18$$
  
 $\hat{\sigma}_{x_2} = 0.9$   
 $\hat{
ho} = 0.6$ 

## Single imputation by the mean

- $\triangleright (x_{i1}, x_{i2}) \underset{i.i.d.}{\sim} \mathcal{N}_2((\mu_{x_1}, \mu_{x_2}), \Sigma_{x_1x_2})$
- $\triangleright$  70 % of missing entries completely at random on  $X_2$
- Estimate parameters on the mean imputed data



Mean imputation deforms joint and marginal distributions

#### Mean imputation should be avoided for estimation



PCA with mean imputation

library(FactoMineR)
PCA(ecolo)
Warning message: Missing
are imputed by the mean
of the variable:
You should use imputePCA
from missMDA

#### EM-PCA

library(missMDA)
imp <- imputePCA(ecolo)
PCA(imp\$comp)</pre>

Ecological data: n = 69000 species - 6 traits. Estimated correlation between Pmass & Rmass  $\approx 0$  (mean imputation) or  $\approx 1$  (EM PCA)<sup>8</sup>

<sup>&</sup>lt;sup>8</sup>J. & Husson. missMDA: Handling Missing Values in Multivariate Data Analysis, *JSS*. 2016.

### Objective: to impute while preserving distribution

Assuming a bivariate gaussian distribution  $x_{i2} = \beta_0 + \beta_1 x_{i1} + \varepsilon_i$ ,  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ 

- ▷ Regression imputation: Estimate  $\beta$  (here with complete data) and impute  $\hat{x}_{i2} = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} \Rightarrow$  variance underestimated and correlation overestimated
- ▷ Stochastic reg. imputation: Estimate  $\beta$  and  $\sigma$  impute from the predictive  $\hat{x}_{i2} \sim \mathcal{N}\left(\beta_0 + \hat{\beta}_1 x_{i1}, \hat{\sigma}^2\right) \Rightarrow$  preserve distributions



#### Assuming a joint distribution

 $\triangleright \; \mathsf{Gaussian} \; \mathsf{model} \; x_i \sim \mathcal{N} \left( \mu, \Sigma 
ight)$ 

- > Low rank :  $X_{n imes d} = \mu_{n imes d} + \varepsilon \ \varepsilon_{ij} \stackrel{iid}{\sim} \mathcal{N}\left(0, \ \sigma^2\right)$  with  $\mu$  of low rank
  - $\Rightarrow$  Powerful in recommendation system: Netflix prize 90% of missing
  - $\Rightarrow$  Use similarities between rows & links between variables + reduct. of dim.
  - $\Rightarrow$  Different regularization depending on noise regime<sup>9,10</sup>
  - $\Rightarrow$  Count data,<sup>11</sup> ordinal data, categorical data, blocks/multilevel data<sup>12</sup>
- ▷ Using optimal transport,<sup>13</sup> deep generative models

<sup>&</sup>lt;sup>9</sup>J. & Sardy. Adaptive Shrinkage of singular values. *Stat & Computing*. 2015.

<sup>&</sup>lt;sup>10</sup>J. & Wager. Stable autoencoding for regularized low-rank matrix estimation. *JMLR*. 2016.

<sup>&</sup>lt;sup>11</sup>Robin, Klopp, J., Moulines, Tibshirani. Main effects and interactions in mixed and incomplete data frames. 2019. *JASA*.

<sup>&</sup>lt;sup>12</sup>J. et al. Imputation of mixed data with multilevel SVD. JCGS. 2018.

<sup>&</sup>lt;sup>13</sup>Muzelec, Cuturi, Boyer, J. Missing Data Imputation using Optimal Transport. *ICML*. 2020.

#### Assuming a joint distribution

 $\triangleright$  Gaussian model  $x_i \sim \mathcal{N}\left(\mu, \Sigma
ight)$ 

- $\triangleright \ \underline{\text{Low rank}}: \ X_{n \times d} = \mu_{n \times d} + \varepsilon \ \varepsilon_{ij} \overset{\text{iid}}{\sim} \mathcal{N}\left(0, \ \sigma^{2}\right) \ \text{with} \ \mu \text{ of low rank}$ 
  - $\Rightarrow$  Powerful in recommendation system: Netflix prize 90% of missing
  - $\Rightarrow$  Use similarities between rows & links between variables + reduct. of dim.
  - $\Rightarrow$  Different regularization depending on noise regime<sup>9,10</sup>
  - $\Rightarrow$  Count data,<sup>11</sup> ordinal data, categorical data, blocks/multilevel data<sup>12</sup>
- ▷ Using optimal transport,<sup>13</sup> deep generative models

#### Iterating conditional models (joint distribution implicitly defined)

with multinomial, Poisson regression (ICE: Imputation by Chained Equations)
 iterative imputation of each variable by random forests

<sup>&</sup>lt;sup>9</sup>J. & Sardy. Adaptive Shrinkage of singular values. *Stat & Computing*. 2015.

<sup>&</sup>lt;sup>10</sup>J. & Wager. Stable autoencoding for regularized low-rank matrix estimation. *JMLR*. 2016.

<sup>&</sup>lt;sup>11</sup>Robin, Klopp, J., Moulines, Tibshirani. Main effects and interactions in mixed and incomplete data frames. 2019. *JASA*.

<sup>&</sup>lt;sup>12</sup>J. et al. Imputation of mixed data with multilevel SVD. JCGS. 2018.

<sup>&</sup>lt;sup>13</sup>Muzelec, Cuturi, Boyer, J. Missing Data Imputation using Optimal Transport. *ICML*. 2020.

#### $\Rightarrow$ Incomplete Traumabase

$X_1$	$X_2$	$X_3$	 Y
NA	20	10	 shock
-6	45	NA	 shock
0	NA	30	 no shock
NA	32	35	 shock
-2	NA	12	 no shock
1	63	40	 shock

$X_1$	$X_2$	<i>X</i> <sub>3</sub>	 Y
NA	20	10	 shock
-6	45	NA	 shock
0	NA	30	 no shock
NA	32	35	 shock
-2	NA	12	 no shock
1	63	40	 shock

#### $\Rightarrow$ Incomplete Traumabase

#### $\Rightarrow \mathsf{Completed} \ \mathsf{Traumabase}$

X1	$X_2$	$X_3$	 Y
3	20	10	 shock
-6	45	6	 shock
0	4	30	 no shock
-4	32	35	 shock
-2	75	12	 no shock
1	63	40	 shock

$X_1$	$X_2$	$X_3$	 Y
NA	20	10	 shock
-6	45	NA	 shock
0	NA	30	 no shock
NA	32	35	 shock
-2	NA	12	 no shock
1	63	40	 shock

 $\Rightarrow$  Incomplete Traumabase

#### $\Rightarrow$ Completed Traumabase

$X_1$	$X_2$	$X_3$	 Y
3	20	10	 shock
-6	45	6	 shock
0	4	30	 no shock
-4	32	35	 shock
-2	75	12	 no shock
1	63	40	 shock

A single value can't reflect the uncertainty of prediction Multiple impute 1) Generate M plausible values for each missing value

$X_1$	$X_2$	<i>X</i> <sub>3</sub>	Y
3	20	10	s
-6	45	6	s
0	4	30	no s
-4	32	35	s
-2	75	12	no s
1	63	40	s

$X_1$	$X_2$	$X_3$	Y
-7	20	10	s
-6	45	9	s
0	12	30	no s
13	32	35	s
-2	10	12	no s
1	63	40	s

$X_1$	$X_2$	<i>X</i> <sub>3</sub>	Y
7	20	10	s
-6	45	12	s
0	-5	30	no s
2	32	35	s
-2	20	12	no s
1	63	40	s

## Visualization of the imputed values<sup>14</sup>

Supplementary projection

Dim 1 (71.33%)

$X_1$	X2	X3	Y
3	20	10	s
-6	45	6	s
0	4	30	no s
-4	32	35	s
-2	15	12	no s
1	63	40	s

N

Ŷ

7

- 6

\_2

Dim 2 (16.94%)

X1	X <sub>2</sub>	X3	Y
-7	20	10	s
-6	45	9	s
0	12	30	no s
13	32	35	s
-2	10	12	no s
1	63	40	s



library(missMDA) MIPCA(traumadata)

Projection of the *M* imputed data on a 'compromise' subspace (PCA with missing values)

Is it possible to handle 30% of missing values? 50%?, etc. Both % of missing values & signal matter (5% of NA can be an issue)

<sup>&</sup>lt;sup>14</sup>J. et al. Multiple imputation in principal component analysis. ADAC. 2011.

#### 1) Generate M plausible values for each missing value

$X_1$	X2	X3	Y
3	20	10	s
-6	45	6	s
0	4	30	no s
-4	32	35	s
1	63	40	s
-2	15	12	no s

<i>x</i> <sub>1</sub>	X2	X3	Y
-7	20	10	s
-6	45	9	s
0	12	30	no s
13	32	35	s
1	63	40	s
-2	10	12	no s

<i>x</i> <sub>1</sub>	X2	X3	Y
7	20	10	s
-6	45	12	s
0	-5	30	no s
2	32	35	s
1	63	40	s
-2	20	12	no s

2) Perform the analysis on each imputed data set:  $\hat{\beta}_m, \widehat{Var}\left(\hat{\beta}_m\right)$ 

3) Combine the results (Rubin's rules):

$$\hat{\beta} = \frac{1}{M} \sum_{m=1}^{M} \hat{\beta}_{m}$$

$$T = \underbrace{\frac{1}{M} \sum_{m=1}^{M} \widehat{Var}\left(\hat{\beta}_{m}\right)}_{\text{Within-imputation variance}} + \underbrace{\left(1 + \frac{1}{M}\right) \frac{1}{M-1} \sum_{m=1}^{M} \left(\hat{\beta}_{m} - \hat{\beta}\right)^{2}}_{\text{Between-imputation variance}}$$

 $\Rightarrow$  Variability of missing values taken into account.

• Methods used in practice are the one implemented in a sustainable way: few implementations of EM strategies

• "Imputation is both seductive & dangerous" (Dempster & Rubin, 1983). Seductive: "can lull the user into the pleasant state of believing that the data are complete

Dangerous: "it lumps together situations where the problem is minor enough to be handled in this way & situations where estimators applied to the imputed data have substantial biases."

• Methods used in practice are the one implemented in a sustainable way: few implementations of EM strategies

• "Imputation is both seductive & dangerous" (Dempster & Rubin, 1983). <u>Seductive</u>: "can lull the user into the pleasant state of believing that the data are complete <u>Dangerous</u>: "it lumps together situations where the problem is minor enough to be handled in this way & situations where estimators applied to the imputed data have substantial biases."

- Single imputation aims at completing data as best as possible
   ⇒ low rank approaches are powerful for heterogeneous data
- Multiple imputation aims at estimating the parameters and their variability taking into account the uncertainty of the missing values

<sup>&</sup>lt;sup>15</sup>Bogdan, J. et al. Adaptive Bayesian SLOPE - High dimensional Model Selection with Missing Values. Journal of Computational and Graphical Statistics. 2020.

• Methods used in practice are the one implemented in a sustainable way: few implementations of EM strategies

• "Imputation is both seductive & dangerous" (Dempster & Rubin, 1983). <u>Seductive</u>: "can lull the user into the pleasant state of believing that the data are complete <u>Dangerous</u>: "it lumps together situations where the problem is minor enough to be handled in this way & situations where estimators applied to the imputed data have substantial biases."

- Single imputation aims at completing data as best as possible
   ⇒ low rank approaches are powerful for heterogeneous data
- Multiple imputation aims at estimating the parameters and their variability taking into account the uncertainty of the missing values
  - How to aggregate lasso regressions? Alternatives EM<sup>15</sup>

<sup>&</sup>lt;sup>15</sup>Bogdan, J. et al. Adaptive Bayesian SLOPE - High dimensional Model Selection with Missing Values. Journal of Computational and Graphical Statistics. 2020.

# Challenges and on-going work with heterogeneous data sources and missing data



Sporadic & systematic (missing variable in one hospital). Due to the pandemic, many patients did not complete their tests

• What to do when you have both MCAR, MAR, MNAR in the data?

# Supervised learning with missing values

#### Prediction with missing values

 $\widetilde{X} = X \odot (1 - M) + \mathbb{NA} \odot M$ . New feature space is  $\widetilde{\mathbb{R}}^d = (\mathbb{R} \cup {\mathbb{NA}})^d$ .

$$\mathbf{Y} = \begin{pmatrix} 4.6\\ 7.9\\ 8.3\\ 4.6 \end{pmatrix} \quad \tilde{\mathbf{X}} = \begin{pmatrix} 9.1 & \text{NA} & 1\\ 2.1 & \text{NA} & 3\\ \text{NA} & 9.6 & 2\\ 4.2 & 5.5 & 6 \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} 9.1 & 8.5 & 1\\ 2.1 & 3.5 & 3\\ 6.7 & 9.6 & 2\\ 4.2 & 5.5 & 6 \end{pmatrix} \quad \mathbf{M} = \begin{pmatrix} 0 & 1 & 0\\ 0 & 1 & 0\\ 1 & 0 & 0\\ 0 & 0 & 0 \end{pmatrix}$$

Find a regression function that minimizes the expected risk

Bayes rule: 
$$f^* \in \underset{f: \tilde{\mathbb{R}}^d \to \mathbb{R}}{\arg \min} \mathbb{E}\left[\left(Y - f(\tilde{X})\right)^2\right]$$

$$f^{*}(\tilde{X}) = \mathbb{E}\left[Y \mid \tilde{X}\right] = \mathbb{E}\left[Y \mid X_{obs(M)}, M\right]$$
$$= \sum_{m \in \{0,1\}^{d}} \mathbb{E}\left[Y \mid X_{obs(m)}, M = m\right] \mathbb{1}_{M=m}$$

 $\Rightarrow$  One model per pattern *m* of missing values (2<sup>*d*</sup> patterns)<sup>16</sup>

 $^{16}$ Rosenbaum & Rubin. (1984). Reducing Bias in Observational Studies Using Subclassification on the Propensity Score. JASA.

#### Prediction with missing values

 $\widetilde{X} = X \odot (1 - M) + \mathbb{N} A \odot M$ . New feature space is  $\widetilde{\mathbb{R}}^d = (\mathbb{R} \cup {\mathbb{N}} A)^d$ .

$$\mathbf{Y} = \begin{pmatrix} 4.6\\ 7.9\\ 8.3\\ 4.6 \end{pmatrix} \quad \tilde{\mathbf{X}} = \begin{pmatrix} 9.1 & \text{NA} & 1\\ 2.1 & \text{NA} & 3\\ \text{NA} & 9.6 & 2\\ 4.2 & 5.5 & 6 \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} 9.1 & 8.5 & 1\\ 2.1 & 3.5 & 3\\ 6.7 & 9.6 & 2\\ 4.2 & 5.5 & 6 \end{pmatrix} \quad \mathbf{M} = \begin{pmatrix} 0 & 1 & 0\\ 0 & 1 & 0\\ 1 & 0 & 0\\ 0 & 0 & 0 \end{pmatrix}$$

Find a regression function that minimizes the expected risk

Bayes rule: 
$$f^* \in \underset{f: \ \widetilde{\mathbb{R}}^d \to \mathbb{R}}{\arg \min} \mathbb{E}\left[\left(Y - f(\tilde{X})\right)^2\right]$$

$$f^{*}(\tilde{X}) = \mathbb{E}\left[Y \mid \tilde{X}\right] = \mathbb{E}\left[Y \mid X_{obs(M)}, M\right]$$
$$= \sum_{m \in \{0,1\}^{d}} \mathbb{E}\left[Y \mid X_{obs(m)}, M = m\right] \mathbb{1}_{M=m}$$

 $\Rightarrow$  One model per pattern *m* of missing values (2<sup>*d*</sup> patterns)<sup>16</sup>

 $^{16}$ Rosenbaum & Rubin. (1984). Reducing Bias in Observational Studies Using Subclassification on the Propensity Score. JASA.
## Prediction with missing values

 $\widetilde{X} = X \odot (1 - M) + \mathbb{N} A \odot M$ . New feature space is  $\widetilde{\mathbb{R}}^d = (\mathbb{R} \cup {\mathbb{N}} A)^d$ .

$$\mathbf{Y} = \begin{pmatrix} 4.6\\ 7.9\\ 8.3\\ 4.6 \end{pmatrix} \quad \tilde{\mathbf{X}} = \begin{pmatrix} 9.1 & \text{NA} & 1\\ 2.1 & \text{NA} & 3\\ \text{NA} & 9.6 & 2\\ 4.2 & 5.5 & 6 \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} 9.1 & 8.5 & 1\\ 2.1 & 3.5 & 3\\ 6.7 & 9.6 & 2\\ 4.2 & 5.5 & 6 \end{pmatrix} \quad \mathbf{M} = \begin{pmatrix} 0 & 1 & 0\\ 0 & 1 & 0\\ 1 & 0 & 0\\ 0 & 0 & 0 \end{pmatrix}$$

Find a regression function that minimizes the expected risk

Bayes rule: 
$$f^* \in \underset{f: \ \widetilde{\mathbb{R}}^d \to \mathbb{R}}{\arg \min} \mathbb{E}\left[\left(Y - f(\tilde{X})\right)^2\right]$$

A learner estimates the regression function from a train set minimizing the empirical risk:  $\hat{f}_{\mathcal{D}_{n,\text{train}}} \in \underset{f: \widetilde{\mathbb{R}}^{d} \to \mathbb{R}}{\operatorname{arg\,min}} \left( \frac{1}{n} \sum_{i=1}^{n} \ell\left(f(\widetilde{X}_{i}), Y_{i}\right) \right)$ 

A new data  $\mathcal{D}_{n,\mathrm{test}}$  to estimate the generalization error rate

• Bayes consistent:  $\mathbb{E}[\ell(\hat{f}_n(\tilde{X}), Y)] \xrightarrow[n \to \infty]{} \mathbb{E}[\ell(f^{\star}(\tilde{X}), Y)]$ 

#### Differences with classical litterature

<u>Aim</u>: predict an outcome Y (not estimate parameters & their variance)

Specificities: both train & test sets with missing values; Otherwise, distributional shift (data generating process (X, Y, M))

 $\Rightarrow$  Is it possible to use previous approaches (EM - impute), consistent?  $\Rightarrow$  Do we need to design new ones?

#### Differences with classical litterature

<u>Aim</u>: predict an outcome Y (not estimate parameters & their variance) <u>Specificities</u>: both train & test sets with missing values; Otherwise, distributional shift (data generating process (X, Y, M))

 $\Rightarrow$  Is it possible to use previous approaches (EM - impute), consistent?  $\Rightarrow$  Do we need to design new ones?

#### Imputation prior to learning: Impute then Regress

Common practice: use off-the-shelf methods for

- 1) imputation of missing values
- 2) supervised-learning on the completed data

Impute train & test sets with the same model. Easy with univariate imputation: compute the means on the observed data  $(\hat{\mu}_1, ..., \hat{\mu}_d)$  of each column of the train set & impute the test set with such means

#### Framework - assumptions

- $$\begin{split} & \mathrel{\mathsf{Regression model:}} Y = f^*(X) + \varepsilon \\ & f^* : \mathbb{R}^d \to \mathbb{R} \text{ a continuous function of the complete data } X \\ & \mathrel{\varepsilon \in \mathbb{R}} \text{ is a centered random noise variable independent of } (X, M_1) \\ & X = (X_1, \dots, X_d) \text{ has a continuous density } g > 0 \text{ on } [0, 1]^d \\ & \|f^*\|_\infty = \sup_{x \in \mathbb{R}^d} |f^*(x)| < \infty \end{split}$$
- $\vdash \text{ Missing data: MAR on } X_1 \text{ with } M_1 \perp X_1 | X_2, \ldots, X_d$  $(x_2, \ldots, x_d) \mapsto \mathbb{P}[M_1 = 1 | X_2 = x_2, \ldots, X_d = x_d] \text{ is continuous}$

 $<sup>^{16}</sup>$ J. et al. Consistency of supervised learning with missing values. 2019.

## Constant (mean) imputation is consistent for prediction<sup>16</sup>

• Constant imputation  $x' = (x'_1, x_2, ..., x_d)$ :  $x'_1 = x_1 \mathbb{1}_{M_1=0} + \alpha \mathbb{1}_{M_1=1}$ 

• Use a **universally consistent algorithm** (for all distribution) to approach the regression function  $f_{impute}^{*}(x') = \mathbb{E}[Y|X = x']$ 

#### Theorem. (J. et al. 2019)

$$\begin{split} f^{\star}_{impute}(x') = & \mathbb{E}[Y|X_2 = x_2, \dots, X_d = x_d, M_1 = 1] \\ & \mathbb{1}_{x'_1 = \alpha} \mathbb{1}_{\mathbb{P}[M_1 = 1|X_2 = x_2, \dots, X_d = x_d] > 0} \\ & + & \mathbb{E}[Y|X = x'] \mathbb{1}_{x'_1 = \alpha} \mathbb{1}_{\mathbb{P}[M_1 = 1|X_2 = x_2, \dots, X_d = x_d] = 0} \\ & + & \mathbb{E}[Y|X = x', M_1 = 0] \mathbb{1}_{x'_1 \neq \alpha}. \end{split}$$

Prediction with constant is equal to the Bayes function almost everywhere

$$f^{\star}_{impute}(X') = f^{\star}(\tilde{X}) = \mathbb{E}[Y|\tilde{X}]$$

Rq: pointwise equality if using a constant out of range.

<sup>&</sup>lt;sup>16</sup>J. et al. Consistency of supervised learning with missing values. 2019.

## Consistency of constant imputation: Rationale

- ▷ Specific value, systematic like a code for missing
- ▷ The learner detects the code and recognizes it at the test time (the imputed data distribution shouldn't differ between train and test)
- > With categorical data, just code "Missing"
- ▷ With continuous data, any constant:
- > De-identified/imputed missing data: recovers from which pattern it comes
- ▷ Need a lot of data (asymptotic result) and a universally consistent learner



Imputing both train & test with the same constant and regress is consistent despite its drawbacks for estimation (useful in practice)

## Consistency of constant imputation: Rationale

- ▷ Specific value, systematic like a code for missing
- ▷ The learner detects the code and recognizes it at the test time (the imputed data distribution shouldn't differ between train and test)
- With categorical data, just code "Missing"
- With continuous data, any constant: out of range
- > De-identified/imputed missing data: recovers from which pattern it comes
- ▷ Need a lot of data (asymptotic result) and a universally consistent learner



Imputing both train & test with the same constant and regress is consistent despite its drawbacks for estimation (useful in practice)

## Bayes optimality of impute-n-regress<sup>17</sup>

• Imputation function:  $\forall m \in \{0,1\}^d$ , let  $\phi^{(m)} \in C_\infty$ :  $\mathbb{R}^{|obs(m)|} \to \mathbb{R}^{|mis(m)|}$ which outputs values for the missing entries based on the observed ones

$$\Phi: (\mathbb{R} \cup \{\mathbb{N}\mathbb{A}\})^d \to \mathbb{R}^d: \forall j \in \llbracket 1, d \rrbracket, \ \Phi_j(\widetilde{X}) = \begin{cases} X_j & \text{if } M_j = 0\\ \phi_j^{(M)}(X_{obs(M)}) & \text{if } M_j = 1 \end{cases}$$

• Regression on imputed data:  $g_{\Phi}^{\star} \in \underset{g:\mathbb{R}^d \mapsto \mathbb{R}}{\operatorname{argmin}} \mathbb{E}\left[\left(Y - g \circ \Phi(\widetilde{X})\right)^2\right]$ , minimizer of the risk on the imputed data

#### Theorem

Assume that the response Y satisfies  $Y = f^*(X) + \varepsilon$ Then, for all missing data mechanisms & almost all imputation functions,  $g_{\Phi}^* \circ \Phi$  is Bayes optimal

 $\Rightarrow$  A universally consistent algorithm trained on the imputed data  $\Phi(\widetilde{X})$  is Bayes consistent

#### Asymptotically, imputing well is not needed to predict well

<sup>&</sup>lt;sup>17</sup>Le Morvan, J. et al. What's a good imputation to predict with missing values? Neurips2021

## Rationale of proof: imputation creates manifolds



## Which imputation function should one choose?



Consistency of impute-then-regress. Ex: 3 regression models, 40% of MCAR in covariates, different imputation methods, then regress with random forests.

- A "better" imputation could create an easier learning problem
- Constant imputation is consistent but introduces strong discontinuities
- $\Rightarrow$  Which imputation and predictor should one use?

#### • Neumiss network:

- $\triangleright$  Classic network with multiplications by the mask nonlinearities  $\odot M$
- ▷ Motivated by linear regression with missing values in the covariates
- Theoritically grounded: approximation of the Bayes predictor (truncated neumiss series to approximate inverses of covariance matrices)
- Couple Neumiss and MLP to jointly learn imputation and regression



<sup>&</sup>lt;sup>18</sup>Le morvan, J. et al. Neumiss networks: differential programming for supervised learning with missing values. *Neurips2020 (Oral)*.

Supervised learning different from inferential aim

### Bayes optimality of Impute then Regress

- Single constant imputation is consistent with a universally consistent learner
- Can even work in MNAR
- Rethinking imputation: a good imputation is the one that makes the prediction easy

#### Implicit and jointly learned Impute-then-Regress strategy

- Neumiss network: new architecture  $\odot M$  nonlinearity
- Tree-based models: Missing Incorporated in Attribute (MIA) (implemented in generalized random forest package grf, partykit)

## Challenges and on-going works with missing values

 Uncertainty quantification with conformal prediction<sup>19</sup> (with Y. Romano & A. Dieuleveut)

Predictive intervals for any predictive algorithm (neural nets, random

forests), in finite samples with no assumption on the data distribution except for the exchangeability

<sup>&</sup>lt;sup>19</sup>Vovk, et al. Algorithmic Learning in a Random World. Springer US, 2005.

## Challenges and on-going works with missing values

 Uncertainty quantification with conformal prediction<sup>19</sup> (with Y. Romano & A. Dieuleveut)

Predictive intervals for any predictive algorithm (neural nets, random forests), in finite samples with no assumption on the data distribution except for the exchangeability

- ▷ Times series with MNAR (prediction with covariates measured regularly over time & static covariates)
- Federated learning with missing values

26

<sup>&</sup>lt;sup>19</sup>Vovk, et al. Algorithmic Learning in a Random World. Springer US, 2005.

<sup>&</sup>lt;sup>20</sup>Wager, J., Doubly robust estimation with incomplete confounders. Ann. Appl. Stat. 2020.

<sup>&</sup>lt;sup>21</sup>Colnet, J. et al. Generalization: sensitivity analysis & missing data. *Journal of causal inf.* 2022.

 $<sup>^{22}\</sup>mbox{Mayer}$  & J. Generalizing treatment effects with incomplete covariates. 2022.

<sup>&</sup>lt;sup>23</sup>Colnet, J. et al. Reweighting RCT for generalization: finite sample analysis & var. select. 2022.

## Challenges and on-going works with missing values

 Uncertainty quantification with conformal prediction<sup>19</sup> (with Y. Romano & A. Dieuleveut)

Predictive intervals for any predictive algorithm (neural nets, random forests), in finite samples with no assumption on the data distribution except for the exchangeability

- ▷ Times series with MNAR (prediction with covariates measured regularly over time & static covariates)
- > Federated learning with missing values
- Causal inference from incomplete heterogeneous sources,<sup>20</sup>,<sup>21</sup>,<sup>22</sup>,<sup>23</sup> (treatment estimation, generalisation of randomized control trial, etc.)

<sup>19</sup>Vovk, et al. Algorithmic Learning in a Random World. Springer US, 2005.

<sup>&</sup>lt;sup>20</sup>Wager, J., Doubly robust estimation with incomplete confounders. Ann. Appl. Stat. 2020.

 <sup>&</sup>lt;sup>21</sup>Colnet, J. et al. Generalization: sensitivity analysis & missing data. *Journal of causal inf.* 2022.
 <sup>22</sup>Maver & J. Generalizing treatment effects with incomplete covariates. 2022.

<sup>&</sup>lt;sup>23</sup>Colnet, J. et al. Reweighting RCT for generalization: finite sample analysis & var. select. 2022.

## Collaborators on missing values

- F. Husson, Prof. Agronomy University. (package missMDA, FactoMineR)
- Gosia Bogdan, Prof. Wroclaw. High dimensional regression
- Claire Boyer, Assoc. Prof. Sorbonne. Signal processing, missing values
- Aymeric Dieuleveut, Asso. Prof. Ecole Polytechnique, Paris. Optimization
- Imke Mayer, Postdoc Charité Institute, Berlin. Causal inference
- Aude Sportisse, Postdoc Inria Nice. Missing values
- Marine Le Morvan, Junior researcher at INRIA, Paris. Supervised learning
- Erwan Scornet, Asso. Prof. Ecole Polytechnique, Paris. Random forests
- Gael Varoquaux, Senior researcher at INRIA, Paris. ML, Scikit-learn
- Margaux Zaffran, PhD student, EDF. Conformal prediction



## Challenges with heterogeneous sources and missing data



## Monitor population & assess wetlands conservation policies

- National agency for wildlife and hunting management (ONCFS) data
- Contingency tables: Water (785 wetland sites) bird (23 species) count data, from 1990-2016 in 5 countries in North Africa
- Side information (17 variables) on sites & years: meteo, altitude, etc.

					Site	Year	Rain	Eco	Country	Agri	
Site	2008	2009	2010		1	2008	163.7	0.8	Algeria	16.2	
1	NA	0	0		2	2008	60.7	0.8	Algeria	16.2	
2	4	50	25	1.00	3	2008	227.9	0.8	Algeria	16.2	
3	NA	0	0		-	2008	174.8	0.8	Algeria	16.2	
4	NA	NA	NA	1	5	2008	163.7	0.8	Algeria	16.2	
5	NA	NA	NA	-	6	2008	230.7	0.8	Algeria	16.2	
6	0	0	0		7	2008	243.5	0.8	Algeria	16.2	
7	5	75	870		8	2008	262.6	0.8	Algeria	16.2	
8	9	34	0		9	2008	197.3	0.8	Algeria	16.2	
9	10	8	30		10	2008	227.9	0.8	Algeria	16.2	

#### Common pochard (canard milouin)

- $\Rightarrow$  Aims: Assess the effect of time on species abundances
- $\Rightarrow$  70% of missing values in contingency tables (drough, war, etc.)<sup>24,25</sup>

<sup>&</sup>lt;sup>24</sup> Robin, J., Moulines Sardy. 2019. Low-rank model with covariates for count data with missing values. *Journal of Multivariate Analysis.* 

<sup>&</sup>lt;sup>25</sup> Robin, Klopp, J., Moulines Tibshirani. Main effects and interactions in mixed and incomplete data frames. 2019. JASA.

## Iterative imputation by random forests versus by low rank (PCA)

	Feat1	Feat2	Feat3	Feat4	Feat5	Feat	1 Fe	at2 Feat3	Feat4	Feat5	Feat1	Feat2	Feat3	Feat4	Feat5
C1	1	1	1	1	1	1	1.0	1.00	1	1	1	1	1	1	1
C2	1	1	1	1	1	1	1.0	1.00	1	1	1	1	1	1	1
C3	2	2	2	2	2	2	2.0	2.00	2	2	2	2	2	2	2
C4	2	2	2	2	2	2	2.0	2.00	2	2	2	2	2	2	2
C5	3	3	3	3	3	3	3.0	3.00	3	3	3	3	3	3	3
C6	3	3	3	3	3	3	3.0	3.00	3	3	3	3	3	3	3
C7	4	4	4	4	4	4	4.0	4.00	4	4	4	4	4	4	4
C8	4	4	4	4	4	4	4.0	4.00	4	4	4	4	4	4	4
C9	5	5	5	5	5	5	5.0	5.00	5	5	5	5	5	5	5
C10	5	5	5	5	5	5	5.0	5.00	5	5	5	5	5	5	5
C11	6	6	6	6	6	6	6.0	6.00	6	6	6	6	6	6	6
C12	6	6	6	6	6	6	6.0	6.00	6	6	6	6	6	6	6
C13	7	7	7	7	7	7	7.0	7.00	7	7	7	7	7	7	7
C14	7	7	7	7	7	7	7.0	7.00	7	7	7	7	7	7	7
Igor	8	NA	NA	8	8	8	6.87	6.87	8	8	8	8	8	8	8
Frank	8	NA	NA	8	8	8	6.87	6.87	8	8	8	8	8	8	8
Bertrand	9	NA	NA	9	9	9	6.87	6.87	9	9	9	9	9	9	9
Alex	9	NA	NA	9	9	9	6.87	6.87	9	9	9	9	9	9	9
Yohann	10	NA	NA	10	10	10	6.87	6.87	10	10	10	10	10	10	10
Jean	10	NA	NA	10	10	10	6.87	6.87	10	10	10	10	10	10	10

#### Missing

missForest

imputePCA

 $\Rightarrow$  Imputation inherits from the method: Random forests (computationaly costly) handles non linear relationships/ PCA linear ones

## Bayes optimality of impute-n-regress (Le morvan et al. 2021)



#### Rationale: Imputation create manifolds to which the learner adapts

- 1. All data points with a missing data pattern m are mapped to a manifold  $\mathcal{M}^{(m)}$  of dimension |obs(m)| (Preimage Theorem)
- The missing data patterns of imputed data points can almost surely be de-identified (Thom transversality Theorem)<sup>26</sup>
- 3. Given 2), we can build prediction functions, independent of *m*, that are Bayes optimal for all missing data patterns

 $^{26}$ Non transverse: the manifolds on which the data with either x1 missing or x2 missing are projected are exactly the same (the same line)

# Jointly learn imputation & prediction: Neumiss

## Linear regression with missing values

#### Linear model:

$$Y = \beta_0 + \langle X, \beta \rangle + \varepsilon, \quad X \in \mathbb{R}^d, \ \varepsilon \text{ Gaussian}$$

Bayes predictor for the linear model:  

$$f^{*}(\tilde{X}) = \mathbb{E}[Y|\tilde{X}] = \mathbb{E}[\beta_{0} + \beta^{\mathsf{T}}X \mid M, X_{obs}(M)]$$

$$= \beta_{0} + \beta^{\mathsf{T}}_{obs}(M) X_{obs}(M) + \beta^{\mathsf{T}}_{mis}(M) \mathbb{E}[X_{mis}(M) \mid M, X_{obs}(M)]$$

$$= \sum_{m \in \{0,1\}^{d}} \beta_{0} + \beta^{\mathsf{T}}_{obs}(m) X_{obs}(m) + \beta^{\mathsf{T}}_{mis}(m) \mathbb{E}[X_{mis}(m) \mid M = m, X_{obs}(m)]$$

#### Assumptions on covariates and missing values (X, M)

1. Gaussian assumption  $X \sim \mathcal{N}(\mu, \Sigma)$  + MCAR and MAR

#### Under Assump. the Bayes predictor is linear per pattern

$$f^{\star}(X_{obs}, M) = \beta_0 + \langle \beta_{obs}, X_{obs} \rangle + \langle \beta_{mis}, \mu_{mis} + \Sigma_{mis,obs} (\Sigma_{obs})^{-1} (X_{obs} - \mu_{obs}) \rangle$$

use of obs instead of obs(M) for lighter notations - Expression for 2.

#### Example

Let  $Y = X_1 + X_2 + \varepsilon$ , where  $X_2 = \exp(X_1) + \varepsilon_1$ . Now, assume that only  $X_1$  is observed. Then, the model can be rewritten as

 $Y = X_1 + \exp(X_1) + \varepsilon + \varepsilon_1,$ 

where  $f(X_1) = X_1 + \exp(X_1)$  is the Bayes predictor. In this example, the submodel for which only  $X_1$  is observed is not linear.

 $\Rightarrow$  There exists a large variety of submodels for a same linear model. Depend on the structure of X and on the missing-value mechanism.

## Neumiss Networks to approximate the covariance matrix

Bayes predictor requires inverting many covariance matrices

$$f^{\star}(X_{obs}, M) = \beta_0^+ \langle \beta_{obs}, X_{obs} \rangle + \langle \beta_{mis}, \mu_{mis} + \Sigma_{mis,obs} (\Sigma_{obs})^{-1} (X_{obs} - \mu_{obs}) \rangle$$

Order- $\ell$  approx of  $(\sum_{obs(m)}^{-1})$  for any m defined recursively:

$$S_{obs(m)}^{(\ell)} = (Id - \Sigma_{obs(m)})S_{obs(m)}^{(\ell-1)} + Id.$$

Neuman Series,  $S^{(0)} = Id$ ,  $\ell = \infty$ :  $(\Sigma_{obs(m)})^{-1} = \sum_{k=0}^{\infty} (Id - \Sigma_{obs(m)})^k$ 

## Neumiss Networks to approximate the covariance matrix

#### Order $\ell$ approx. of the Bayes predictor)

 $f_{\ell}^{\star}(X_{obs}, M) = \langle \beta_{obs}, X_{obs} \rangle + \langle \beta_{mis}, \mu_{mis} + \Sigma_{mis,obs} \frac{S_{obs}^{(\ell)}}{S_{obs}^{obs}(m)} (X_{obs} - \mu_{obs}) \rangle$ 

Order- $\ell$  approx of  $(\sum_{obs(m)}^{-1})$  for any m defined recursively:

$$S_{obs(m)}^{(\ell)} = (Id - \Sigma_{obs(m)})S_{obs(m)}^{(\ell-1)} + Id.$$

Neuman Series,  $S^{(0)} = Id$ ,  $\ell = \infty$ :  $(\Sigma_{obs(m)})^{-1} = \sum_{k=0}^{\infty} (Id - \Sigma_{obs(m)})^k$ 

## Neumiss Networks to approximate the covariance matrix

#### Order $\ell$ approx. of the Bayes predictor)

 $f_{\ell}^{\star}(X_{obs}, M) = \langle \beta_{obs}, X_{obs} \rangle + \langle \beta_{mis}, \mu_{mis} + \Sigma_{mis,obs} S_{obs(m)}^{(\ell)}(X_{obs} - \mu_{obs}) \rangle$ 

Order- $\ell$  approx of  $(\sum_{obs(m)}^{-1})$  for any m defined recursively:

$$S_{obs(m)}^{(\ell)} = (Id - \Sigma_{obs(m)})S_{obs(m)}^{(\ell-1)} + Id.$$

Neuman Series,  $S^{(0)} = Id$ ,  $\ell = \infty$ :  $(\Sigma_{obs(m)})^{-1} = \sum_{k=0}^{\infty} (Id - \Sigma_{obs(m)})^k$ 

#### $\Rightarrow$ Neural network architecture to approximate the Bayes predictor



**Figure 1:** Depth of 3,  $\bar{m} = 1 - m$ . Each weight matrix  $W^{(k)}$  corresponds to a simple transformation of the covariance matrix indicated in blue.

- Implementing a network with the matrix weights  $W^{(k)} = (I \Sigma_{obs(m)})$ masked differently for each sample can be challenging
- Masked weights is equivalent to masking input & output vector. Let v a vector,  $\overline{m} = 1 - m$ .  $(W \odot \overline{m} \overline{m}^{\top})v = (W(v \odot \overline{m})) \odot \overline{m}$

Classic network with multiplications by the mask nonlinearities  $\odot M$ 



#### Jointly learn imputation and regression

<sup>&</sup>lt;sup>27</sup> Le morvan, J. et al. Neumiss networks: differential programming for supervised learning with missing values. *Neurips2020 (Oral)*.

• Oracles regression  $f^*$  & conditional imputation  $\mathbb{E}[X_{mis}|X_{obs}, M]$ :  $f^* \circ \Phi^{CI}$ 

#### Proposition (excess of risk)

Assum PSD matrices  $\overline{H}^+$  &  $\overline{H}^-$  s.t. for all  $X \in S$ ,  $\overline{H}^- \leq H(X) \leq \overline{H}^+$ , H(X) the Hessian of  $f^*$  at X (min. & max. curvatures of  $f^*$  in any direction are uniformly bounded over the entire space)

$$\mathcal{R}\left(f^{\star}\circ\Phi^{\textit{Cl}}\right)-\mathcal{R}^{\star}\leq \frac{1}{4}\mathbb{E}_{\textit{M}}[\mathsf{max}\left(\mathsf{tr}\left(\bar{H}^{-}_{\textit{mis},\textit{mis}}\boldsymbol{\Sigma}_{\textit{mis}|\textit{obs},\textit{M}}\right)^{2},\mathsf{tr}\left(\bar{H}^{+}_{\textit{mis},\textit{mis}}\boldsymbol{\Sigma}_{\textit{mis}|\textit{obs},\textit{M}}\right)^{2}\right)]$$

High excess risk if both 1) the curvature of  $f^*$  is high and 2) the variance of the missing data given the observed one is high (linear regression consistent)

• Oracles regression  $f^*$  & conditional imputation  $\mathbb{E}[X_{mis}|X_{obs}, M]$ :  $f^* \circ \Phi^{CI}$ 

#### Proposition (excess of risk)

Assum PSD matrices  $\overline{H}^+$  &  $\overline{H}^-$  s.t. for all  $X \in S$ ,  $\overline{H}^- \leq H(X) \leq \overline{H}^+$ , H(X) the Hessian of  $f^*$  at X (min. & max. curvatures of  $f^*$  in any direction are uniformly bounded over the entire space)

$$\mathcal{R}\left(f^{\star}\circ\Phi^{\textit{Cl}}\right)-\mathcal{R}^{\star}\leq \frac{1}{4}\mathbb{E}_{\textit{M}}[\mathsf{max}\left(\mathsf{tr}\left(\bar{H}^{-}_{\textit{mis},\textit{mis}}\boldsymbol{\Sigma}_{\textit{mis}|\textit{obs},\textit{M}}\right)^{2},\mathsf{tr}\left(\bar{H}^{+}_{\textit{mis},\textit{mis}}\boldsymbol{\Sigma}_{\textit{mis}|\textit{obs},\textit{M}}\right)^{2}\right)]$$

High excess risk if both 1) the curvature of  $f^*$  is high and 2) the variance of the missing data given the observed one is high (linear regression consistent)

• Is there a <u>continuous</u> function g, s.t.  $g \circ \Phi^{CI}$  is Bayes optimal? *No*.

## CART (Breiman, 1984)

Built recursively by splitting the current cell into two children: Find the feature  $j^*$ , the threshold  $z^*$  which minimises the (quadratic) loss

$$(j^{\star}, z^{\star}) \in \underset{(j,z)\in\mathcal{S}}{\operatorname{arg\,min}} \mathbb{E}\Big[ \left(Y - \mathbb{E}[Y|X_j \le z]\right)^2 \cdot \mathbb{1}_{X_j \le z} + \left(Y - \mathbb{E}[Y|X_j > z]\right)^2 \cdot \mathbb{1}_{X_j > z}\Big].$$



root

## CART (Breiman, 1984)

Built recursively by splitting the current cell into two children: Find the feature  $j^*$ , the threshold  $z^*$  which minimises the (quadratic) loss

$$(j^{\star}, z^{\star}) \in \underset{(j,z)\in\mathcal{S}}{\operatorname{arg\,min}} \mathbb{E}\Big[ \left(Y - \mathbb{E}[Y|X_j \leq z]\right)^2 \cdot \mathbb{1}_{X_j \leq z} + \left(Y - \mathbb{E}[Y|X_j > z]\right)^2 \cdot \mathbb{1}_{X_j > z} \Big].$$



## CART (Breiman, 1984)

Built recursively by splitting the current cell into two children: Find the feature  $j^*$ , the threshold  $z^*$  which minimises the (quadratic) loss

$$(j^{\star}, z^{\star}) \in \underset{(j,z)\in\mathcal{S}}{\operatorname{arg\,min}} \mathbb{E}\Big[ \left(Y - \mathbb{E}[Y|X_j \leq z]\right)^2 \cdot \mathbb{1}_{X_j \leq z} + \left(Y - \mathbb{E}[Y|X_j > z]\right)^2 \cdot \mathbb{1}_{X_j > z}\Big].$$



## CART with missing values

root

	$X_1$	<i>X</i> <sub>2</sub>	Y
1			
2	NA		
3	NA		
4			

## CART with missing values



1) Select variable and threshold on observed values (1 & 4 for  $X_1$ )  $\mathbb{E}\Big[(Y - \mathbb{E}[Y|X_j \le z, M_j = 0])^2 \cdot \mathbb{1}_{X_j \le z, M_j = 0} + (Y - \mathbb{E}[Y|X_j > z, M_j = 0])^2 \cdot \mathbb{1}_{X_j > z, M_j = 0}\Big].$ 

## CART with missing values



1) Select variable and threshold on observed values (1 & 4 for  $X_1$ )  $\mathbb{E}\Big[(Y - \mathbb{E}[Y|X_j \le z, M_j = 0])^2 \cdot \mathbb{1}_{X_j \le z, M_j = 0} + (Y - \mathbb{E}[Y|X_j > z, M_j = 0])^2 \cdot \mathbb{1}_{X_j > z, M_j = 0}\Big].$ 

2) Propagate observations (2 & 3) with missing values?

- Probabilistic split:  $Bernoulli(\frac{\#L}{\#L+\#R})$  (Rweeka)
- Block: Send all to a side by minimizing the error (xgboost, lightgbm)
- Surrogate split: Search another variable that gives a close partition (rpart)

One step: select the variable, the threshold and propagate missing values

1. 
$$\{\widetilde{X}_j \leq z \text{ or } \widetilde{X}_j = \mathbb{N}A\} \text{ vs } \{\widetilde{X}_j > z\}$$
  
2.  $\{\widetilde{X}_j \leq z\} \text{ vs } \{\widetilde{X}_j > z \text{ or } \widetilde{X}_j = \mathbb{N}A\}$   
3.  $\{\widetilde{X}_i \neq \mathbb{N}A\} \text{ vs } \{\widetilde{X}_i = \mathbb{N}A\}.$ 

- $\triangleright$  The splitting location z depends on the missing values
- $\triangleright$  Missing values treated like a category (well to handle  $\mathbb{R} \cup \mathtt{NA}$ )
- $\triangleright$  Good for informative pattern (*M* explains *Y*)

Targets one model per pattern:

$$\mathbb{E}\left[Y\Big|\tilde{X}\right] = \sum_{m \in \{0,1\}^d} \mathbb{E}\left[Y|X_{obs(m)}, M = m\right] \mathbb{1}_{M=m}$$

▷ Implementation<sup>28</sup>: grf package, scikit-learn, partykit

#### $\Rightarrow$ Extremely good performances in practice for any mechanism

 $<sup>^{28}</sup>$  implementation trick, J. Tibshirani, duplicate the incomplete columns, and replace the missing entries once by  $+\infty$  and once by  $-\infty$
## Causal inference from (incomplete) heterogeneous sources

• Estimate causal effect: Example on trauma brain patients<sup>29</sup> "tranexamic acid" (treatment) impact on "28 days mortality" (outcome)

<u>Estimator</u>: Augmented Inverse Propensity Weighting uses non-parametric regression models Treatment  $\sim$  covariates & Outcome  $\sim$  covariates  $\Rightarrow$  Extended to handle missing values (implemented in grf package)

<sup>&</sup>lt;sup>29</sup>Wager, J., Doubly robust estimation with incomplete confounders. Ann. Appl. Stat. 2020.

# Causal inference from (incomplete) heterogeneous sources

• Estimate causal effect: Example on trauma brain patients<sup>29</sup> "tranexamic acid" (treatment) impact on "28 days mortality" (outcome)

<u>Estimator</u>: Augmented Inverse Propensity Weighting uses non-parametric regression models Treatment  $\sim$  covariates & Outcome  $\sim$  covariates  $\Rightarrow$  Extended to handle missing values (implemented in grf package)

• Generalization of randomized control trial's findings toward a target pop., with distributional shift (data fusion, recovery from selection biais)<sup>30,31,32,33</sup>

			Covariates			Treat	Outcomes
	Set	S	<i>x</i> <sub>1</sub>	X2	$X_3$	W	Y
1	RCT	1	1.1	20	NA	1	24.1
	RCT	1	-6	NA	NA	0	26.3
п	RCT	1	0	15	NA	1	23.5
n + 1	Target	?	NA	35	7.1		
n + 2	Target	?	-2	52	2.4		
	Target	?					
n + m	Target	?	-2	22	NA		

<sup>29</sup>Wager, J., Doubly robust estimation with incomplete confounders. Ann. Appl. Stat. 2020.
 <sup>30</sup>Colnet, J. et al. (2021). Causal inference for combining RCT & obs. studies: a review.
 <sup>31</sup>Mayer & J. Generalizing treatment effects with incomplete covariates. 2022.
 <sup>32</sup>Colnet, J. et al. Generalization: sensitivity analysis & missing data. Journal of causal inf. 2022.
 <sup>33</sup>Colnet, J. et al. Reweighting RCT for generalization: finite sample analysis & var. select. 2022.

## Traumabase project: decision support for trauma patients

- ▷ 30 000 French trauma patients<sup>34</sup>
- ▷ 250 features from the accident site to the hospital discharge
- ▷ 30 hospitals
- $\triangleright$  4000 new patients/ year

Center	Accident	Age	Sex	Weight	Lactactes	BP	TXA.	Y
Beaujon	fall	54	m	85	NM	180	treated	0
Pitie	gun	26	m	NR	NA	131	untreated	1
Beaujon	moto	63	m	80	3.9	145	treated	1
Pitie	moto	30	W	NR	Imp	107	untreated	0
HEGP	knife	16	m	98	2.5	118	treated	1
:								•.

 $\Rightarrow$  Estimate causal effect: Administration of the treatment "tranexamic acid (TXA)" given within 3 hours of the accident, on the outcome 28 days intra hospital mortality (Y) for trauma brain patients.<sup>35</sup>

<sup>&</sup>lt;sup>34</sup>www.traumabase.eu - https://www.traumatrix.fr/

<sup>&</sup>lt;sup>35</sup>Mayer, Wager, J. Doubly robust treatment effect estimation with incomplete confounders. *Annals Of Applied Statistics*. 2020.

### Effect of tranexamic acid on in-ICU mortality

Model Treatment on Covariates  $e(x) = \mathbb{P}(W_i = 1 | X_i = x)$ Model Outcome on Covariates  $\mu_{(w)}(x) = \mathbb{E}[Y_i(w) | X_i = x]$ 

Augmented Inverse Propensity Weighting - double robust

$$\hat{\tau}_{AIPW} = rac{1}{n} \sum_{i=1}^{n} \left( \hat{\mu}_{(1)}(X_i) - \hat{\mu}_{(0)}(X_i) + W_i rac{Y_i - \hat{\mu}_{(1)}(X_i)}{\hat{e}(X_i)} - (1 - W_i) rac{Y_i - \hat{\mu}_{(0)}(X_i)}{1 - \hat{e}(X_i)} 
ight)$$

 $\hat{\tau}_{AIPW} \text{ is } \sqrt{n} \text{-consistent, asympt. normal with semi parametric variance given:} \\ \mathbb{E} \left[ \left( \hat{e}(X_i)^{(-i)} - e(X_i) \right)^2 \right]^{\frac{1}{2}} \times \mathbb{E} \left[ \left( \hat{\mu}_{(W)} \left( X_i \right)^{(-i)} - \mu_{(W)} \left( X_i \right) \right)^2 \right]^{\frac{1}{2}} = o \left( \frac{1}{\sqrt{n}} \right)$ 



x-axis: Estimat. of the Average Treatment Effect ( $\times 100$ ), bootstrap CI

AIPW with missing implemented in generalized random forest package grf

 $\Rightarrow$  RCT gold standart to estimate treatment effect

Trial sample can be different from the population eligible for treatment

 $\Rightarrow$  Leveraging RCT and covariates from target population to transport the treatment effect estimated from the RCT to another population with a distributional shift (data fusion, recovery from selection biais),<sup>36</sup>,<sup>37,38</sup>

- ightarrow Reduce drug approval times and costs for patients who could benefit
- $\rightarrow$  Prices depend on efficiency

			Covariates			Treat	Outcomes
	Set	S	$X_1$	X2	$X_3$	W	Y
1	$\mathcal{RCT}$	1	1.1	20	NA	1	24.1
	$\mathcal{RCT}$	1	-6	45	NA	0	26.3
п	$\mathcal{RCT}$	1	0	15	NA	1	23.5
n + 1	Obs	?	-1	35	7.1		
n + 2	Obs	?	-2	52	2.4		
	OLS	?					
n + m	Obs	?	-2	22	3.4		

 $<sup>^{36}</sup>$  Mayer & J. Generalizing treatment effects with incomplete covariates. Archiv. 2022.

<sup>&</sup>lt;sup>37</sup>Colnet, J. et al. Generalizing a causal effect: sensitivity analysis and missing covariates. *Journal of causal inference*. 2022.

<sup>&</sup>lt;sup>38</sup>Colnet, J. et al. Reweighting RCT for generalization: finite sample analysis & var. select. 2022.

<sup>&</sup>lt;sup>39</sup> Colnet, J. et al. (2021). Causal inference for combining RCT & obs. studies: a review.

### MNAR data: identifiability issues, few solutions in practice

Before estimation, we should prove the identifiability of the parameters Example: Credit: Ilya Shpitser  $X^{NA} = [1, NA, 0, 1, NA, 0]$ 

▷ **Case 1:** X missing only if X = 1.

 $X = [1, 1, 0, 1, 1, 0], \mathbb{P}(X = 1) = 2/3$ 

 $\triangleright$  Case 2: X missing only if X = 0.

 $X = [1, 0, 0, 1, 0, 0], \mathbb{P}(X = 1) = 1/3$ 

⇒ Start from 2 equal observed distribution. It leads to different parameters of the data distribution  $\mathbb{P}(X = 1)$ <u>Identifiability</u>: the parameters of (X, M) are uniquely determined from available information (X, M = 0)

Estimation: restrictive setting (few variables, only missing values on the outcome, simple models)  $^{40\,41\,42}$ 

<sup>&</sup>lt;sup>40</sup>Ibrahim, et al. Missing covariates in glm when the mechanism is non-ignorable. JRSSB. 1999.
<sup>41</sup>Tang. Statistical inference for nonignorable missing-data. Statistic. theory & rel. fields. 2018.
<sup>42</sup>Mohan, Thoemmes, Pearl. Estimation with incomplete data: The linear case. IJCAI. 2018.

#### Notations

#### • Random Variables:

 $\triangleright \ X \in \mathbb{R}^d$ : the complete unvailable data

 $\triangleright \ \widetilde{X} \in \{\mathbb{R} \cup \{\mathbb{N}A\}\}^d$ : incomplete data (observed), NA: Not Available

 $\triangleright \ M \in \{0,1\}^d$ : the missing-data pattern, the mask

obs(M) (resp. mis(M)) indices of the observed (resp. missing) entries.

• Realizations:

$$\begin{aligned} x &= (1.1, 2.3, 3.1, 8, 5.27) \\ \widetilde{x} &= (1.1, \text{NA}, -3.1, 8, \text{NA}) \\ m &= (0, 1, 0, 0, 1) \\ x_{\text{obs}(m)} &= (1.1, 3.1, 8), \\ x_{\text{mis}(m)} &= (2.3, 5.27) \end{aligned}$$

**MCAR**<sup>43</sup>: For all  $m \in \{0, 1\}^d$ ,  $P(M = m \mid X) = P(M = m)$ **MAR**<sup>44</sup>: For all  $m \in \{0, 1\}^d$ ,  $P(M = m \mid X) = P(M = m \mid X_{obs(m)})$ 

 $<sup>^{43}</sup>$ Michel, Naf, Spohn, Â<sup>°</sup> Meinshausen. 2021. PKLM: a flexible mcar test using classification.  $^{44}$ What Is Meant by "Missing at Random"? Seaman, et al. Statistical Science. 2013.