

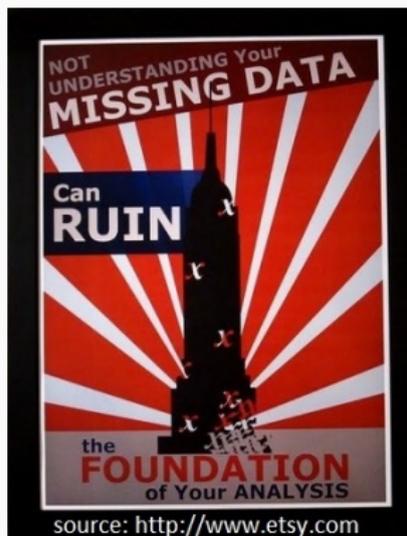
A missing value tour in R

Julie Josse

Ecole Polytechnique, INRIA

10 July 2019

useR!2019, Toulouse



1. Introduction
2. Handling missing values (inferential framework)
3. Supervised learning with missing values
4. Discussion - challenges

Introduction

Collaborators

- PhD students: W. Jiang, I. Mayer, N. Prost, G. Robin, A. Sportisse
- Colleagues: C. Boyer (LPSM), G. Bogdan (Wroclaw), F. Husson (Agrocampus) - (package [missMDA](#)), J-P Nadal (EHESS), E. Scornet (X), G. Varoquaux (INRIA), S. Wager (Stanford)
- Traumabase (hospital): T. Gauss, S. Hamada, J-D Moyer/ Capgemini



Traumabase

- 20000 patients
- 250 continuous and categorical variables: **heterogeneous**
- 11 hospitals: **multilevel data**
- 4000 new patients/ year

Center	Accident	Age	Sex	Weight	Lactactes	BP	shock	...
Beaujon	fall	54	m	85	NM	180	yes	
Pitie	gun	26	m	NR	NA	131	no	
Beaujon	moto	63	m	80	3.9	145	yes	
Pitie	moto	30	w	NR	Imp	107	no	
HEGP	knife	16	m	98	2.5	118	no	
⋮								⋮

Traumabase

- 20000 patients
- 250 continuous and categorical variables: **heterogeneous**
- 11 hospitals: **multilevel data**
- 4000 new patients/ year

Center	Accident	Age	Sex	Weight	Lactactes	BP	shock	...
Beaujon	fall	54	m	85	NM	180	yes	
Pitie	gun	26	m	NR	NA	131	no	
Beaujon	moto	63	m	80	3.9	145	yes	
Pitie	moto	30	w	NR	Imp	107	no	
HEGP	knife	16	m	98	2.5	118	no	
								...

⇒ **Estimate causal effect:** Administration of the **treatment** "tranexamic acid" (within 3 hours after the accident) on the **outcome** mortality for traumatic brain injury patients

Traumabase

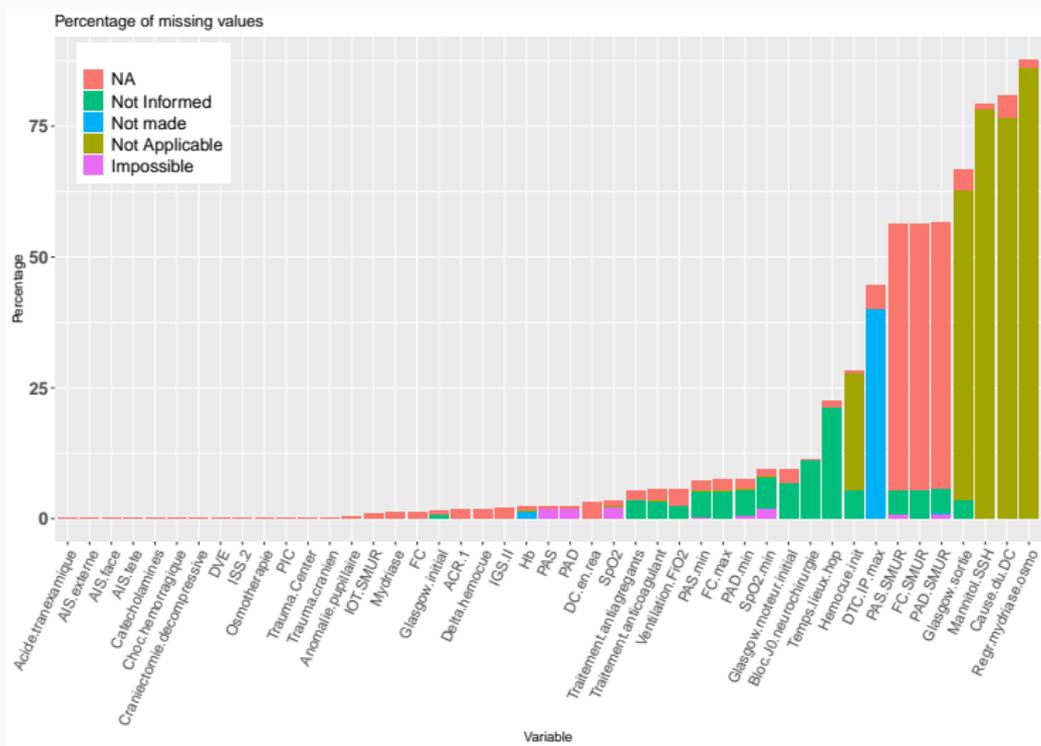
- 20000 patients
- 250 continuous and categorical variables: **heterogeneous**
- 11 hospitals: **multilevel data**
- 4000 new patients/ year

Center	Accident	Age	Sex	Weight	Lactactes	BP	shock	...
Beaujon	fall	54	m	85	NM	180	yes	
Pitie	gun	26	m	NR	NA	131	no	
Beaujon	moto	63	m	80	3.9	145	yes	
Pitie	moto	30	w	NR	Imp	107	no	
HEGP	knife	16	m	98	2.5	118	no	
⋮								⋮

⇒ **Predict** the risk of hemorrhagic shock given pre-hospital features

Ex random forests/logistic regression with covariates with missing values

Missing values

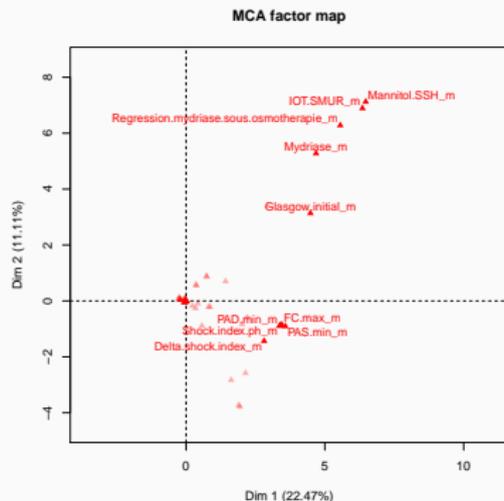
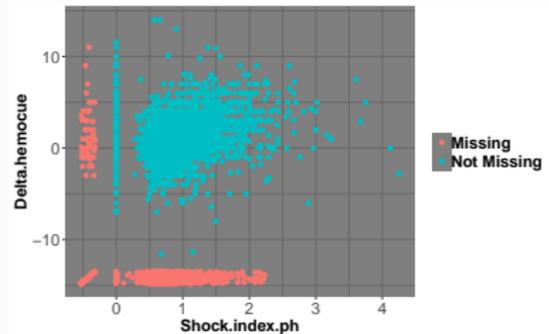


Multilevel data/ data integration: Systematic missing variable in one hospital

Visualization

The first thing to do with missing values (as for any analysis) is descriptive statistics: Visualization of patterns to get hints on how and why they occur

VIM (M. Templ), [naniar](#) (N. Tierney), [FactoMineR](#) (Husson *et al.*)



Right: *PAS_m* close to *PAD_m*: Often missing on both *PAS* & *PAD*

IOT: nested questions. Q1: yes/no, if yes Q2 - Q4, if no Q2 - Q4 "missing"

Note: Crucial **before** starting any treatment of missing values and **after**

Handling missing values (inferential framework)

Solutions to handle missing values

Books: Schafer (2002), Little & Rubin (2002); Kim & Shao (2013); Carpenter & Kenward (2013); van Buuren (2018), etc.

Modify the estimation process to deal with missing values

Maximum likelihood: **EM algorithm** to obtain point estimates +
Supplemented EM (Meng & Rubin, 1991) / Louis formulae for their variability
Ex logistic regression: EM to get $\hat{\beta}$ + Louis to get $\hat{V}(\hat{\beta})$

Aim: **Estimate parameters & their variance** from an incomplete data
⇒ Inferential framework

Solutions to handle missing values

Books: Schafer (2002), Little & Rubin (2002); Kim & Shao (2013); Carpenter & Kenward (2013); van Buuren (2018), etc.

Modify the estimation process to deal with missing values

Maximum likelihood: **EM algorithm** to obtain point estimates +
Supplemented EM (Meng & Rubin, 1991) / Louis formulae for their variability
Ex logistic regression: EM to get $\hat{\beta}$ + Louis to get $\hat{V}(\hat{\beta})$

Cons: Difficult to establish - not many softwares even for simple models
One specific algorithm for each statistical method...

Aim: **Estimate parameters & their variance** from an incomplete data
⇒ Inferential framework

Solutions to handle missing values

Books: Schafer (2002), Little & Rubin (2002); Kim & Shao (2013); Carpenter & Kenward (2013); van Buuren (2018), etc.

Modify the estimation process to deal with missing values

Maximum likelihood: **EM algorithm** to obtain point estimates +
Supplemented EM (Meng & Rubin, 1991) / Louis formulae for their variability
Ex logistic regression: EM to get $\hat{\beta}$ + Louis to get $\hat{V}(\hat{\beta})$

Cons: Difficult to establish - not many softwares even for simple models
One specific algorithm for each statistical method...

Imputation (multiple) to get a complete data set

Any analysis can be performed

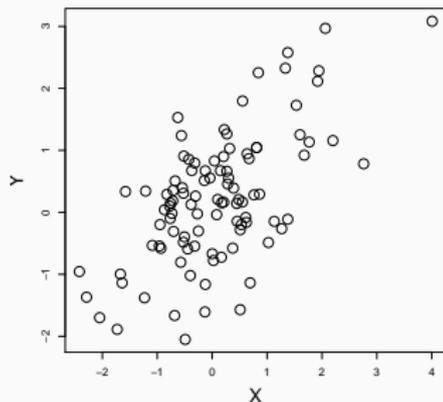
Ex logistic regression: Impute and apply logistic model to get $\hat{\beta}$, $\hat{V}(\hat{\beta})$

Aim: **Estimate parameters & their variance** from an incomplete data
⇒ Inferential framework

Mean imputation

- $(x_i, y_i) \underset{\text{i.i.d.}}{\sim} \mathcal{N}_2((\mu_x, \mu_y), \Sigma_{xy})$

X	Y
-0.56	-1.93
-0.86	-1.50
.....	...
2.16	0.7
0.16	0.74



$$\mu_y = 0$$

$$\sigma_y = 1$$

$$\rho = 0.6$$

$\hat{\mu}_y = -0.01$

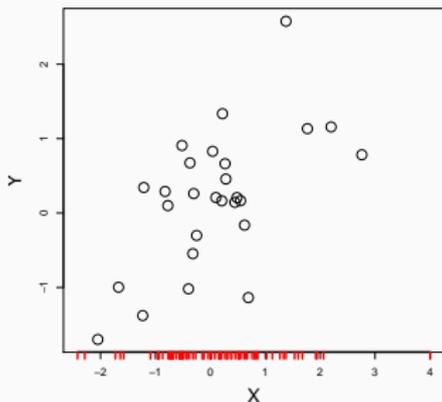
$\hat{\sigma}_y = 1.01$

$\hat{\rho} = 0.66$

Mean imputation

- $(x_i, y_i) \underset{\text{i.i.d.}}{\sim} \mathcal{N}_2((\mu_x, \mu_y), \Sigma_{xy})$
- 70 % of missing entries completely at random on Y

X	Y
-0.56	NA
-0.86	NA
.....	...
2.16	0.7
0.16	NA

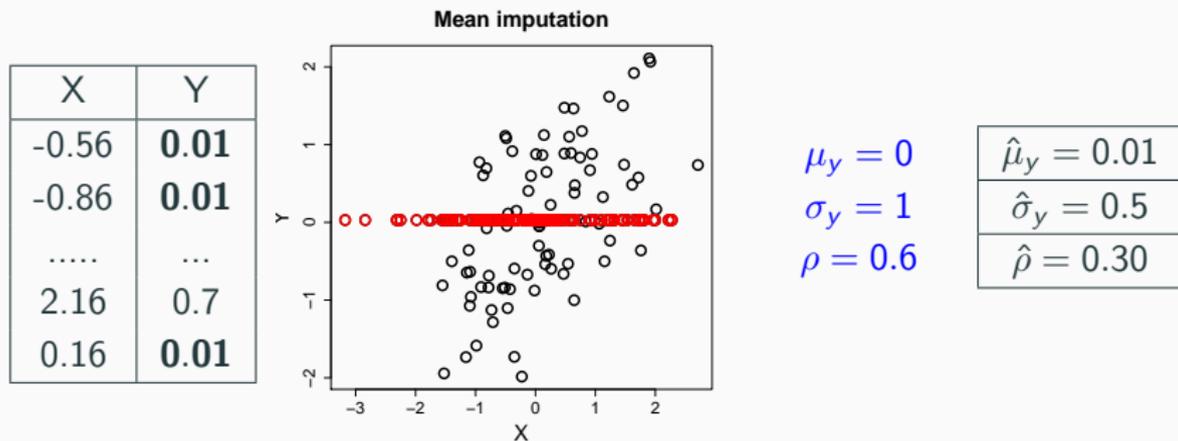


$$\begin{aligned}\mu_y &= 0 \\ \sigma_y &= 1 \\ \rho &= 0.6\end{aligned}$$

$\hat{\mu}_y = 0.18$
$\hat{\sigma}_y = 0.9$
$\hat{\rho} = 0.6$

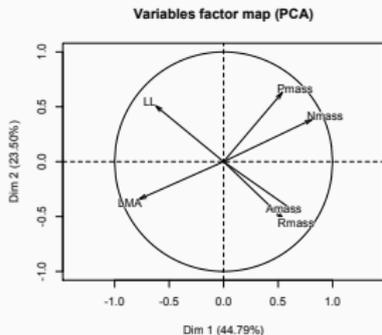
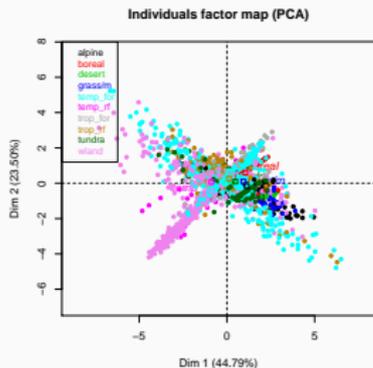
Mean imputation

- $(x_i, y_i) \underset{\text{i.i.d.}}{\sim} \mathcal{N}_2((\mu_x, \mu_y), \Sigma_{xy})$
- 70 % of missing entries completely at random on Y
- Estimate parameters on the mean imputed data

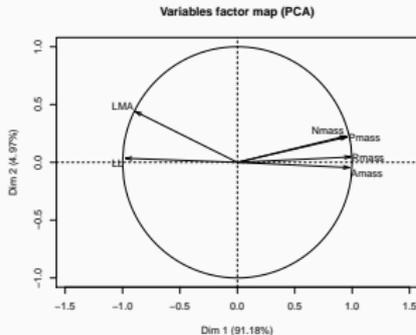
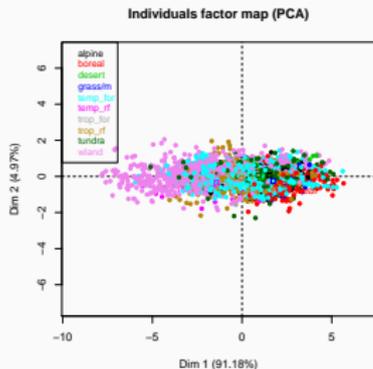


Mean imputation deforms joint and marginal distributions

Mean imputation is bad for estimation



```
library(FactoMineR)
PCA(eco1)
Warning message: Missing
are imputed by the mean
of the variable:
You should use imputePCA
from missMDA
```



```
library(missMDA)
imp <- imputePCA(eco1)
PCA(imp$comp)
```

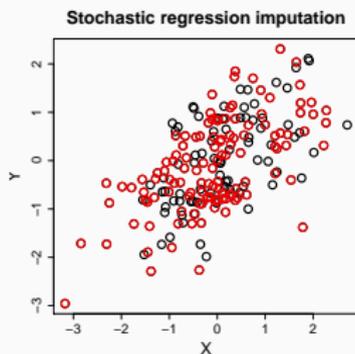
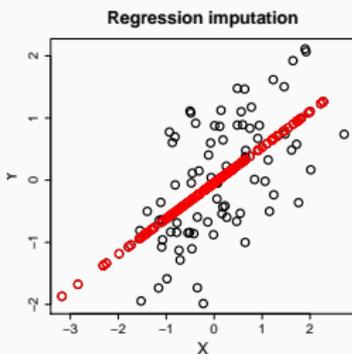
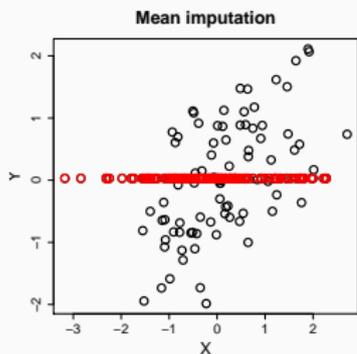
Ecological data: ¹ $n = 69000$ species - 6 traits. Estimated correlation between P_{mass} & $R_{mass} \approx 0$ (mean imputation) or ≈ 1 (EM PCA)

¹Wright, I. et al. (2004). The worldwide leaf economics spectrum. *Nature*.

Imputation methods

- by regression takes into account the relationship: Estimate β - impute $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \Rightarrow$ variance underestimated and correlation overestimated
- by stochastic reg: Estimate β and σ - impute from the predictive $y_i \sim \mathcal{N}(x_i \hat{\beta}, \hat{\sigma}^2) \Rightarrow$ preserve distributions

Here $\hat{\beta}, \hat{\sigma}^2$ estimated with complete data, but MLE can be obtained with EM



$$\mu_y = 0$$

$$\sigma_y = 1$$

$$\rho = 0.6$$

0.01
0.5
0.30

0.01
0.72
0.78

0.01
0.99
0.59

Imputation methods for multivariate data

Assuming a joint model

- Gaussian distribution: $x_j \sim \mathcal{N}(\mu, \Sigma)$ ([Amelia](#) Honaker, King, Blackwell)
- low rank: $X_{n \times d} = \mu_{n \times d} + \varepsilon \varepsilon_{ij}^{\text{iid}} \sim \mathcal{N}(0, \sigma^2)$ with μ of low rank k ([softimpute](#) Hastie & Mazuder; [missMDA](#) J. & Husson)
- latent class - nonparametric Bayesian ([dpmpm](#) Reiter)
- deep learning using variational autoencoders (MIWAE, Mattei, 2018)

Using conditional models (joint implicitly defined)

- with logistic, multinomial, poisson regressions ([mice](#) van Buuren)
- iterative impute each variable by random forests ([missForest](#) Stekhoven)

Imputation for categorical, mixed, blocks/multilevel data ², etc.

⇒ [Missing values taskview](#)³ J., Mayer., Tierney, Vialaneix

²J., Husson, Robin & Narasimhan. (2018). Imputation of mixed data with multilevel SVD.

³<https://cran.r-project.org/web/views/MissingData.html>

Random forests versus PCA

	Feat1	Feat2	Feat3	Feat4	Feat5...	Feat1	Feat2	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

⇒ Imputation inherits from the method: RF (computationally costly)
good for non linear relationships / PCA good for linear relationships

Single imputation: Underestimation of the variability

⇒ Incomplete Traumbase

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

Single imputation: Underestimation of the variability

⇒ 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

⇒ 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

Single imputation: Underestimation of the variability

⇒ 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

⇒ 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	X ₃	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	X ₃	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

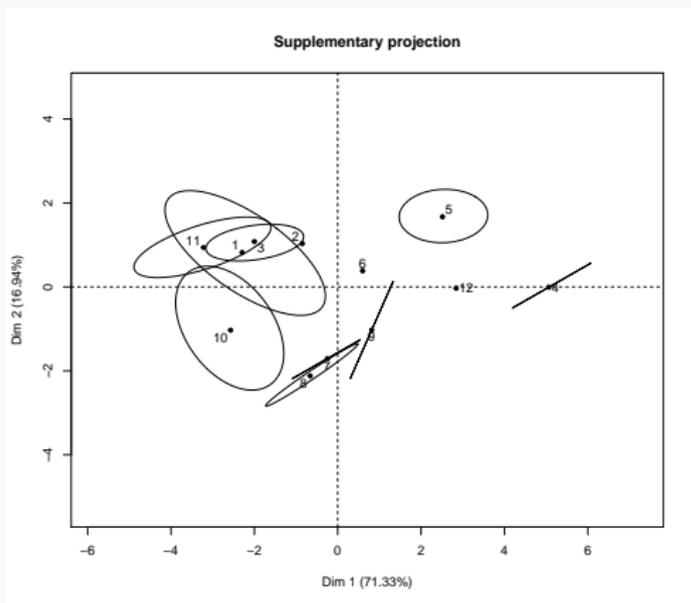
```
library(mice); mice(traumadata)
library(missMDA); MIPCA(traumadata)
```

Visualization of the imputed values

X ₁	X ₂	X ₃	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

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	X ₃	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



library(missMDA)
MIPCA(traumadata)
library(Amelia)
?compare.density

Percentage of NA?

Multiple imputation

1) Generate M plausible values for each missing value

X ₁	X ₂	X ₃	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

X ₁	X ₂	X ₃	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

X ₁	X ₂	X ₃	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: $\hat{\beta}_m, \widehat{Var}(\hat{\beta}_m)$

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

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

$$T = \frac{1}{M} \sum_{m=1}^M \widehat{Var}(\hat{\beta}_m) + \left(1 + \frac{1}{M}\right) \frac{1}{M-1} \sum_{m=1}^M (\hat{\beta}_m - \hat{\beta})^2$$

```
imp.mice <- mice(traumadata)
lm.mice.out <- with(imp.mice, glm(Y ~ ., family = "binomial"))
```

⇒ Variability of missing values taken into account

Logistic regression with missing covariates: Parameter estimation, model selection and prediction (Jiang, J., Lavielle, Gauss, Hamada, 2018)

$x = (x_{ij})$ a $n \times d$ matrix of quantitative covariates

$y = (y_i)$ an n -vector of binary responses $\{0, 1\}$

Logistic regression model: $\mathbb{P}(y_i = 1|x_i; \beta) = \frac{\exp(\beta_0 + \sum_{j=1}^d \beta_j x_{ij})}{1 + \exp(\beta_0 + \sum_{j=1}^d \beta_j x_{ij})}$

Covariables: $x_i \underset{\text{i.i.d.}}{\sim} \mathcal{N}_d(\mu, \Sigma)$

Log-likelihood with $\theta = (\mu, \Sigma, \beta)$:

$$\mathcal{LL}(\theta; x, y) = \sum_{i=1}^n \left(\log(p(y_i|x_i; \beta)) + \log(p(x_i; \mu, \Sigma)) \right).$$

Logistic regression with missing covariates: Parameter estimation, model selection and prediction (Jiang, J., Lavielle, Gauss, Hamada, 2018)

$x = (x_{ij})$ a $n \times d$ matrix of quantitative covariates

$y = (y_i)$ an n -vector of binary responses $\{0, 1\}$

Logistic regression model: $\mathbb{P}(y_i = 1 | x_i; \beta) = \frac{\exp(\beta_0 + \sum_{j=1}^d \beta_j x_{ij})}{1 + \exp(\beta_0 + \sum_{j=1}^d \beta_j x_{ij})}$

Covariables: $x_i \underset{\text{i.i.d.}}{\sim} \mathcal{N}_d(\mu, \Sigma)$

Log-likelihood with $\theta = (\mu, \Sigma, \beta)$:

$$\mathcal{LL}(\theta; x, y) = \sum_{i=1}^n \left(\log(p(y_i | x_i; \beta)) + \log(p(x_i; \mu, \Sigma)) \right).$$

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

Likelihood inference with Missing At Random values

$$\mathcal{LL}(\theta; x, y) = \sum_{i=1}^n \left(\log(p(y_i|x_i; \beta)) + \log(p(x_i; \mu, \Sigma)) \right)$$

X_1	X_2	X_3	...	M_1	M_2	M_3	...	Y
NA	20	10	...	1	0	0	...	shock
-6	45	NA	...	0	0	1	...	shock
0	NA	30	...	0	1	0	...	no shock
NA	32	35	...	1	0	0	...	shock

$m = (m_{ij})$ a $n \times d$ matrix $m_{ij} = 0$ if x_{ij} is observed and 1 otherwise

$(y_i, x_i, m_i) \underset{\text{i.i.d.}}{\sim} \{p_\theta(x, y)q_\phi(m|x, y)\}$ data & missing values mechanism

Likelihood inference with Missing At Random values

$$\mathcal{L}(\theta; x, y) = \sum_{i=1}^n \left(\log(p(y_i|x_i; \beta)) + \log(p(x_i; \mu, \Sigma)) \right)$$

X_1	X_2	X_3	...	M_1	M_2	M_3	...	Y
NA	20	10	...	1	0	0	...	shock
-6	45	NA	...	0	0	1	...	shock
0	NA	30	...	0	1	0	...	no shock
NA	32	35	...	1	0	0	...	shock

$m = (m_{ij})$ a $n \times d$ matrix $m_{ij} = 0$ if x_{ij} is observed and 1 otherwise

$(y_i, x_i, m_i) \underset{\text{i.i.d.}}{\sim} \{p_\theta(x, y)q_\phi(m|x, y)\}$ data & missing values mechanism

Ex: Income & Age with missing values on income

MAR: depends only on observed values, i.e. on age (not income)

Ignorable mechanism $\mathcal{L}_{obs}(\theta) \triangleq \prod_{i=1}^n \int p_\theta(x_i, y_i) dx_{i,mis}$

Stochastic Approximation EM - package misaem

$$\arg \max \mathcal{LL}(\theta; x_{\text{obs}}, y) = \int \mathcal{LL}(\theta; x, y) dx_{\text{mis}}$$

- **E-step:** Evaluate the quantity

$$\begin{aligned} Q_k(\theta) &= \mathbb{E}[\mathcal{LL}(\theta; x, y) | x_{\text{obs}}, y; \theta_{k-1}] \\ &= \int \mathcal{LL}(\theta; x, y) p(x_{\text{mis}} | x_{\text{obs}}, y; \theta_{k-1}) dx_{\text{mis}} \end{aligned}$$

- **M-step:** $\theta_k = \arg \max_{\theta} Q_k(\theta)$

\Rightarrow *Unfeasible computation of expectation*

MCEM (Wei & Tanner, 1990): Generate samples of missing data from $p(x_{\text{mis}} | x_{\text{obs}}, y; \theta_{k-1})$ and replace the expectation by an empirical mean

\Rightarrow *Require a huge number of samples*

SAEM (Lavielle, 2014) almost sure convergence to MLE

Unbiased estimates: $\hat{\beta}_1, \dots, \hat{\beta}_d - \hat{V}(\hat{\beta}_1), \dots, \hat{V}(\hat{\beta}_d)$ - good coverage

Supervised learning with missing values

On the consistency of supervised learning with missing values. (2019). J., Prost, Scornet & Varoquaux

- A feature matrix \mathbf{X} and a response vector Y
- Find a prediction function that minimizes the expected risk

$$\text{Bayes rule: } f^* \in \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}[\ell(f(\mathbf{X}), Y)]; \quad f^*(\mathbf{X}) = \mathbb{E}[Y|\mathbf{X}]$$

- Empirical risk: $\hat{f}_{\mathcal{D}_{n,\text{train}}} \in \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \left(\frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{X}_i), Y_i) \right)$

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

- Bayes consistent: $\mathbb{E}[\ell(\hat{f}_n(\mathbf{X}), Y)] \xrightarrow{n \rightarrow \infty} \mathbb{E}[\ell(f^*(\mathbf{X}), Y)]$

On the consistency of supervised learning with missing values. (2019). J., Prost, Scornet & Varoquaux

- A feature matrix \mathbf{X} and a response vector Y
- Find a prediction function that minimizes the expected risk

$$\text{Bayes rule: } f^* \in \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}[\ell(f(\mathbf{X}), Y)]; \quad f^*(\mathbf{X}) = \mathbb{E}[Y|\mathbf{X}]$$

- Empirical risk: $\hat{f}_{\mathcal{D}_{n,\text{train}}} \in \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \left(\frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{X}_i), Y_i) \right)$

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

- Bayes consistent: $\mathbb{E}[\ell(\hat{f}_n(\mathbf{X}), Y)] \xrightarrow{n \rightarrow \infty} \mathbb{E}[\ell(f^*(\mathbf{X}), Y)]$

Differences with classical literature

- explicitly consider the response variable Y - Aim: Prediction
 - two data sets (out of sample) with missing values: Train & test sets
- ⇒ Is it possible to use previous approaches (EM - impute), consistent?
- ⇒ Do we need to design new ones?

EM and out-of sample prediction - package misaem

$$\mathbb{P}(y_i = 1|x_i; \beta) = \frac{\exp(x_i\beta)}{1+\exp(x_i\beta)} \quad \text{After EM: } \hat{\theta}_n = (\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_d, \hat{\mu}, \hat{\Sigma})$$

New obs: $x_{n+1} = (x_{(n+1)1}, \text{NA}, \text{NA}, x_{(n+1)4}, \dots, x_{(n+1)d})$

Predict Y on a **test set with missing entries** $x_{\text{test}} = (x_{\text{obs}}, x_{\text{miss}})$

EM and out-of sample prediction - package misaem

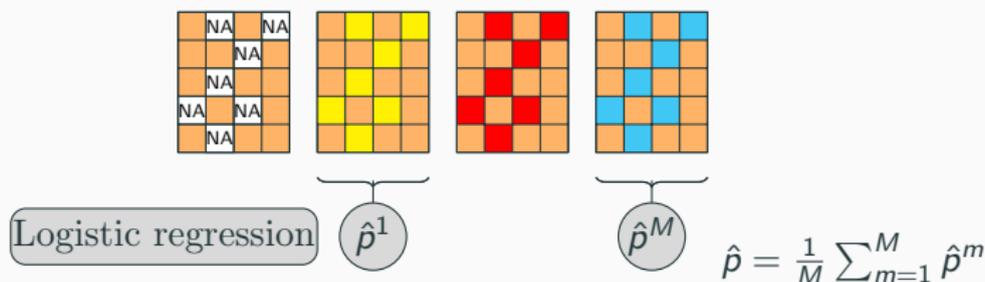
$$\mathbb{P}(y_i = 1|x_i; \beta) = \frac{\exp(x_i\beta)}{1+\exp(x_i\beta)} \quad \text{After EM: } \hat{\theta}_n = (\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_d, \hat{\mu}, \hat{\Sigma})$$

New obs: $x_{n+1} = (x_{(n+1)1}, \text{NA}, \text{NA}, x_{(n+1)4}, \dots, x_{(n+1)d})$

Predict Y on a test set with missing entries $x_{\text{test}} = (x_{\text{obs}}, x_{\text{miss}})$

$$\begin{aligned} \hat{y} &= \arg \max_y p_{\hat{\theta}}(y|x_{\text{obs}}) = \arg \max_y \int p_{\hat{\theta}}(y|x) p_{\hat{\theta}}(x_{\text{miss}}|x_{\text{obs}}) dx_{\text{miss}} \\ &= \arg \max_y \mathbb{E}_{p_{x_m|x_o=x_o}} p_{\hat{\theta}_n}(y|X_m, x_o) \approx \arg \max_y \sum_{m=1}^M p_{\hat{\theta}_n}(y|x_{\text{obs}}, x_{\text{miss}}^{(m)}) \end{aligned}$$

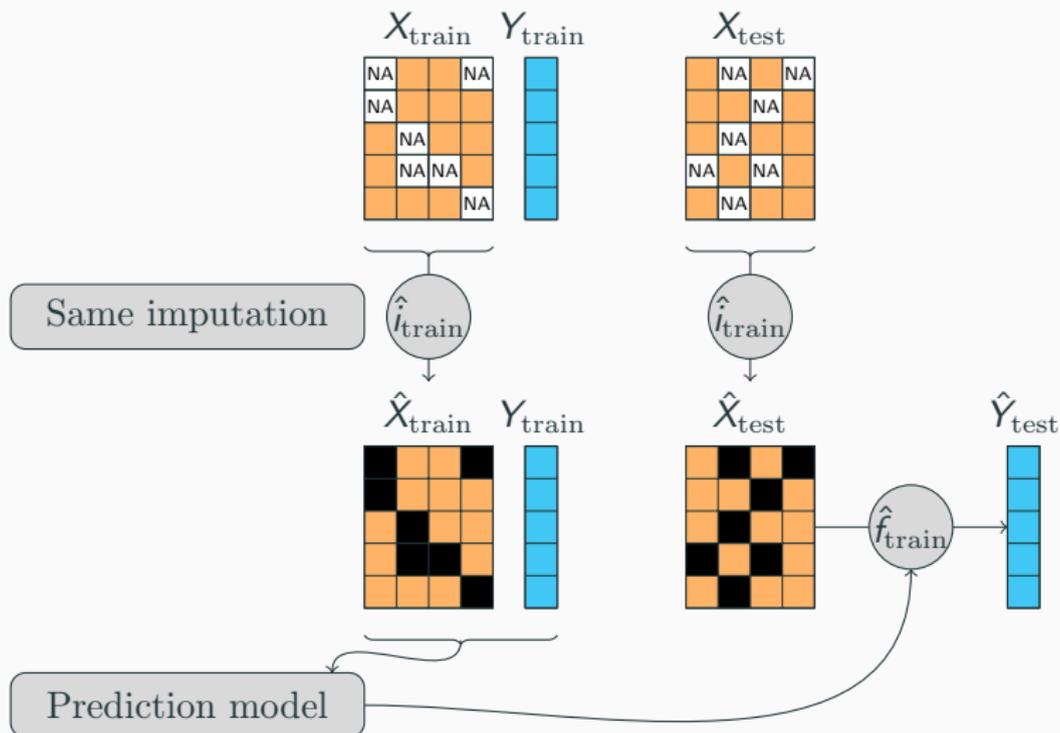
\approx Multiple imputation: Draw M values from $X_{\text{miss}}|X_{\text{obs}}$



Imputation prior to learning

Impute the train with \hat{I}_{train} learn a model \hat{f}_{train} with $\hat{X}_{train}, Y_{train}$

Impute the test with the same imputation \hat{I}_{train} - predict \hat{X}_{test} with \hat{f}_{train}



Imputation prior to learning

Imputation with the same model

Easy to implement for univariate imputation: The means ($\hat{\mu}_1, \dots, \hat{\mu}_d$) of each column of the train. Also OK for Gaussian imputation.

Issue: Many methods are "black-boxes" and take as an input the incomplete data and output the completed data (`mice`, `missForest`)

Separate imputation

Impute train and test separately (with a different model)

Issue: Depends on the size of the test set? one observation?

Group imputation/ semi-supervised

Impute train and test simultaneously but the predictive model is learned only on the training imputed data set

Issue: Sometimes no training set at test time

Imputation with the same model: Mean imputation consistent

Learn on the mean-imputed training data, impute the test set with the **same means** and predict is optimal if the missing data are MAR and the **learning algorithm is universally consistent**

Framework - assumptions

- $Y = f(\mathbf{X}) + \varepsilon$
- $\mathbf{X} = (X_1, \dots, X_d)$ has a continuous density $g > 0$ on $[0, 1]^d$
- $\|f\|_\infty < \infty$
- Missing data MAR on X_1 with $M_1 \perp\!\!\!\perp X_1 | X_2, \dots, X_d$.
- $(x_2, \dots, x_d) \mapsto \mathbb{P}[M_1 = 1 | X_2 = x_2, \dots, X_d = x_d]$ is continuous
- ε is a centered noise independent of (\mathbf{X}, M_1)

(remains valid when missing values occur for variables X_1, \dots, X_j)

Imputation with the same model: Mean imputation consistent

Learn on the mean-imputed training data, impute the test set with the **same means** and predict is optimal if the missing data are MAR and the **learning algorithm is universally consistent**

Mean imputed entry $\mathbf{x}' = (x'_1, x_2, \dots, x_d)$: $x'_1 = x_1 \mathbb{1}_{M_1=0} + \mathbb{E}[X_1] \mathbb{1}_{M_1=1}$

Note the data: $\tilde{\mathbf{X}} = \mathbf{X} \odot (\mathbf{1} - \mathbf{M}) + \text{NA} \odot \mathbf{M}$ (takes value in $\mathbb{R} \cup \{\text{NA}\}$)

Theorem

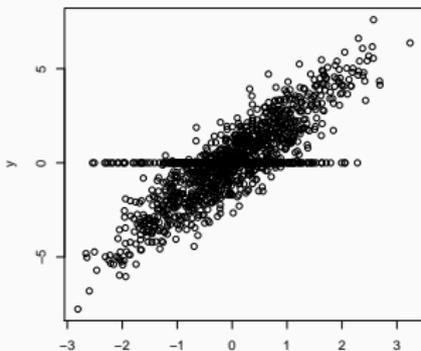
Prediction with mean is equal to the Bayes function almost everywhere

$$f_{\text{impute}}^*(x') = \tilde{f}^*(\tilde{\mathbf{X}}) = \mathbb{E}[Y | \tilde{\mathbf{X}} = \tilde{\mathbf{x}}]$$

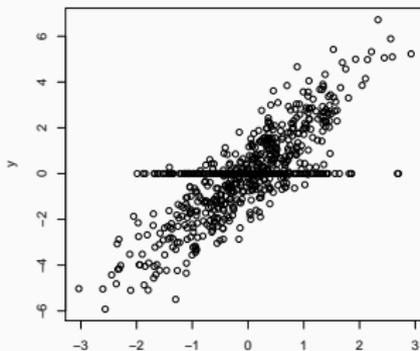
Other values than the mean are OK but use the same value for the train and test sets, otherwise the algorithm may fail as the distributions differ

Consistency of supervised learning with NA: Rationale

- Specific value, systematic like a code for missing
- The learner detects the code and recognizes it at the test time
- With categorical data, just code "Missing"
- With continuous data, any constant:
- Need a lot of data (asymptotic result) and a super powerful learner



Train

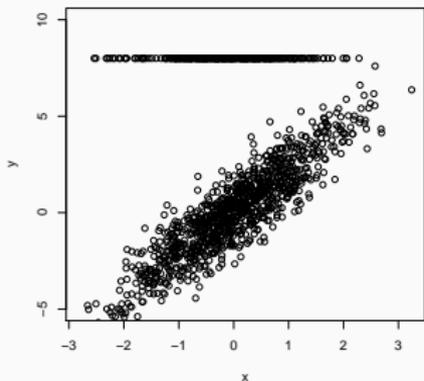


Test

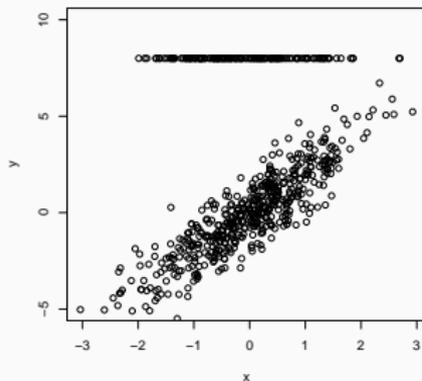
Mean imputation not bad for prediction; it is consistent; despite its drawbacks for estimation - Useful in practice!

Consistency of supervised learning with NA: Rationale

- Specific value, systematic like a code for missing
- The learner detects the code and recognizes it at the test time
- With categorical data, just code "Missing"
- With continuous data, any constant: out of range
- Need a lot of data (asymptotic result) and a super powerful learner



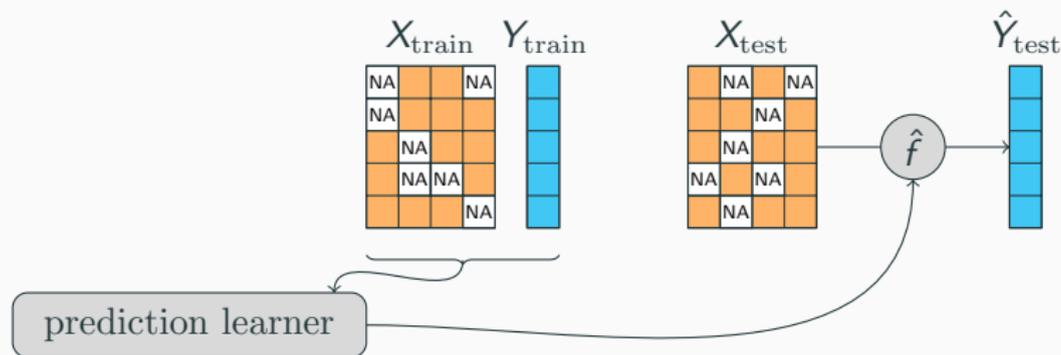
Train



Test

Mean imputation not bad for prediction; it is consistent; despite its drawbacks for estimation - Useful in practice!

End-to-end learning with missing values

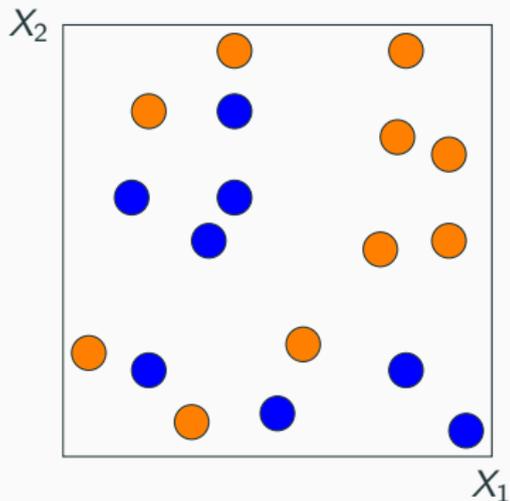


- Trees well suited for empirical risk minimization with missing values: Handle half discrete data $\tilde{\mathbf{X}}$ that takes values in $\mathbb{R} \cup \{\text{NA}\}$
- Random forests powerful learner

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

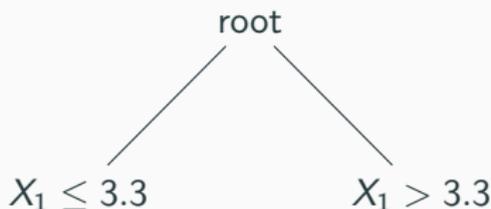
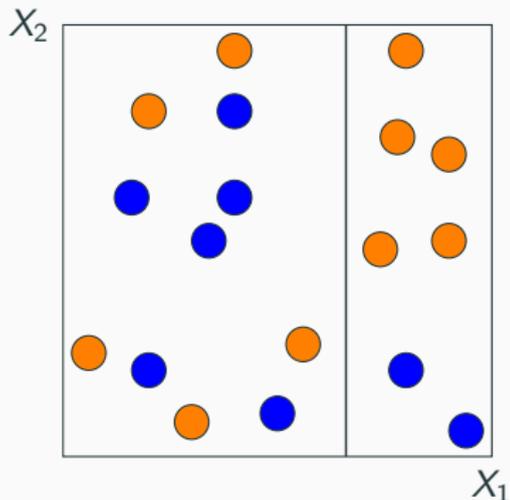
$$(j^*, z^*) \in \arg \min_{(j,z) \in \mathcal{S}} \mathbb{E} \left[(Y - \mathbb{E}[Y|X_j \leq z])^2 \cdot \mathbb{1}_{X_j \leq z} + (Y - \mathbb{E}[Y|X_j > z])^2 \cdot \mathbb{1}_{X_j > z} \right].$$



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

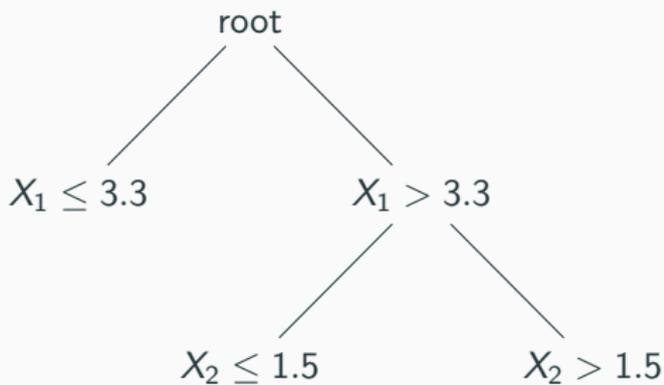
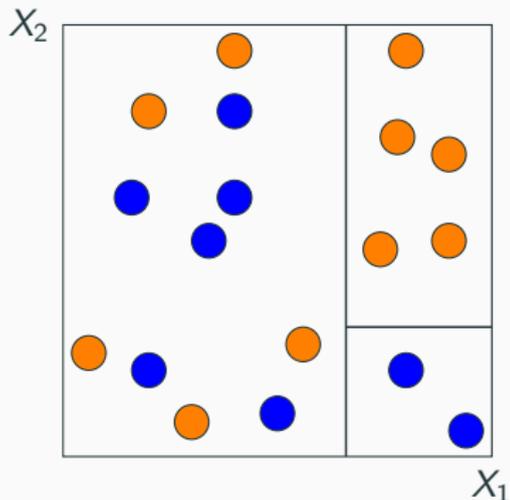
$$(j^*, z^*) \in \arg \min_{(j,z) \in \mathcal{S}} \mathbb{E} \left[(Y - \mathbb{E}[Y|X_j \leq z])^2 \cdot \mathbb{1}_{X_j \leq z} + (Y - \mathbb{E}[Y|X_j > z])^2 \cdot \mathbb{1}_{X_j > z} \right].$$



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

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



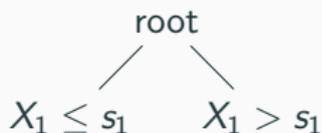
CART with missing values

root

	X_1	X_2	Y
1			
2	NA		
3	NA		
4			

CART with missing values

	X_1	X_2	Y
1			
2	NA		
3	NA		
4			



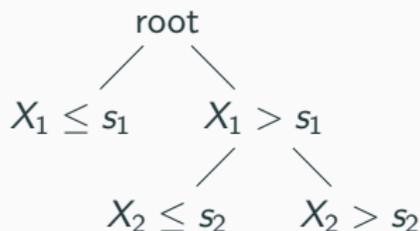
1) Select variable and threshold on observed data ⁴

$$\mathbb{E} \left[(Y - \mathbb{E}[Y|X_j \leq z, M_j = 0])^2 \cdot \mathbb{1}_{X_j \leq z, M_j = 0} + (Y - \mathbb{E}[Y|X_j > z, M_j = 0])^2 \cdot \mathbb{1}_{X_j > z, M_j = 0} \right].$$

⁴ Variable selection bias (not a problem to predict): `ctree` function, `partykit` package, Hothorn, Hornik & Zeileis.

CART with missing values

	X_1	X_2	Y
1			
2	NA		
3	NA		
4			



1) Select variable and threshold on observed data ⁴

$$\mathbb{E} \left[(Y - \mathbb{E}[Y|X_j \leq z, M_j = 0])^2 \cdot \mathbb{1}_{X_j \leq z, M_j = 0} + (Y - \mathbb{E}[Y|X_j > z, M_j = 0])^2 \cdot \mathbb{1}_{X_j > z, M_j = 0} \right].$$

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

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

⁴ Variable selection bias (not a problem to predict): `ctree` function, `partykit` package, Hothorn, Hornik & Zeileis.

Missing incorporated in attribute (Twala et al. 2008)

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

$$f^* \in \arg \min_{f \in \mathcal{P}_{c,miss}} \mathbb{E} \left[(Y - f(\tilde{\mathbf{X}}))^2 \right],$$

where $\mathcal{P}_{c,miss} = \mathcal{P}_{c,miss,L} \cup \mathcal{P}_{c,miss,R} \cup \mathcal{P}_{c,miss,sep}$ with

1. $\mathcal{P}_{c,miss,L} \rightarrow \{ \{ \tilde{X}_j \leq z \vee \tilde{X}_j = \text{NA} \}, \{ \tilde{X}_j > z \} \}$
2. $\mathcal{P}_{c,miss,R} \rightarrow \{ \{ \tilde{X}_j \leq z \}, \{ \tilde{X}_j > z \vee \tilde{X}_j = \text{NA} \} \}$
3. $\mathcal{P}_{c,miss,sep} \rightarrow \{ \{ \tilde{X}_j \neq \text{NA} \}, \{ \tilde{X}_j = \text{NA} \} \}$.

- Missing values treated like a category (well to handle $\mathbb{R} \cup \text{NA}$)
- Good for informative pattern (\mathbf{M} explains Y)
- Implementation: Duplicate the incomplete columns, and replace the missing entries once by $+\infty$ and once by $-\infty$ (J. Tibshirani)

Implemented for conditional trees and forests [partykit package](#)

\Rightarrow Target one model per pattern (2^d):

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

Simulations: 20% missing values

Quadratic: $Y = X_1^2 + \varepsilon$, $x_i. \in \mathcal{N}(\mu, \Sigma_{4 \times 4})$, $\rho = 0.5$, $n = 1000$

$$\tilde{d}_n = \begin{bmatrix} 2 & 3 & \text{NA} & 0 & 15 \\ 1 & \text{NA} & 3 & 5 & 13 \\ 9 & 4 & 2 & \text{NA} & 18 \\ 7 & 6 & \text{NA} & \text{NA} & 10 \end{bmatrix}$$

$$\tilde{d}_n + \text{mask} = \begin{bmatrix} 2 & 3 & \text{NA} & 0 & 0 & 0 & 1 & 0 & 15 \\ 1 & \text{NA} & 3 & 5 & 0 & 1 & 0 & 0 & 13 \\ 9 & 4 & 2 & \text{NA} & 0 & 0 & 0 & 1 & 18 \\ 7 & 6 & \text{NA} & \text{NA} & 0 & 0 & 1 & 1 & 10 \end{bmatrix}$$

Imputation (mean, Gaussian) + prediction with trees

Imputation (mean, Gaussian) + mask + prediction with trees

Trees MIA

Simulations: 20% missing values

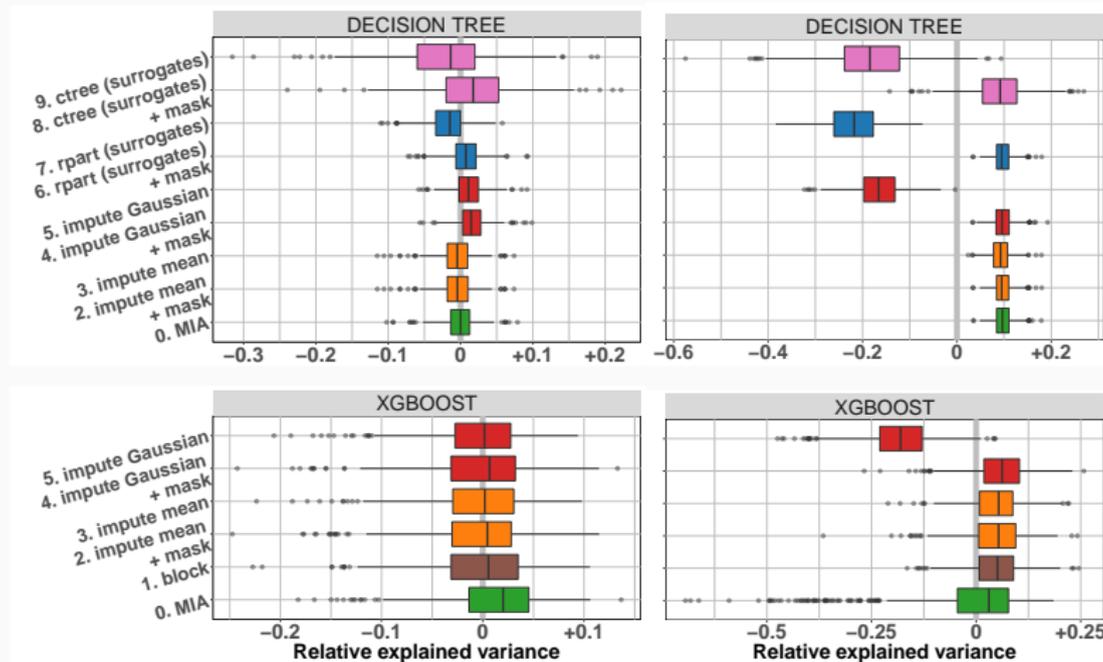
Quadratic: $Y = X_1^2 + \varepsilon$, $x_i \in \mathcal{N}(\mu, \Sigma_{4 \times 4})$, $\rho = 0.5$, $n = 1000$

MCAR (MAR)

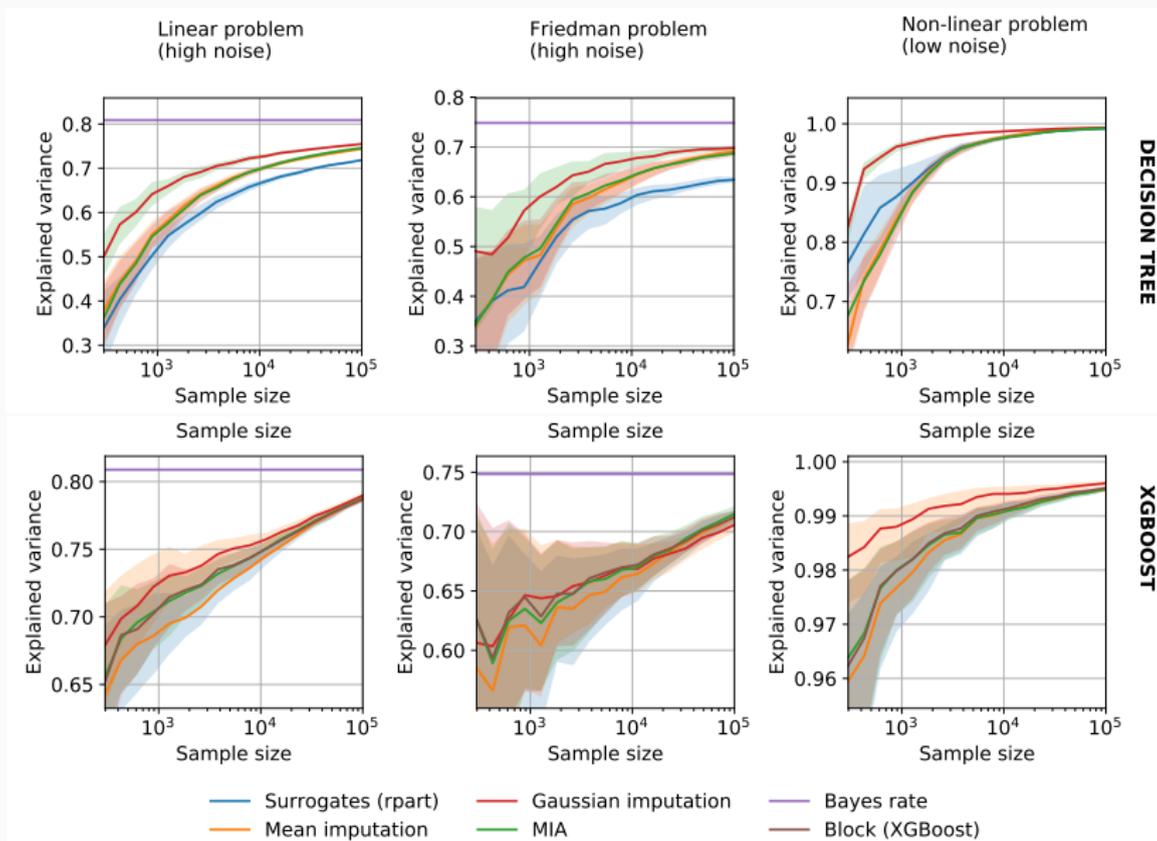
$$M_{i,1} \sim \mathcal{B}(\rho)$$

MNAR - Predictive

$$M_{i,1} = \mathbb{1}_{X_{i,1} > [X_1]_{(1-\rho)n}} - Y = X_1^2 + 3M_1 + \varepsilon$$



Consistency: 40% missing values MCAR



Discussion - challenges

Take home message EM/imputation

- Few implementation of EM strategies

“The idea of imputation is both seductive and dangerous”. *It is seductive because it can lull the user into the pleasurable state of believing that the data are complete after all, and it is dangerous because it lumps together situations where the problem is sufficiently minor that it can be legitimately handled in this way and situations where standard estimators applied to the imputed data have substantial biases.”* (Dempster & Rubin, 1983)

- Single imputation aims at completing a dataset as best as possible
- **Multiple imputation** aims at estimating the parameters and their variability taking into account the uncertainty of the missing values
- Single imputation can be appropriate for point estimates
- Both % of NA & structure matter (5% of NA can be an issue)

Principal component methods powerful for single & multiple imputation of quanti & categorical data: Dimensionality reduction and capture similarities between observations and variables. `missMDA` package

Take-home message supervised learning

- Incomplete train and test → **same imputation model**
- **Single mean imputation is consistent given a powerful learner**
- Empirically, good imputation methods reduce the number of samples required to reach good prediction

Tree-based models :

- **Missing Incorporated in Attribute** optimizes not only the split but also the handling of the missing values
- Informative missing data: **Adding the mask** helps imputation - MIA

To be done

- Nonasymptotic results
- Uncertainty associated with the prediction
- Distributional shift: No missing values in the test set?
- Prove the usefulness of methods in MNAR

Still an active area of research! Join this exciting field!

Current works

- Variable selection in high dimension Adaptive bayesian SLOPE with missing values. 2019. Jiang, Bogdan, J., Miasojedow, Rockova & TraumaBase
- **MNAR missing values**
 - Contribution of causality for missing data
 - Graphical Models for Processing Missing Data. 2019. Mohan, Pearl.
 - Estimation and imputation in Probabilistic Principal Component Analysis with Missing Not At Random data. 2019. Sportisse, Boyer, J.
 - Contribution of neural nets J., Prost, Scornet, Varoquaux

Other challenges

- MI theory: Good theory for regression parameters but others? Theory with other asymptotic small n , large p ?, etc.
- Practical imputation issues: Imputation not in agreement (X & X^2), imputation out of range? problems of logical bounds (> 0), etc.

[R-miss-tastic](https://rmissstastic.netlify.com/R-miss-tastic) <https://rmissstastic.netlify.com/R-miss-tastic>

J., I. Mayer, N. Tierney & N. Vialaneix

Project funded by the R consortium (Infrastructure Steering Committee)⁵

Aim: a reference platform on the theme of missing data management

- list existing packages
- available literature
- tutorials
- analysis workflows on data
- main actors

⇒ Federate the community

⇒ Contribute!

⁵<https://www.r-consortium.org/projects/call-for-proposals>

Examples:

- Lecture ⁶ - General tutorial : Statistical Methods for Analysis with Missing Data (Mauricio Sadinle)
- Lecture - Multiple Imputation: mice by Nicole Erler ⁷
- Longitudinal data, Time Series Imputation (Steffen Moritz - very active contributor of r-miss-tastic), Principal Component Methods⁸

⁶<https://rmissstastic.netlify.com/lectures/>

⁷https://rmissstastic.netlify.com/tutorials/erler_course_multipleimputation_2018/erler_practical_mice_2018

⁸https://rmissstastic.netlify.com/tutorials/Josse_slides_imputation_PCA_2018.pdf

Thank you

