

Learning with missing values: from estimation to prediction

Erwan Scornet
Lecturer at Sorbonne University

October 2025



Collaborators

- Alexis Ayme, Post-doc ENS Ulm, Paris. Linear models, Optim.
- Claire Boyer, Professor Paris-Saclay. Signal, Optim.
- Aymeric Dieuleveut, Professor at IPP, Paris. Optim.
- Julie Josse, Senior researcher, INRIA, Montpellier. Causality
- Marine Le Morvan, Junior researcher, INRIA, Paris. Supervised learning
- Christophe Muller, PhD student, Oxford
- Jeffrey Naf, Assist. Professor, GSEM, Geneva. Distributional Prediction
- Angel Reyero Lobo, PhD Student, INRIA, Toulouse. Variable Importance
- Gael Varoquaux, Senior researcher, INRIA, Paris. ML, Scikit-learn



Thanks for the slides too!

Traumabase: an observational French registry²

3 / 99

- ▷ 40000 trauma patients
- ▷ 300 heterogeneous features from pre-hospital and in-hospital settings
- ▷ 40 trauma centers, 4000 new patients per year

Center	Accident	Age	Sex	Lactate	Blood Pres.	Shock	Platelet	...
Beaujon	fall	54	m	NM	180	yes	292000	
Pitie	gun	26	m	NA	131	no	323000	
Beaujon	moto	63	m	3.9	NR	yes	318000	
Pitie	moto	30	w	Imp	107	no	211000	
⋮								

⇒ **Explain and Predict** hemorrhagic shock, need for neurosurgery and need for a trauma center given pre-hospital features.

Ex: logistic regression/ random forests + **Quantify uncertainty**¹

¹Zaffran, J., Dieuleveut, Romano. Conformal Prediction with Missing Values. *ICML 2023*.

²www.traumabase.eu - <https://www.traumatrix.fr/>

Missing values^{3, 4, 5}

4 / 99

Missing values are everywhere: unanswered questions in a survey, lost data, damaged plants, machines that fail...



"The best thing to do with missing values is not to have any"

Gertrude Mary Cox (1900-1978)

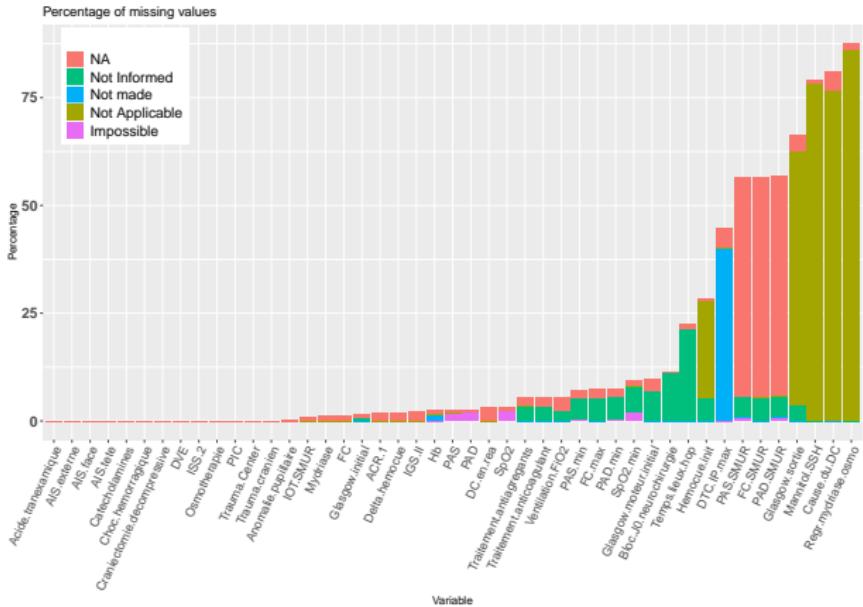
⇒ Still an issue in the "big data" area (data from different sources)

³Little & Rubin (2019). Statistical Analysis with Missing Data, Third Edition, Wiley.

⁴Van Buuren (2018). Flexible Imputation of Data. Second Edition, Chapman & Hall.

⁵Schafer (1997). Analysis of Incomplete Multivariate Data, Chapman & Hall.

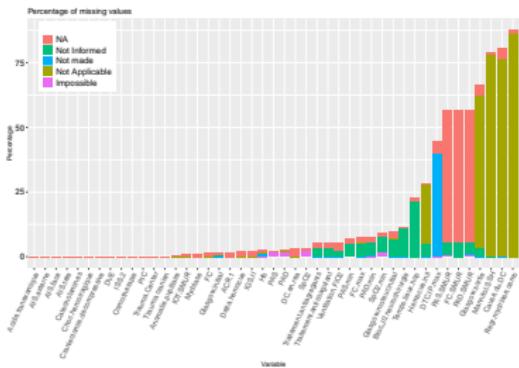
Missing data: important bottleneck in statistical practice



Different types of missing values

- ▷ Not informed: not recorded
- ▷ Not made: possibly due to patient status
- ▷ Not applicable: not supposed to be measured
- ▷ Impossible
- ▷ NA: unknown

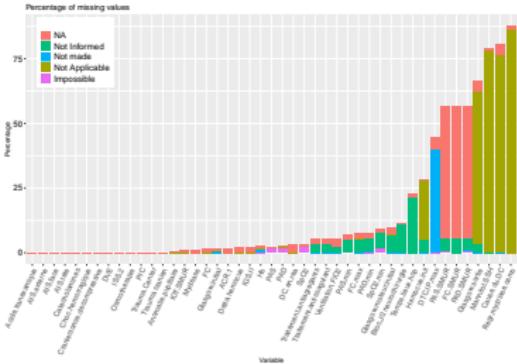
Missing data: important bottleneck in statistical practice



"One of the ironies of Big Data is that missing data play an ever more significant role"⁶

⁶Zhu, Wang, Samworth. High-dimensional PCA with heterogeneous missingness. *JRSSB*. 2022.

Missing data: important bottleneck in statistical practice 6 / 99



"One of the ironies of Big Data is that missing data play an ever more significant role"⁶

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

$d = 5 \quad \Rightarrow \quad \approx 95\% \text{ of rows kept}$
 $d = 300 \quad \Rightarrow \quad \approx 5\% \text{ of rows kept}$

⁶Zhu, Wang, Samworth. High-dimensional PCA with heterogeneous missingness. *JRSSB*. 2022.

Linear model

$$Y = X^T \beta^* + \text{noise}$$

- ▷ $Y \in \mathbb{R}$ (regression) outcome is always observed
- ▷ $X \in \mathbb{R}^d$ contains missing values!

Linear model

$$Y = X^T \beta^* + \text{noise}$$

- ▷ $Y \in \mathbb{R}$ (regression) outcome is always observed
- ▷ $X \in \mathbb{R}^d$ contains missing values!

Three different tasks: imputation, estimation, prediction.

1. **Imputation** - Replace missing values to obtain a complete data set, on which any classical analysis can be performed.
2. **Estimation** - Provide an estimate of β^* - allows predicting outputs of complete data.

Linear model

$$Y = X^T \beta^* + \text{noise}$$

- ▷ $Y \in \mathbb{R}$ (regression) outcome is always observed
- ▷ $X \in \mathbb{R}^d$ contains missing values!

Three different tasks: imputation, estimation, prediction.

1. **Imputation** - Replace missing values to obtain a complete data set, on which any classical analysis can be performed.
2. **Estimation** - Provide an estimate of β^* - allows predicting outputs of complete data.
3. **Prediction** - Predict Y for a new X with missing entries

Warning: A good estimate of β^* does not lead to a prediction of Y

$$X = (\text{na}, 5, \text{na}, -6) \quad X^T \beta^* = ??$$

Solutions to handle missing values in the covariates

8 / 99

Abundant literature: Creation of **Rmistaic platform**⁷ (**> 150 packages**)

- ▷ **Imputation:** (Single/Multiple) imputation to get a/several complete data set(s). Ex: (M)ICE
- ▷ **Estimation:** Modify the estimation process to deal with missing values
 - Maximum likelihood inference: Expectation Maximization algorithms⁸
- ▷ **Prediction:** Predict an outcome with missing data in covariates⁹¹⁰.
Solutions: using deterministic (e.g. constant) imputation or Missing Incorporated in Attributes for trees based methods (**grf package**)

⁷Mayer, J. et al. A unified platform for missing values methods and workflows. *R journal*. 2022.

⁸Jiang, J. et al. Logistic Regression with Missing Covariates *CSDA*. 2019. - **misaem package**

⁹J. et al. Consistency of supervised learning with missing values. *Stats papers*. 2018-2024.

¹⁰Le morvan, J. et al. What's a good imputation to predict with missing values? *Neurips2021*.

1. Missing values mechanism
2. Single Imputation
3. Multiple Imputation
4. Imputation quality
5. Supervised Learning with Missing values
 - Decision trees as PbP predictors
 - Impute-then-regress procedures with consistent predictors
6. Linear models
 - Linear regression: A pattern-by-pattern approach
 - Linear regression: Impute-then-regress procedures via zero-imputation
 - Classification with missing values
7. Conclusion

1. Missing values mechanism

2. Single Imputation

3. Multiple Imputation

4. Imputation quality

5. Supervised Learning with Missing values

Decision trees as PbP predictors

Impute-then-regress procedures with consistent predictors

6. Linear models

Linear regression: A pattern-by-pattern approach

Linear regression: Impute-then-regress procedures via zero-imputation

Classification with missing values

7. Conclusion

Missing values mechanism: Rubin's taxonomy^{11, 12}

11 / 99

- Random Variables:

- ▷ $X^* \in \mathbb{R}^d$: complete unavailable data, $X \in \mathbb{R}^d$: observed data with NA
- ▷ $M \in \{0, 1\}^d$: missing pattern, or mask, $M_j = 1$ if and only if X_j is missing

- Realizations: For a pattern m , $o(x, m) = (x_j)_{j \in \{1, \dots, d\}: m_j=0}$ the observed elements of x and while $o^c(x, m) = (x_j)_{j \in \{1, \dots, d\}: m_j=1}$, the missing elements.

$$x^* = (1, 2, 3, 8, 5)$$

$$x = (1, \text{NA}, 3, 8, \text{NA})$$

$$m = (0, 1, 0, 0, 1)$$

$$o(x, m) = (1, 3, 8), \quad o^c(x^*, m) = (2, 5)$$

¹¹Rubin. Inference and missing data. *Biometrika*. 1976.

¹²What Is Meant by "Missing at Random"? Seaman, et al. *Statistical Science*. 2013.

Missing values mechanism: Rubin's taxonomy^{11, 12}

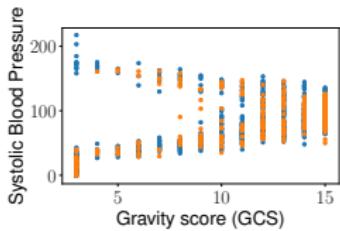
11 / 99

- Random Variables:

- ▷ $X^* \in \mathbb{R}^d$: complete unavailable data, $X \in \mathbb{R}^d$: observed data with NA
- ▷ $M \in \{0, 1\}^d$: missing pattern, or mask, $M_j = 1$ if and only if X_j is missing

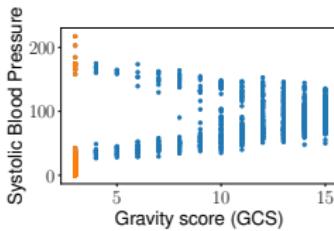
For a pattern m , $o(x, m) = (x_j)_{j \in \{1, \dots, d\}: m_j=0}$ the observed elements of x and while $o^c(x, m) = (x_j)_{j \in \{1, \dots, d\}: m_j=1}$, the missing elements.

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



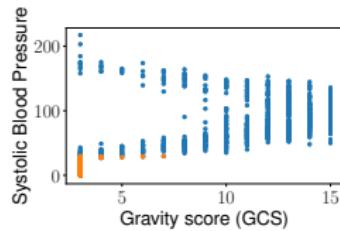
**Missing Completely at Random
(MCAR)**

$$m \in \mathcal{M}, x \in \mathcal{X},$$
$$\mathbb{P}(M = m|x) = \mathbb{P}(M = m)$$



**Missing at Random
(MAR)**

$$\forall m \in \mathcal{M}, x \in \mathcal{X}$$
$$\mathbb{P}(M = m|x)$$
$$= \mathbb{P}(M = m|o(x, m))$$



**Missing Not At Random
(MNAR)**

If not MAR: it is MNAR

¹¹Rubin. Inference and missing data. *Biometrika*. 1976.

¹²What Is Meant by "Missing at Random"? Seaman, et al. *Statistical Science*. 2013.

Two views to model the joint distribution of (X, M)

12 / 99

- ▷ Selection Model¹³: $p^*(M = m, x) = \mathbb{P}(M = m | x)p^*(x)$

Definition: SM-MAR

$$\mathbb{P}(M = m | x) = \mathbb{P}(M = m | o(x, m)) \text{ for all } m \in \mathcal{M}, x \in \mathcal{X}.$$

The proba. of any m occurring only depends on the obs part of x .

- ▷ Pattern Mixture Model¹⁴: $p^*(M = m, x) = p^*(x | M = m)\mathbb{P}(M = m)$

Definition: PMM-MAR

$$p^*(o^c(x, m) | o(x, m), M = m) = \textcolor{red}{p^*(o^c(x, m) | o(x, m))}.$$

for all $m \in \mathcal{M}, x \in \mathcal{X}$. The conditional distrib. of missing given obs. in pattern m is equal to the unconditional one.^a

^aMolenberghs et al. Every MNAR model has a MAR counterpart with equal fit. *JRSSB*. 2008

- Proposition: SM-MAR is equivalent to PMM-MAR

¹³Heckman. Sample selection bias as a specification error. *Econometrica*. 1979

¹⁴Little. Pattern-mixture models for multivariate incomplete data. *JASA*. 1993

Testing the missing values mechanism

13 / 99

- ▷ Can we observe the missing value mechanism from the sample?

Unfortunately, the general answer is **no**

¹⁵ Little. *A Test of Missing Completely at Random for Multivariate Data with Missing Values.* 1988

¹⁶ Michel, Naf, Spohn, Meinshausen. PKLM: a flexible MCAR test using classification, *Psychometrika*. 2025

¹⁷ Berrett, Samworth. *Optimal nonparametric testing of missing completely at random and its connections to compatibility*, AoS. 2023

- ▷ Can we observe the missing value mechanism from the sample?

Unfortunately, the general answer is **no**

MCAR vs MAR in Gaussian setting

- ▷ If we assume MAR is true we can test H_0 : MCAR vs H_A : MAR.
- ▷ A classical test is the Little test¹⁵ that operates under the assumption of Gaussianity.

¹⁵ Little. *A Test of Missing Completely at Random for Multivariate Data with Missing Values*. 1988

¹⁶ Michel, Naf, Spohn, Meinshausen. PKLM: a flexible MCAR test using classification, *Psychometrika*. 2025

¹⁷ Berrett, Samworth. *Optimal nonparametric testing of missing completely at random and its connections to compatibility*, *AoS*. 2023

- ▷ Can we observe the missing value mechanism from the sample?

Unfortunately, the general answer is **no**

MCAR vs MAR in Gaussian setting

- ▷ If we assume MAR is true we can test H_0 : MCAR vs H_A : MAR.
- ▷ A classical test is the Little test¹⁵ that operates under the assumption of Gaussianity.

Nonparametric tests

- ▷ One of the very few (if not only) useable nonparametric test is our PKLM Test¹⁶
- ▷ There is also interesting theoretical work¹⁷

¹⁵ Little. *A Test of Missing Completely at Random for Multivariate Data with Missing Values*. 1988

¹⁶ Michel, Naf, Spohn, Meinshausen. PKLM: a flexible MCAR test using classification, *Psychometrika*. 2025

¹⁷ Berrett, Samworth. *Optimal nonparametric testing of missing completely at random and its connections to compatibility*, AoS. 2023

1. Missing values mechanism

2. Single Imputation

3. Multiple Imputation

4. Imputation quality

5. Supervised Learning with Missing values

Decision trees as PbP predictors

Impute-then-regress procedures with consistent predictors

6. Linear models

Linear regression: A pattern-by-pattern approach

Linear regression: Impute-then-regress procedures via zero-imputation

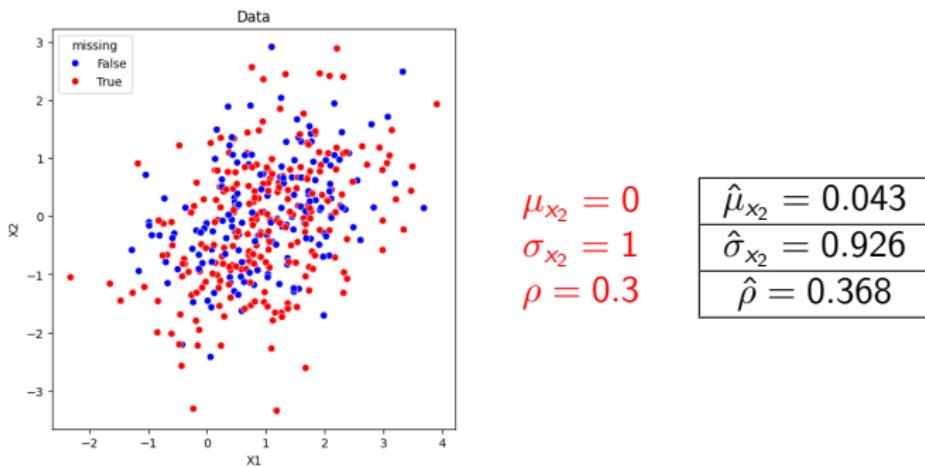
Classification with missing values

7. Conclusion

Generative setting

- ▷ $(X_1, X_2) \sim \mathcal{N}((\mu_{x_1}, \mu_{x_2}), \Sigma); n = 400$
- ▷ $(\mu_{x_1}, \mu_{x_2}) = (1, 0)$ and $\Sigma = ((1, 0.3), (0.3, 1))$
- ▷ MCAR missing values on X_2 only with probability $p = 0.6$.

Discard incomplete observations and then estimate parameters

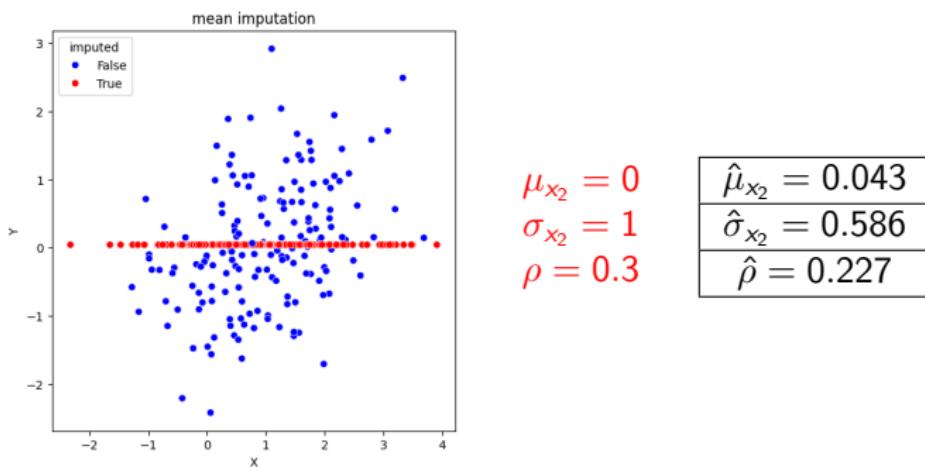


¹⁸The code to reproduce the plots is available in [Rmיסטатик](#)

Generative setting

- ▷ $(X_1, X_2) \sim \mathcal{N}((\mu_{x_1}, \mu_{x_2}), \Sigma); n = 400$
- ▷ $(\mu_{x_1}, \mu_{x_2}) = (1, 0)$ and $\Sigma = ((1, 0.3), (0.3, 1))$
- ▷ MCAR missing values on X_2 only with probability $p = 0.6$.

Impute by the mean and then estimate parameters



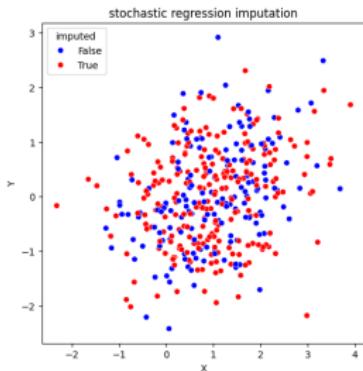
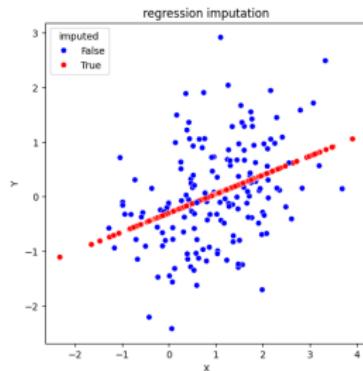
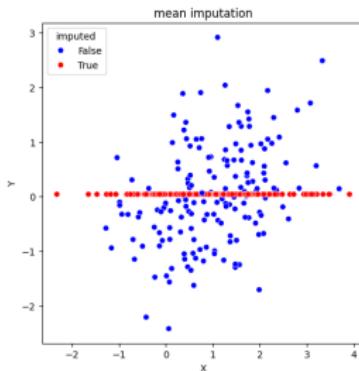
Mean imputation deforms joint and marginal distributions

Objective: to impute while preserving distribution

16 / 99

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 β (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 β and σ - impute from the predictive $\hat{x}_{i2} \sim \mathcal{N}(\hat{\beta}_0 + \hat{\beta}_1 x_{i1}, \hat{\sigma}^2) \Rightarrow$ preserve distributions



$$\begin{aligned}\mu_{x_2} &= 0 \\ \sigma_{x_2} &= 1 \\ \rho &= 0.3\end{aligned}$$

0.043
0.926
0.368

0.038
0.647
0.539

0.037
0.909
0.275

Impute while preserving distribution. Multivariate case^{17 / 99}

- ▷ Assuming a joint distribution
 - ◊ Gaussian model $x_i \sim \mathcal{N}(\mu, \Sigma)$
 - ◊ Low rank : $X_{n \times d} = \mu_{n \times d} + \varepsilon$ $\varepsilon_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$ with μ of low rank
 - ⇒ Different regularization depending on noise regime¹⁸
 - ⇒ Count data¹⁹, ordinal data, categorical data, blocks/multilevel data
 - ◊ Optimal transport²⁰, deep generative models: GAIN²¹, MIWAE²², etc.²³
²⁴
- ▷ Iterating conditional models (joint distribution implicitly defined)
 - ◊ with parametric regression (M)ICE: (Multiple) Imput. by Chained Equations²⁵
 - ◊ iterative imputation of each variable by random forests²⁶

¹⁸J. & Wager. Stable autoencoding for regularized low-rank matrix estimation. *JMLR*. 2016.

¹⁹Robin, Klopp, J., Moulines, Tibshirani. Main effects & interac. in mixed data. *JASA*. 2019.

²⁰Muzelec, Cuturi, Boyer, J. Missing Data Imputation using Optimal Transport. *ICML*. 2020.

²¹Yoon et al. GAIN: Missing data imputation using generative adversarial nets. *ICML*. 2018.

²²Mattei & Frellsen. Miwae: Deep generative model & imput. of inc. data. *ICML*. 2018.

²³Deng et al. Extended missing data imput. via gans. *DMKD*. 2022.

²⁴Fang, Bao. Fragmgan: gan for fragmentary data imputation. *STRF* 2023.

²⁵van Buuren, S. Flexible Imputation of Missing Data. Chapman & Hall/CRC Press. 2018.

²⁶Stekhoven, Bühlmann. MissForest—non-parametric imputation for mixed data. *Bioinfo*. 2012.

Imputation by Chained Equations (ICE)

18 / 99

Init.

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F

Imputation by Chained Equations (ICE)

18 / 99

Impute via mean/mode

Init.

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F

Age	Inc.	Gen.
34	13	F
18	12	F
26	14	M
26	13	F

Imputation by Chained Equations (ICE)

18 / 99

Impute via mean/mode

Init.

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F

Age	Inc.	Gen.
34	13	F
18	12	F
26	14	M
26	13	F

Set values of Age
originally missing
as unknown

1st step
Age

Age	Inc.	Gen.
34	13	F
18	12	F
?	14	M
?	13	F

Imputation by Chained Equations (ICE)

18 / 99

Impute via mean/mode

Init.

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F

Age	Inc.	Gen.
34	13	F
18	12	F
26	14	M
26	13	F

Set values of Age
originally missing
as unknown

Fit a predictive model on
complete observation
to predict Age

1st step
Age

Age	Inc.	Gen.
34	13	F
18	12	F
?	14	M
?	13	F

Age	Inc.	Gen.
34	13	F
18	12	F
?	14	M
?	13	F

Imputation by Chained Equations (ICE)

18 / 99

Impute via mean/mode

Init.

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F

Age	Inc.	Gen.
34	13	F
18	12	F
26	14	M
26	13	F

Set values of Age
originally missing
as unknown

Fit a predictive model on
complete observation
to predict Age

Use the fitted model
to impute ?

1st step
Age

Age	Inc.	Gen.
34	13	F
18	12	F
?	14	M
?	13	F

Age	Inc.	Gen.
34	13	F
18	12	F
?	14	M
?	13	F

Age	Inc.	Gen.
34	13	F
18	12	F
50	14	M
34	13	F

Imputation by Chained Equations (ICE)

18 / 99

Impute via mean/mode

Init.

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F

Age	Inc.	Gen.
34	13	F
18	12	F
26	14	M
26	13	F

Set values of Inc.
originally missing
as unknown

'Inc.'
step

Age	Inc.	Gen.
34	?	F
18	12	F
50	14	M
34	?	F

Imputation by Chained Equations (ICE)

18 / 99

Impute via mean/mode

Init.

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F

Age	Inc.	Gen.
34	13	F
18	12	F
26	14	M
26	13	F

Set values of Inc.
originally missing
as unknown

Fit a predictive model on
complete observation
to predict Inc.

'Inc.'
step

Age	Inc.	Gen.
34	?	F
18	12	F
50	14	M
34	?	F

Age	Inc.	Gen.
34	?	F
18	12	F
50	14	M
34	?	F

Imputation by Chained Equations (ICE)

18 / 99

Impute via mean/mode

Init.

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F

Age	Inc.	Gen.
34	13	F
18	12	F
26	14	M
26	13	F

Set values of Inc.
originally missing
as unknown

Fit a predictive model on
complete observation
to predict Inc.

Use the fitted model
to impute ?

'Inc.'
step

Age	Inc.	Gen.
34	?	F
18	12	F
50	14	M
34	?	F

Age	Inc.	Gen.
34	?	F
18	12	F
50	14	M
34	?	F

Age	Inc.	Gen.
34	12	F
18	12	F
50	14	M
34	12	F

Imputation by Chained Equations (ICE)

18 / 99

Impute via mean/mode

Init.

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F

Age	Inc.	Gen.
34	13	F
18	12	F
26	14	M
26	13	F

Set values of Gen.
originally missing
as unknown

'Gen.'
step

Age	Inc.	Gen.
34	12	F
18	12	?
50	14	M
34	12	F

Imputation by Chained Equations (ICE)

18 / 99

Impute via mean/mode

Init.

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F

Age	Inc.	Gen.
34	13	F
18	12	F
26	14	M
26	13	F

Set values of Gen.
originally missing
as unknown

Fit a predictive model on
complete observation
to predict Gen.

'Gen.'
step

Age	Inc.	Gen.
34	12	F
18	12	?
50	14	M
34	12	F

Age	Inc.	Gen.
34	12	F
18	12	?
50	14	M
34	12	F

Imputation by Chained Equations (ICE)

18 / 99

Impute via mean/mode

Init.

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F

Age	Inc.	Gen.
34	13	F
18	12	F
26	14	M
26	13	F

Set values of Gen.
originally missing
as unknown

Fit a predictive model on
complete observation
to predict Gen.

Use the fitted model
to impute ?

'Gen.'
step

Age	Inc.	Gen.
34	12	F
18	12	?
50	14	M
34	12	F

Age	Inc.	Gen.
34	12	F
18	12	?
50	14	M
34	12	F

Age	Inc.	Gen.
34	12	F
18	12	F
50	14	M
34	12	F

Imputation by Chained Equations (ICE)

19 / 99

	Age	Inc.	Gen.		Age	Inc.	Gen.		Age	Inc.	Gen.
Init.	34	NA	F		34	13	F		34	13	F
	18	12	NA		18	12	F		18	12	F
	NA	14	M		26	14	M		50	14	M
	NA	NA	F		26	13	F		34	13	F
'Age' step	34	13	F		34	13	F		34	13	F
	18	12	F		18	12	F		18	12	F
	?	14	M		?	14	M		50	14	M
	?	13	F		?	13	F		34	13	F
'Inc.' step	34	?	F		34	?	F		34	12	F
	18	12	F		18	12	F		18	12	F
	50	14	M		50	14	M		50	14	M
	34	?	F		34	?	F		34	12	F
'Gen.' step	34	12	F		Age	Inc.	Gen.		Age	Inc.	Gen.
	18	12	?		34	12	F		34	12	F
	50	14	M		18	12	?		18	12	?
	34	12	F		50	14	M		50	14	M
					34	12	F		34	12	F

Hyperparameters - R implementation

20 / 99

- ▷ Initialization
- ▷ Number of cycles
- ▷ Ordering of variables: same order, random order...

- ▷ Initialization
- ▷ Number of cycles
- ▷ Ordering of variables: same order, random order...
- ▷ Predictive models
 - ◊ Predictive mean matching (numeric data)²⁷
 - ◊ Logistic regression imputation (binary data)²⁸
 - ◊ Multinomial regression imputation (unordered categorical data)
 - ◊ Proportional odds model (ordered categorical data) ²⁹

²⁷<https://stefvanbuuren.name/fimd/sec-pmm.html>

²⁸<https://www.rdocumentation.org/packages/mice/versions/3.17.0/topics/mice.impute.logreg>

²⁹<https://online.stat.psu.edu/stat504/lesson/8/8.4>

- ▷ Initialization
- ▷ Number of cycles
- ▷ Ordering of variables: same order, random order...
- ▷ Predictive models
 - ◊ Predictive mean matching (numeric data)
 - ◊ Logistic regression imputation (binary data)
 - ◊ Multinomial regression imputation (unordered categorical data)
 - ◊ Proportional odds model (ordered categorical data)

Logistic regression imputation - Bayesian logistic regression

- ▷ Fit a logistic model on the data
- ▷ Construct $\hat{\beta}$ and an estimation of its covariance matrix $\hat{\Sigma}$.
- ▷ Draw $\tilde{\beta} \sim \mathcal{N}(\hat{\beta}, \hat{\Sigma})$.
- ▷ Compute the predicted score as $\sigma(X^\top \tilde{\beta})$.
- ▷ Impute by drawing a Bernoulli with parameter $\sigma(X^\top \tilde{\beta})$.

- ▷ Initialization
- ▷ Number of cycles
- ▷ Ordering of variables: same order, random order...
- ▷ Predictive models
 - ◊ Predictive mean matching (numeric data)
 - ◊ Logistic regression imputation (binary data)
 - ◊ Multinomial regression imputation (unordered categorical data)
 - ◊ Proportional odds model (ordered categorical data)

Predictive mean matching

- ▷ Fit a linear model on the data
- ▷ Construct $\hat{\beta}$ and an estimation of its covariance matrix $\hat{\Sigma}$.
- ▷ Draw $\tilde{\beta} \sim \mathcal{N}(\hat{\beta}, \hat{\Sigma})$.
- ▷ Compute the predicted scores as $X^\top \tilde{\beta}$.
- ▷ Find the $k = 5$ observations for which $X_i^\top \hat{\beta}$ is the closest to $X^\top \tilde{\beta}$
- ▷ Impute by drawing uniformly at random one observations among the k selected observations.

- ▷ Initialization
- ▷ Number of cycles
- ▷ Ordering of variables: same order, random order...
- ▷ Predictive models
 - ◊ Predictive mean matching (numeric data)
 - ◊ Logistic regression imputation (binary data)
 - ◊ Multinomial regression imputation (unordered categorical data)
 - ◊ Proportional odds model (ordered categorical data)

Random forests - Mice.RF

- ▷ Fit a random forest on the data
- ▷ For a given 'missing' observation, put it down each tree and collect all observations in all leaves
- ▷ Impute by drawing at random an observation among the previous set

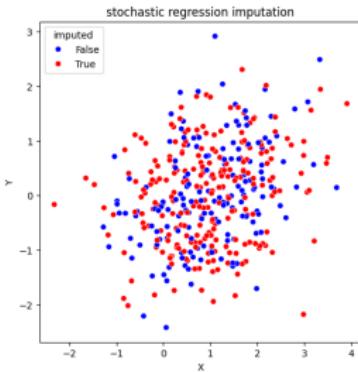
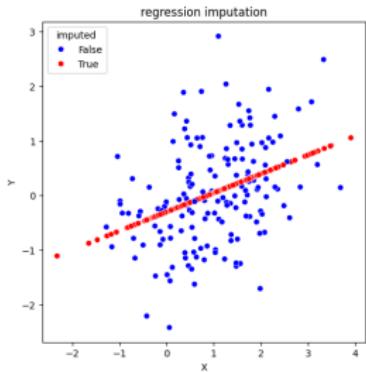
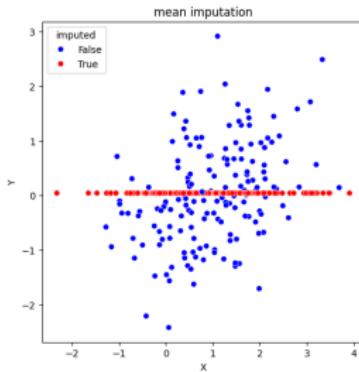
- ▷ Initialization
- ▷ Number of cycles
- ▷ Ordering of variables: same order, random order...
- ▷ Predictive models
 - ◊ Predictive mean matching (numeric data)
 - ◊ Logistic regression imputation (binary data)
 - ◊ Multinomial regression imputation (unordered categorical data)
 - ◊ Proportional odds model (ordered categorical data)

Random forests - MissForest

- ▷ Fit a random forest on the data
- ▷ Impute by predicting the value output by the RF

1. Missing values mechanism
2. Single Imputation
3. Multiple Imputation
4. Imputation quality
5. Supervised Learning with Missing values
 - Decision trees as PbP predictors
 - Impute-then-regress procedures with consistent predictors
6. Linear models
 - Linear regression: A pattern-by-pattern approach
 - Linear regression: Impute-then-regress procedures via zero-imputation
 - Classification with missing values
7. Conclusion

Single imputation methods



$$\mu_y = 0$$

$$\sigma_y = 1$$

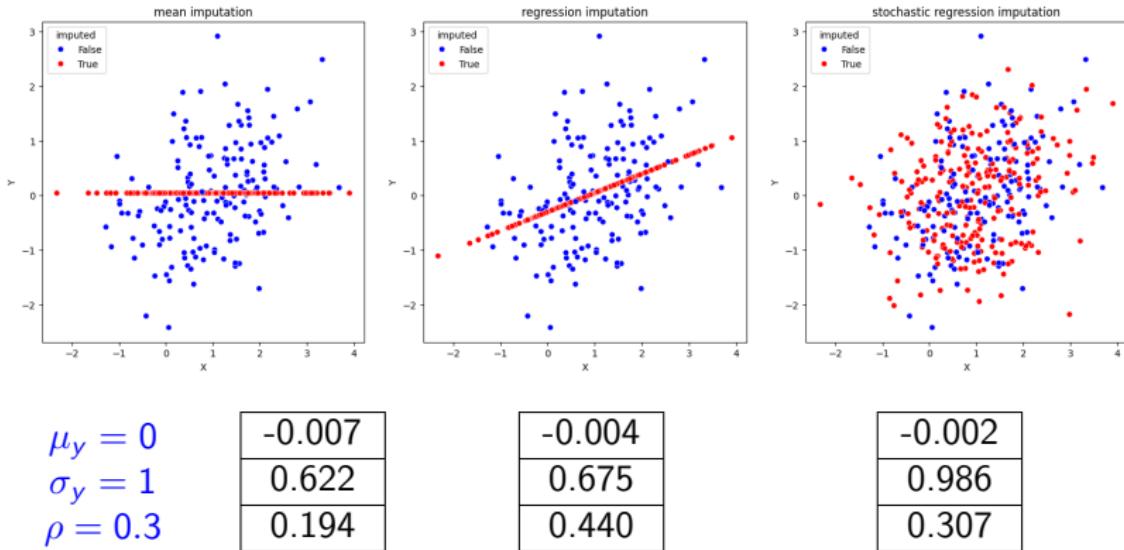
$$\rho = 0.3$$

-0.007
0.622
0.194

-0.004
0.675
0.440

-0.002
0.986
0.307

Single imputation methods



How to build confidence intervals for μ_y ?

Let $Y = (Y_1, \dots, Y_n)'$ be i.i.d. independent Gaussian $\mathcal{N}(\mu_y, \sigma_y^2)$.

▷ Unknown variance:

$$\frac{\hat{\mu}_y - \mu_y}{\hat{\sigma}_{\hat{\mu}_y}} \sim T(n-1)$$

▷ Unknown variance:

$$\sqrt{n} \left(\frac{\hat{\mu}_y - \mu_y}{\hat{\sigma}_y} \right) \sim T(n-1)$$

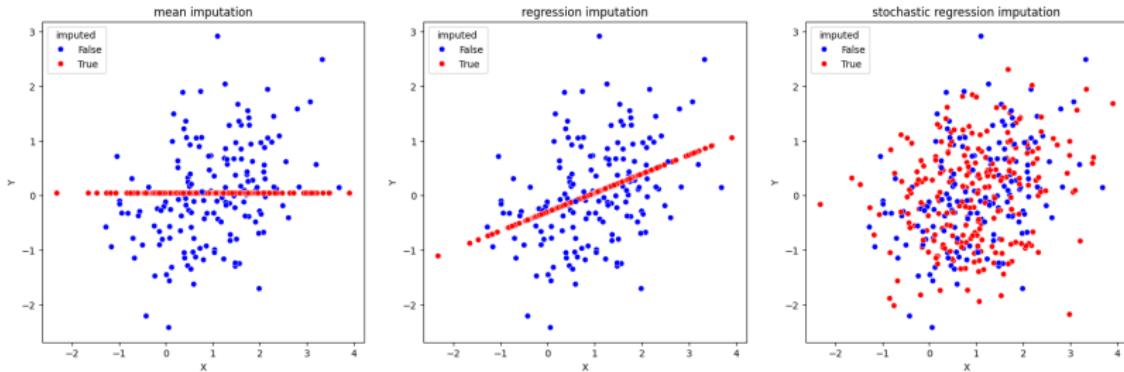
▷ CI for μ_y at level α : $\left[\hat{\mu}_y - \frac{\hat{\sigma}_y}{\sqrt{n}} qt_{1-\alpha/2}(n-1), \hat{\mu}_y + \frac{\hat{\sigma}_y}{\sqrt{n}} qt_{1-\alpha/2}(n-1) \right]$

Simulation - Computing coverage

1. Generate bivariate Gaussian data ($\mu_y = 0, \sigma_y = 1, \rho = 0.6$)
2. Put MCAR missing values on y and impute missing entries
3. Compute the confidence interval of μ_y
4. Count if the true value $\mu_y = 0$ is in the confidence interval
5. Repeat the steps 1-4, 10000 times

Code available on Rmיסטатик.

Single imputation methods: Danger!



$$\mu_y = 0$$

$$\sigma_y = 1$$

$$\rho = 0.3$$

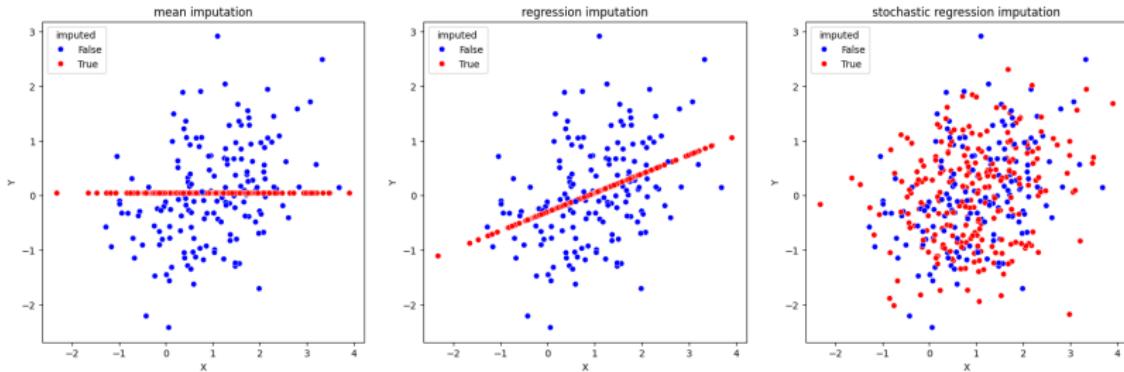
$$CI\mu_y 95\%$$

-0.007
0.622
0.194

-0.004
0.675
0.440

-0.002
0.986
0.307

Single imputation methods: Danger!



$$\mu_y = 0$$

$$\sigma_y = 1$$

$$\rho = 0.3$$

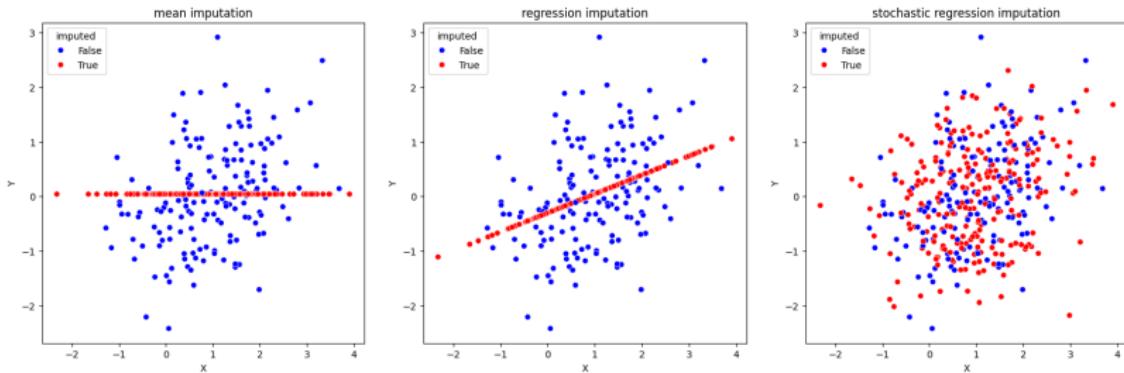
$$CI \mu_y 95\%$$

-0.007
0.622
0.194
55.0

-0.004
0.675
0.440
60.3

-0.002
0.986
0.307
73.3

Single imputation methods: Danger!



$$\mu_y = 0$$

$$\sigma_y = 1$$

$$\rho = 0.3$$

$$CI \mu_y 95\%$$

-0.007
0.622
0.194
55.0

-0.004
0.675
0.440
60.3

-0.002
0.986
0.307
73.3

⇒ Standard errors $\hat{\sigma}_{\hat{\mu}_y}$ based on the imputed data set are underestimated

The idea of imputation is both seductive and dangerous (Dempster and Rubin, 1983)

Asymptotic confidence interval for μ_y : $\left[\hat{\mu}_y - z_{\alpha/2} \frac{\hat{\sigma}_y}{\sqrt{n}}; \hat{\mu}_y + z_{1-\alpha/2} \frac{\hat{\sigma}_y}{\sqrt{n}} \right]$

Consider MCAR values and

- ▷ Impute missing values on via (stochastic) linear regression
- ▷ $\hat{\mu}_y$ is the average of y computed on the imputed data set

Asymptotic variance (Little & Rubin, 2019. p158)

$$\text{Var}[\hat{\mu}_y - \mu_y] \simeq \frac{\hat{\sigma}_y^2}{n_{full}} \left(1 - \hat{\rho}^2 \frac{n - n_{full}}{n} \right),$$

where $\hat{\sigma}_y$ is estimated on the complete observations only and n_{full} the number of complete observations.

- ▷ If there are few missing data ($n_{full} \sim (n)$), then $\text{Var}[\hat{\mu}_y - \mu_y] \sim \hat{\sigma}_y^2/n$, the ACI has the correct asymptotic coverage (Idem if $\rho = 1$).
- ▷ But, in general, **coverage of single imputation is too low**: need to take into account the uncertainty associated to the predictions.

Multiple imputation: correct standard errors

- 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)²⁷:

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

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

²⁷ see Chapter 14 of Semiparametric Theory and Missing Data. A.A. Tsiatis. 2006.

MI based on stochastic regression

1. Generate M imputed data sets: for $m = 1, \dots, M$,
 - ▷ draw \hat{y}_i from $\mathcal{N}(x_i\hat{\beta}, \hat{\sigma}^2)$
2. Performe the analysis on each imputed data set
3. Compute the variance (= within + between imputation variance)

²⁸Code available on Rmisticic.

MI based on stochastic regression

1. Generate M imputed data sets: for $m = 1, \dots, M$,
 - ▷ draw \hat{y}_i from $\mathcal{N}(x_i\hat{\beta}, \hat{\sigma}^2)$
2. Performe the analysis on each imputed data set
3. Compute the variance (= within + between imputation variance)

	$M = 1$	$M = 50$
$\mu_y = 0$	-0.002	-0.02
$\sigma_y = 1$	0.986	0.936
$\rho = 0.3$	0.307	0.314
$CI_{\mu_y 95\%}$	73.3	82.0

²⁸Code available on Rmisticat.

MI based on stochastic regression

1. Generate M imputed data sets: for $m = 1, \dots, M$,
 - ▷ draw \hat{y}_i from $\mathcal{N}(x_i\hat{\beta}, \hat{\sigma}^2)$
2. Performe the analysis on each imputed data set
3. Compute the variance (= within + between imputation variance)

	$M = 1$	$M = 50$
$\mu_y = 0$	-0.002	-0.02
$\sigma_y = 1$	0.986	0.936
$\rho = 0.3$	0.307	0.314
$CI_{\mu_y} 95\%$	73.3	82.0

- ▷ Variability of the parameters is missing: "improper" imputation
- ▷ Prediction variance = estimation variance plus noise

²⁸Code available on Rmisticat.

MI based on stochastic regression

1. Generate M imputed data sets: for $m = 1, \dots, M$,
 - ▷ Generate $\hat{\beta}^1, \dots, \hat{\beta}^M$ by bootstrap or via posterior distribution (Data Augmentation, Tanner & Wong, 1987))
 - ▷ Impute missing values \hat{y}_i^m by drawing $\mathcal{N}(x_i \hat{\beta}^m, (\hat{\sigma}^2)^m)$
2. Performe the analysis on each imputed data set
3. Compute the variance (= within + between imputation variance)

	$M = 1$	$M = 50$	$M = 50$ with boot.
$\mu_y = 0$	-0.002	-0.02	-0.006
$\sigma_y = 1$	0.986	0.936	1.036
$\rho = 0.3$	0.307	0.314	0.295
$CI\mu_y^{95\%}$	73.3	82.0	98.0

²⁸Code available on Rmisticat.

⇒ Aim: provide an estimation of all parameters with their estimated variance.

Parametric Multiple imputation

1. Generating M imputed data sets, taking into account:
 - ▷ structural noise (e.g. σ^2 via stochastic regression)
 - ▷ parameter variance (e.g. via bootstrapping)
2. Performing the analysis on each imputed data set^a,
3. Compute the variance (= within + between imputation variance)

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

^aThe analysis model may be "in agreement" with the imputation model: congeniality.

⇒ Aim: provide an estimation of all parameters with their estimated variance.

NonParametric Multiple imputation

1. Generating M imputed data sets, taking into account:
 - ▷ structural noise (e.g. σ^2 via stochastic regression)
 - ▷ parameter variance (e.g. via bootstrapping)
2. Performing the analysis on each imputed data set^a,
3. Aggregate the result of each analysis (e.g. taking the mean of predicted output values)

^aThe analysis model may be "in agreement" with the imputation model: congeniality.

Multiple Imputation with joint modeling

\Rightarrow Hypothesis $x_i \sim \mathcal{N}(\mu, \Sigma)$

Expectation Maximization Bootstrap

1. Bootstrap rows: X^1, \dots, X^M
2. EM algorithm: $(\hat{\mu}^1, \hat{\Sigma}^1), \dots, (\hat{\mu}^M, \hat{\Sigma}^M)$
3. Imputation: $\hat{x}_{i,miss}^m$ drawn from $\mathcal{N}\left(\hat{\mu}_{miss|obs}^m, \hat{\Sigma}_{miss|obs}^m\right)$

Easy to parallelized. Implemented in **Amelia** ([website](#))



Amelia Earhart



James Honaker



Gary King



Matt Blackwell

- Impute variables 1 by 1 using all other variables as inputs (round-robin)
- One model/variable: flexible for different types of variables
- Cycle through variables: iteratively refining imputations

MICE

1. Initial imputation: mean imputation
2. For a variable j
 - $(\hat{\beta}_{-j}, \hat{\sigma}_{-j})$ drawn from a **Bootstrap**: $(\hat{\beta}_{-j}^1, \hat{\sigma}_{-j}^1), \dots, (\hat{\beta}_{-j}^M, \hat{\sigma}_{-j}^M)$
 - Impute X_j^m via **stochastic regression** $\mathcal{N}\left((x_{i,-j})' \hat{\beta}_{-j}^m, \hat{\sigma}_{-j}^m\right)$
3. Cycling through variables

⇒ With continuous variables & regression/variable: gibbs $\mathcal{N}(\mu, \Sigma)$ ^{30 31}

"There is no clear-cut method for determining whether MICE has converged"

Implemented in R package **mice** & **IterativeImputer** from scikitlearn (default iterative ridge regression)



Stef van Buuren

³⁰ Monte Carlo statistical methods (Robert, Casella, 2004) (p344),

³¹ The EM algorithm and extensions (McLachlan, et al. 1998) (p243)

³² van Buuren. 2018. Flexible Imputation of Missing Data. Second Edition. CRC Press

Conditional modeling takes the lead?

- ▷ Flexible: one model/variable. Easy to deal with interactions and variables of different nature (binary, ordinal, categorical...)
- ▷ Many statistical models are conditional models
- ▷ Tailor to your data - Super powerful in practice
- ⇒ Drawbacks: one model/variable. **Computational costly^a**

^aImprovement on mice pmm for large sample size, see mice github repo - still costly for large d

What to do with high correlation or when $n < p$

- ▷ JM shrinks the covariance $\Sigma + k\mathbb{I}$ (selection of k ?)
- ▷ CM: ridge regression or predictors selection/variable

Challenges with multiple imputation

- ▷ MI in high dimension? Theory with small n , large p ?
- ▷ Aggregating lasso regressions? clustering?

1. Missing values mechanism
2. Single Imputation
3. Multiple Imputation
4. Imputation quality
5. Supervised Learning with Missing values
 - Decision trees as PbP predictors
 - Impute-then-regress procedures with consistent predictors
6. Linear models
 - Linear regression: A pattern-by-pattern approach
 - Linear regression: Impute-then-regress procedures via zero-imputation
 - Classification with missing values
7. Conclusion

How to evaluate imputation quality?

33 / 99

- ▷ Aim: imputed data must resemble complete data.

Original data set

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F
34	NA	M
22	28	F
29	10	NA
34	NA	F
80	NA	NA
68	15	F

How to evaluate imputation quality?

- ▷ Aim: imputed data must resemble complete data.

Original data set

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F
34	NA	M
22	28	F
29	10	NA
34	NA	F
80	NA	NA
68	15	F

Imputed data set

Age	Inc.	Gen.
34	13	F
18	12	NA
30	14	M
30	13	F
34	13	M
22	28	F
29	10	NA
34	13	F
80	13	F
68	15	F

How to evaluate imputation quality?

33 / 99

- ▷ Aim: imputed data must resemble complete data.

Original data set

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F
34	NA	M
22	28	F
29	10	NA
34	NA	F
80	NA	NA
68	15	F

Imputed data set

Age	Inc.	Gen.
34	13	F
18	12	NA
30	14	M
30	13	F
34	13	M
22	28	F
29	10	NA
34	13	F
80	13	F
68	15	F

What is the quality of data imputation?

How to evaluate imputation quality?

33 / 99

- ▷ Aim: imputed data must resemble complete data.

Original data set

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F
34	NA	M
22	28	F
29	10	NA
34	NA	F
80	NA	NA
68	15	F

How to evaluate imputation quality?

33 / 99

- ▷ Aim: imputed data must resemble complete data.

Original data set

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F
34	NA	M
22	28	F
29	10	NA
34	NA	F
80	NA	NA
68	15	F

How to evaluate imputation quality?

- ▷ Aim: imputed data must resemble complete data.

Original data set

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F
34	NA	M
22	28	F
29	10	NA
34	NA	F
80	NA	NA
68	15	F

Additional missing values

Age	Inc.	Gen.
34	NA	F
NA	NA	NA
NA	14	NA
NA	NA	F
34	NA	M
22	NA	F
NA	10	NA
34	NA	F
80	NA	NA
68	NA	NA

How to evaluate imputation quality?

- ▷ Aim: imputed data must resemble complete data.

Original data set

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F
34	NA	M
22	28	F
29	10	NA
34	NA	F
80	NA	NA
68	15	F

Additional missing values

Age	Inc.	Gen.
34	NA	F
NA	NA	NA
NA	14	NA
NA	NA	F
34	NA	M
22	NA	F
NA	10	NA
34	NA	F
80	NA	NA
68	NA	NA

Imputed missing values

Age	Inc.	Gen.
34	12	F
46	12	NA
46	14	M
46	12	F
34	12	M
22	12	F
46	10	NA
34	12	F
80	12	NA
68	12	F

How to evaluate imputation quality?

- ▷ Aim: imputed data must resemble complete data.

Original data set

Age	Inc.	Gen.
34	NA	F
18	12	NA
NA	14	M
NA	NA	F
34	NA	M
22	28	F
29	10	NA
34	NA	F
80	NA	NA
68	15	F

Additional missing values

Age	Inc.	Gen.
34	NA	F
NA	NA	NA
NA	14	NA
NA	NA	F
34	NA	M
22	NA	F
NA	10	NA
34	NA	F
80	NA	NA
68	NA	NA

Imputed missing values

Age	Inc.	Gen.
34	12	F
46	12	NA
46	14	M
46	12	F
34	12	M
22	12	F
46	10	NA
34	12	F
80	12	F
68	12	F

Compared initial vs imputed values via predictive metrics (MSE, MAE...)

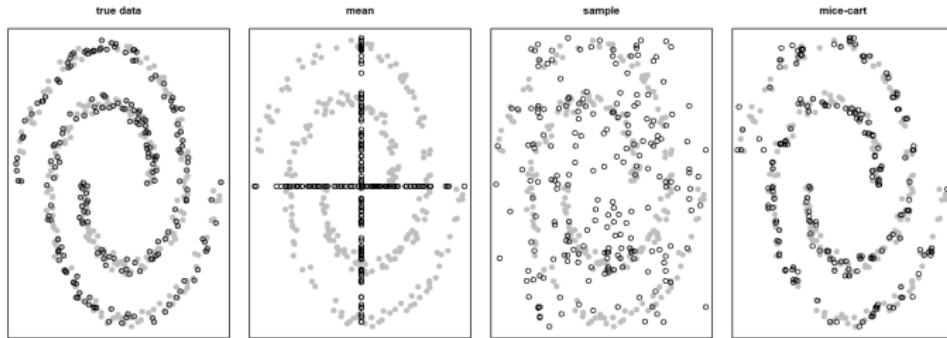
Measures related to imputation quality

34 / 99

Pointwise predictive measure such as MSE rank highest imputation close to the conditional expectation

▷ Favor imputation with small variability

Imputation is a distributional task so one should use distributional measures³³³⁴ to assess its quality.



³³Székely & Rizzo. Energy statistics *Journal of stat. planning & inference*. 2013

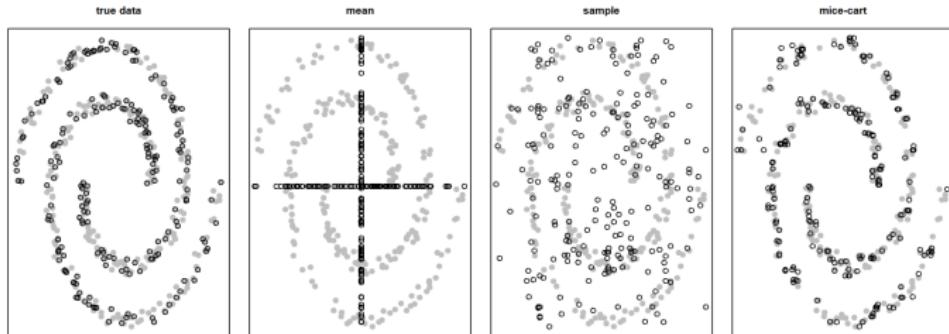
³⁴Gneiting, Raftery, Strictly Proper Scoring Rules, Prediction, and Estimation, *JASA*, 2007

Measures related to imputation quality

Pointwise predictive measure such as MSE rank highest imputation close to the conditional expectation

- ▷ Favor imputation with small variability

Imputation is a distributional task so one should use distributional measures³³³⁴ to assess its quality.



Imputation method	Mean	Sample	Mice-CART
Renormalized RMSE	0	-0.18	-0.22

³³Székely & Rizzo. Energy statistics *Journal of stat. planning & inference*. 2013

³⁴Gneiting, Raftery, Strictly Proper Scoring Rules, Prediction, and Estimation, *JASA*, 2007

- ▷ Energy score (distribution vs a point)

$$es(H, x) = \frac{1}{2} \mathbb{E}_{X, X' \sim H} [\|X - X'\|_{\mathbb{R}^d}] - \mathbb{E}_{X \sim H} [\|X - x\|_{\mathbb{R}^d}]$$

- ▷ The **energy** score can be used to score **distributional prediction/imputation**

Controlled simulation setting

- ▷ Generate complete data
- ▷ Mask some data according to MCAR/MAR/MNAR mechanism
- ▷ Learn a distributional imputation method H
- ▷ For any $x \in \mathbb{R}^d$, sample imputed values from H to estimate $es(H, x)$
- ▷ Average over $X \sim P^*$ (complete data distribution) to estimate

$$S(H, P^*) := \mathbb{E}_{Y \sim P^*} [es(H, Y)]$$

- ▷ The question of how to evaluate imputation methods becomes much harder when the **true underlying values are not available**.

A new procedure

36 / 99

- ▷ Consider a distribution κ on the subsets of $\{1, \dots, d\}$
- ▷ For each $A \subset \{1, \dots, d\}$, we let P_A^M be the marginal distribution of M on A . We denote $M_A \sim P_A^M$.
- ▷ We also let $H_A|M_A = m_A$, i.e. the distribution of an imputation H , given the missingness pattern m_A on the projection A .

A new procedure

- ▷ Consider a distribution κ on the subsets of $\{1, \dots, d\}$
- ▷ For each $A \subset \{1, \dots, d\}$, we let P_A^M be the marginal distribution of M on A . We denote $M_A \sim P_A^M$.
- ▷ We also let $H_A|M_A = m_A$, i.e. the distribution of an imputation H , given the missingness pattern m_A on the projection A .

Imputation score of imputation H

$$S_{NA}^*(H, P) = \mathbb{E}_{A \sim \kappa, M_A \sim P_A^M, X_A \sim H_{M_A}} \left[\log \left(\frac{p_A(X_A | M_A = \mathbf{0})}{h_{M_A}(X_A)} \right) \right].$$

Group observations into J groups according to their missing data pattern M_1, \dots, M_J .

Procedure

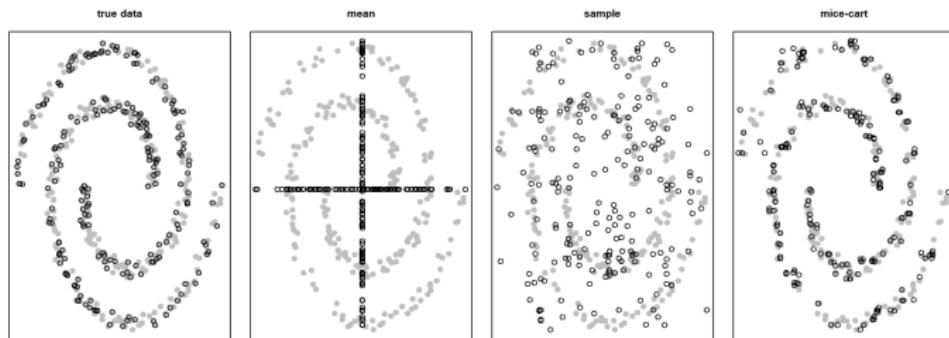
For each missing pattern m among M_1, \dots, M_J

1. Choose `num.proj` projections on $\{1, \dots, d\}$ such that each projection contains at least one observed and one missing component.
2. Obtain the imputed data from pattern m , denoted by \hat{X}_i . Split them into two halves \hat{X}_i^0 and \hat{X}_i^1 .
3. For each projection A_k ($k = 1, \dots, \text{num.proj}$),
 - a) Get the complete data $X_{A_k}^{\text{comp}}$ from the projected data X_{A_k}
 - b) Get the projected imputed data \hat{X}_{i,A_k}^0
 - c) Fit a forest with `num.trees.per.proj` to discriminate $X_{A_k}^{\text{comp}}$ from \hat{X}_{i,A_k}^0 (ensuring balanced classes).
4. Aggregate all forests and let $\hat{g}_A(x)$ be the probability output by the forest at x .
5. Compute the individual scores $\log \hat{g}_A(x)$ for $x \in \hat{X}_i^1$
6. Average all scores across all observations, missing patterns and imputed data sets (multiple imputation) to get the final imputation score.

Measures related to imputation quality

37 / 99

Imputation is a distributional task so one should use distributional measures³⁵³⁶ to assess its quality.



Imputation method	Mean	Sample	Mice-CART
Renormalized RMSE	0	-0.18	-0.22
Renormalized Energy score			

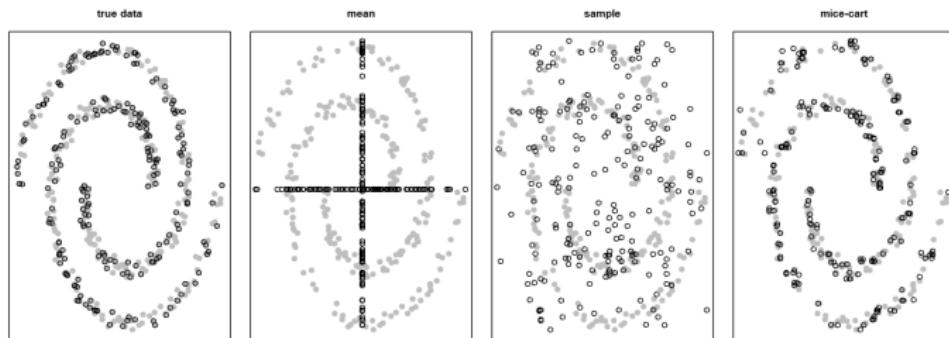
³⁵Székely & Rizzo. Energy statistics *Journal of stat. planning & inference*. 2013

³⁶Gneiting, Raftery, Strictly Proper Scoring Rules, Prediction, and Estimation, *JASA*, 2007

Measures related to imputation quality

37 / 99

Imputation is a distributional task so one should use distributional measures³⁵³⁶ to assess its quality.



Imputation method	Mean	Sample	Mice-CART
Renormalized RMSE	0	-0.18	-0.22
Renormalized Energy score	-22.4	-1.39	0

³⁵Székely & Rizzo. Energy statistics *Journal of stat. planning & inference*. 2013

³⁶Gneiting, Raftery, Strictly Proper Scoring Rules, Prediction, and Estimation, *JASA*, 2007

Major characteristics of imputations

Imputation should

- (1) be a distributional regression method,
- (2) be able to capture nonlinearities in the data,
- (3) be able to deal with distributional shifts in the observed variables,

- ▷ Conditional and marginal **distribution shifts** can occur for different patterns under MAR
- ▷ Conditional shifts are handled with FCS

Method	(1)	(2)	(3)
missForest (Stekhoven & Bühlmann, 2011)		✓	
mice-cart (Burgette & Reiter, 2010)	✓	✓	
mice-RF (Doove et al., 2014)	✓	✓	
mice-DRF (Näf et al., 2024)	✓	✓	
mice-norm.nob (Gaussian)	✓		✓
mice-norm.predict (Regression)			✓

MAR with shift in cond. distribution between patterns³⁹/⁹⁹

- Example: two patterns $m_1 = (0, 0)$ and $m_2 = (1, 0)$, with $\Sigma = ((2, 1), (1, 1))$ and **a shift**:

$$X \mid M = m_1 \sim N((0, 0), \Sigma)$$

$$X \mid M = m_2 \sim N((5, 5), \Sigma).$$

MAR with shift in cond. distribution between patterns^{39 / 99}

- Example: two patterns $m_1 = (0, 0)$ and $m_2 = (1, 0)$, with $\Sigma = ((2, 1), (1, 1))$ and a **shift**:

$$X \mid M = m_1 \sim N((0, 0), \Sigma)$$

$$X \mid M = m_2 \sim N((5, 5), \Sigma).$$

- A special case of MAR: conditional distributions are the same across patterns:

$$X_1 \mid X_2, M = m_1 = X_1 \mid X_2, M = m_2.$$

Definition (Conditional indep. MAR - CIMAR)

For all $m, m' \in \mathcal{M}, x \in \mathcal{X}$,

$$p^*(o^c(x, m) \mid o(x, m), M = m') = \textcolor{red}{p^*(o^c(x, m) \mid o(x, m))}.$$

MAR with shift in cond. distribution between patterns^{39 / 99}

- Example: two patterns $m_1 = (0, 0)$ and $m_2 = (1, 0)$, with $\Sigma = ((2, 1), (1, 1))$ and a **shift**:

$$X \mid M = m_1 \sim N((0, 0), \Sigma)$$

$$X \mid M = m_2 \sim N((5, 5), \Sigma).$$

- A special case of MAR: conditional distributions are the same across patterns:

$$X_1 \mid X_2, M = m_1 = X_1 \mid X_2, M = m_2.$$

Definition (Conditional indep. MAR - CIMAR)

For all $m, m' \in \mathcal{M}, x \in \mathcal{X}$,

$$p^*(o^c(x, m) \mid o(x, m), M = m') = \textcolor{red}{p^*(o^c(x, m) \mid o(x, m))}.$$

Beware! Even in this case, the joint distribution varies across pattern, since the marginal distribution of X_2 changes

Forests generalize poorly outside of the training set

40 / 99

- Example: two patterns $m_1 = (0, 0)$ and $m_2 = (1, 0)$, with $\Sigma = ((2, 1), (1, 1))$ and a shift $X | M = m_1 \sim N((0, 0), \Sigma)$, $X | M = m_2 \sim N((5, 5), \Sigma)$.

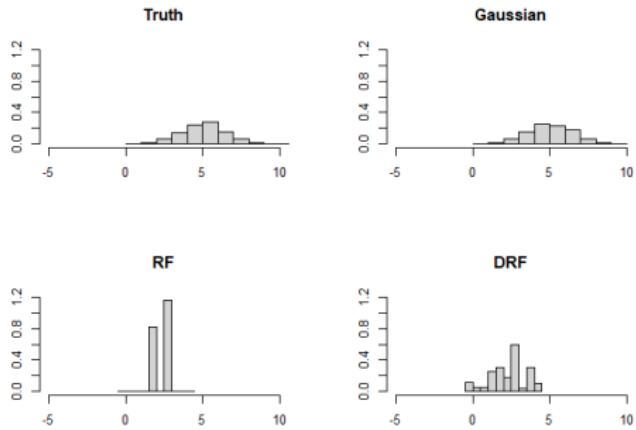


Figure: True distribution against a draw from different imputation methods.

- DRF, a distributional method, fails to deal with covariate shift
- ▷ Imputation should be centered around 5.

MAR with shifts in cond. distribution between patterns

Consider $X \in \mathbb{R}^3$ with three different missing patterns:

$$m_1 = (0, 0, 0), \quad m_2 = (1, 0, 0) \quad \text{and} \quad m_3 = (1, 1, 0).$$

MCAR: No change allowed.

For all $m, m' \in \mathcal{M}, x \in \mathcal{X}$, $p^*(x) = p^*(x | M = m) = p^*(x | M = m')$

CIMAR: No conditional changes allowed

$p^*(x_1, x_2 | x_3, M = m_1) = p^*(x_1, x_2 | x_3, M = m_2) = p^*(x_1, x_2 | x_3, M = m_3) =$
 $\color{red}{p^*(x_1, x_2 | x_3)}$

Distrib. of $X_1, X_2 | X_3$ is not allowed to change from one pattern to another, though the marginal distrib. of X_3 can change.

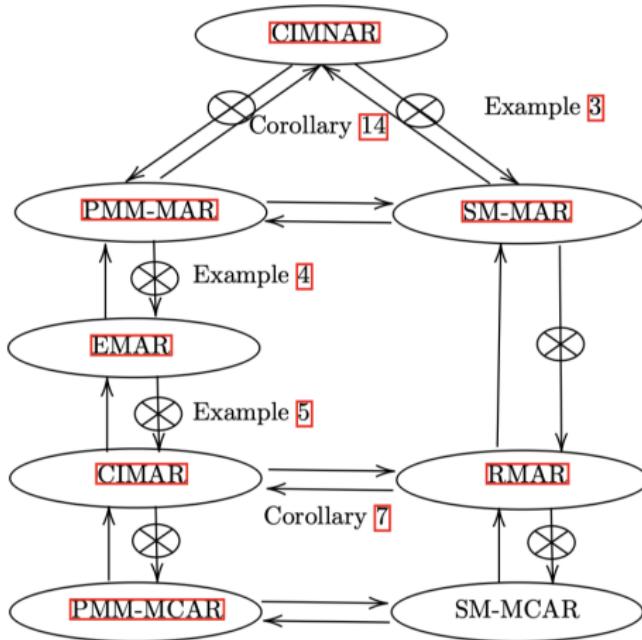
PMM-MAR: many changes allowed

$p^*(x_1, x_2 | x_3, M = m_3) = \color{red}{p^*(x_1, x_2 | x_3)}$

Both distrib. of observed variables and conditional ones can change from pattern to pattern.

Relationships between the M(N)AR conditions³⁷

42 / 99

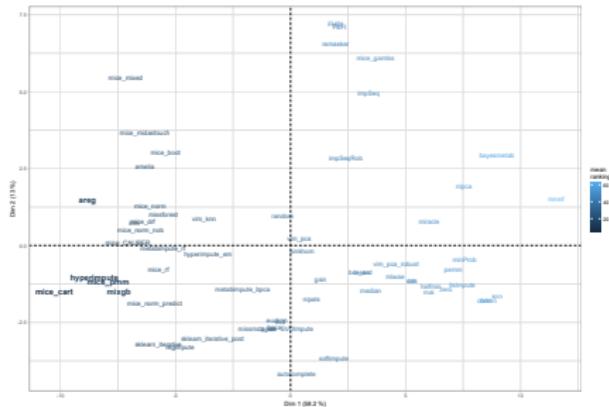


³⁷Naf, Scornet J.. (2024). What is a good imputation under MAR. Submitted.

Benchmarking imputation methods

43 / 99

- ▷ 65 methods (R & Python)
- ▷ 14 datasets: 100-50000 observations and 3-400 features
- ▷ 10-30 % NA MCAR, MAR, Standardized energy distance



- ▷ Mice-cart³⁸, aregImpute (close to mice+splines+pmm)³⁹, Hyperimpute (mice + model selection RF, XGBoost, Logistic Reg., etc)⁴⁰, Mice mixed⁴¹

³⁸ Buuren & Groothuis-O. (2011). Multivariate imputation by chained equations in R. *JSS*.

³⁹ Harrell & Dupont (2018). Hmisc: Harrell miscellaneous. R package. *Stat. Comput.*

⁴⁰ Jarrett et al. (2022). Hyperimpute: Gen. iter. imput. with automatic model selection. *ICML*.

⁴¹ Varga (2020). missCompare: Intuitive Missing Data Imputation. R package. *Stat. Comput.*

Take home message on inference & imputation

44 / 99

- ▷ Different missing data scenario designed for likelihood inference (e.g. EM algorithm) but that can be very complex (distribution shift in MAR).
- ▷ Use single imputation only for point estimates
- ▷ In general, look for an imputation that preserve the joint distribution of the data
- ▷ Compare imputation methods with distributional metrics like energy distance
- ▷ **Multiple imputation** aims at estimating the parameters and their variability taking into account **the uncertainty of the missing values**
- ▷ Use Multiple imputation to get confidence intervals
- ▷ mice-DRF promising (code available) - mice-Engression⁴²

⁴²Shen & Meinshausen (2024). Engression: extrapolation through the lens of distributional regression. *JRSS B*.

1. Missing values mechanism
2. Single Imputation
3. Multiple Imputation
4. Imputation quality
5. Supervised Learning with Missing values
 - Decision trees as PbP predictors
 - Impute-then-regress procedures with consistent predictors
6. Linear models
 - Linear regression: A pattern-by-pattern approach
 - Linear regression: Impute-then-regress procedures via zero-imputation
 - Classification with missing values
7. Conclusion

Formalizing the problem

- ▷ **Assumption** - The response Y is a function of the (unavailable) complete data plus some noise:

$$Y = f^*(\textcolor{orange}{X}) + \varepsilon, \quad X \in \mathbb{R}^d, \quad Y \in \mathbb{R}.$$

- ▷ Optimization problem:

$$\min_{f: (\mathbb{R} \cup \{\text{NA}\})^d \mapsto \mathbb{R}} \mathcal{R}(f) := \mathbb{E} \left[(Y - f(\tilde{X}))^2 \right]$$

- ▷ A Bayes predictor is a minimizer of the risk. It is given by:

$$\tilde{f}^*(\tilde{X}) := \mathbb{E} [Y | X_{obs(M)}, M] = \mathbb{E} [f(X) | X_{obs(M)}, M]$$

where $M \in \{0, 1\}^d$ is the missingness indicator.

- ▷ The Bayes rate \mathcal{R}^* is the risk of the Bayes predictor: $\mathcal{R}^* = \mathcal{R}(\tilde{f}^*)$.
- ▷ A Bayes optimal function f achieves the Bayes rate, i.e., $\mathcal{R}(f) = \mathcal{R}^*$.

Supervised learning with missing values

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

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

Supervised learning with missing values

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

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

Finding the Bayes predictor.

$$f^* \in \operatorname{argmin}_{f: \tilde{\mathbb{R}}^d \rightarrow \mathbb{R}} \mathbb{E} \left[(Y - f(\tilde{X}))^2 \right].$$

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

\Rightarrow One model per pattern (2^d) (Rubin, 1984, generalized propensity score)

Bayes predictor.

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

- ▷ Difficulty due to the half nature of the input space
- ▷ Worst case: 2^d models to learn

Two common strategies:

- ▷ **Impute-then-regress strategies** - impute the data then learn on the imputed data set
 - ◊ Computationally efficient but possibly inconsistent
- ▷ **Pattern-by-pattern strategies** - use a different predictor for each missing pattern
 - ◊ Consistent by design but intractable in most situations

1. Missing values mechanism
2. Single Imputation
3. Multiple Imputation
4. Imputation quality
5. Supervised Learning with Missing values
 - Decision trees as PbP predictors
 - Impute-then-regress procedures with consistent predictors
6. Linear models
 - Linear regression: A pattern-by-pattern approach
 - Linear regression: Impute-then-regress procedures via zero-imputation
 - Classification with missing values
7. Conclusion

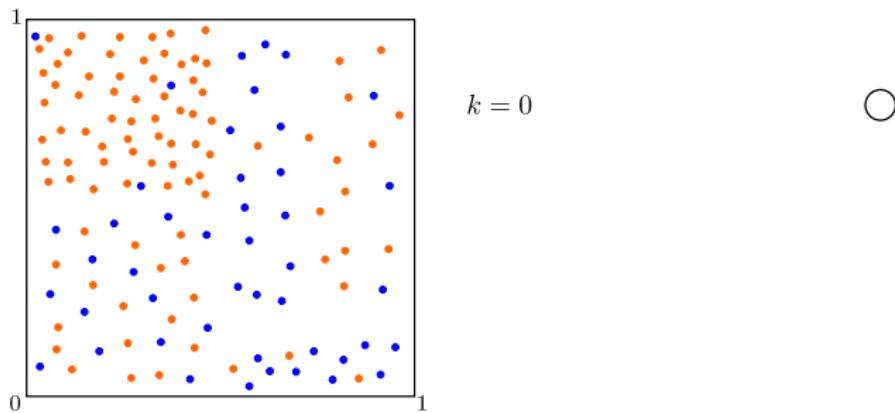
CART (Classification And Regression Tree, 1984)

50 / 99

Built by recursively splitting cells until some stopping criterion is satisfied.

Find the feature j^* , the threshold z^* which minimises the loss

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



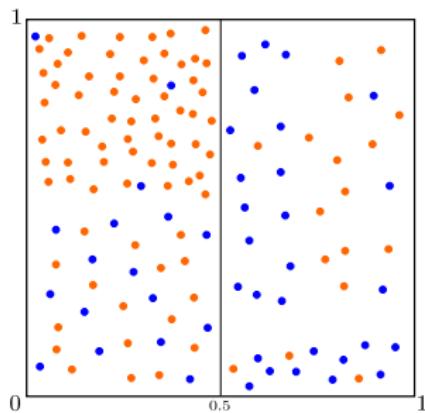
CART (Classification And Regression Tree, 1984)

50 / 99

Built by recursively splitting cells until some stopping criterion is satisfied.

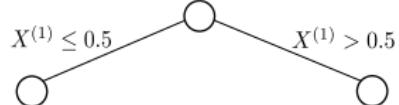
Find the feature j^* , the threshold z^* which minimises the loss

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



$k = 0$

$k = 1$



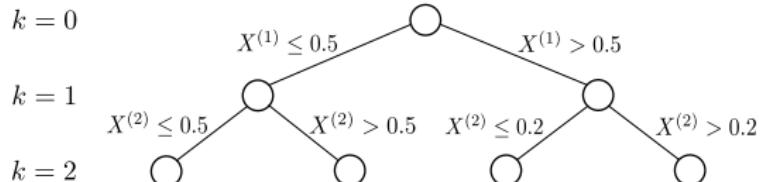
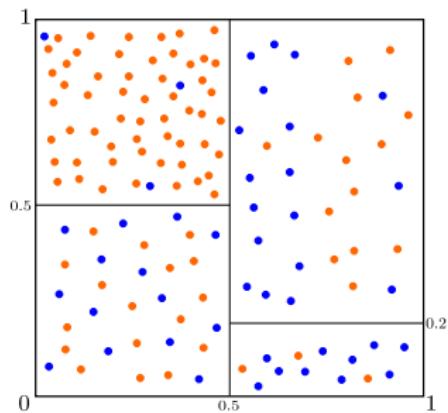
CART (Classification And Regression Tree, 1984)

50 / 99

Built by recursively splitting cells until some stopping criterion is satisfied.

Find the feature j^* , the threshold z^* which minimises the loss

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



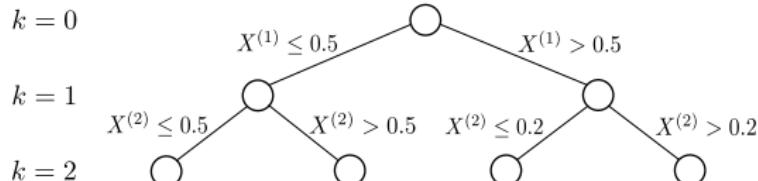
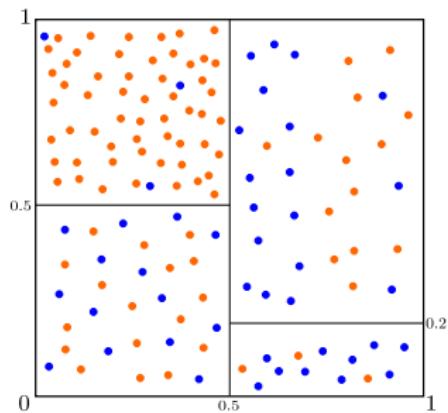
CART (Classification And Regression Tree, 1984)

50 / 99

Built by recursively splitting cells until some stopping criterion is satisfied.

Find the feature j^* , the threshold z^* which minimises the loss

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



Two difficulties with missing data

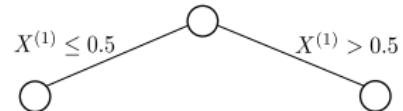
- ▷ How to find the best split?
- ▷ How to propagate missing data down the tree?

CART with missing values

	X_1	X_2	Y
1			
2	NA		
3	NA		
4			

$$k = 0$$

$$k = 1$$

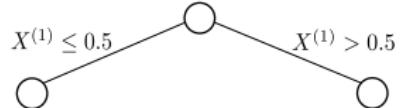


CART with missing values

	X_1	X_2	Y
1			
2	NA		
3	NA		
4			

$$k = 0$$

$$k = 1$$



Two steps:

1. For each variable, compute the splitting criterion on observed values only (e.g., 1 & 4 for X_1)

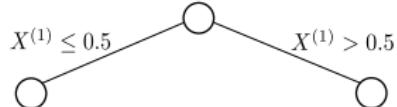
$$\mathbb{E} \left[(Y - \mathbb{E}[Y|X_j \leq z, M_j = 0])^2 \cdot \mathbb{1}_{X_j \leq z, M_j = 0} + (Y - \mathbb{E}[Y|X_j > z, M_j = 0])^2 \cdot \mathbb{1}_{X_j > z, M_j = 0} \right].$$

CART with missing values

	X_1	X_2	Y
1			
2	NA		
3	NA		
4			

$$k = 0$$

$$k = 1$$



Two steps:

1. For each variable, compute the splitting criterion on observed values only (e.g., 1 & 4 for X_1)

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

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

- ▷ Probabilistic split: $Bernoulli(\#L/(\#L + \#R))$ (C4.5)
- ▷ Block: Send all to a side by minimizing the error (lightgbm)
- ▷ Surrogate split: Search another variable that gives a close partition (rpart)

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

1. $\{\tilde{X}_j \leq z \text{ or } \tilde{X}_j = \text{NA}\}$ vs $\{\tilde{X}_j > z\}$
2. $\{\tilde{X}_j \leq z\}$ vs $\{\tilde{X}_j > z \text{ or } \tilde{X}_j = \text{NA}\}$
3. $\{\tilde{X}_j \neq \text{NA}\}$ vs $\{\tilde{X}_j = \text{NA}\}$.

- ▷ The splitting location z depends on the missing values
- ▷ **Missing values treated like a category** (well to handle $\mathbb{R} \cup \text{NA}$)
- ▷ Good for informative pattern, target one model per pattern:

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

- ▷ Implementations **grf/partykit package, XGBoost**
- ▷ Extremely **good performances** in practice **for any mechanism**

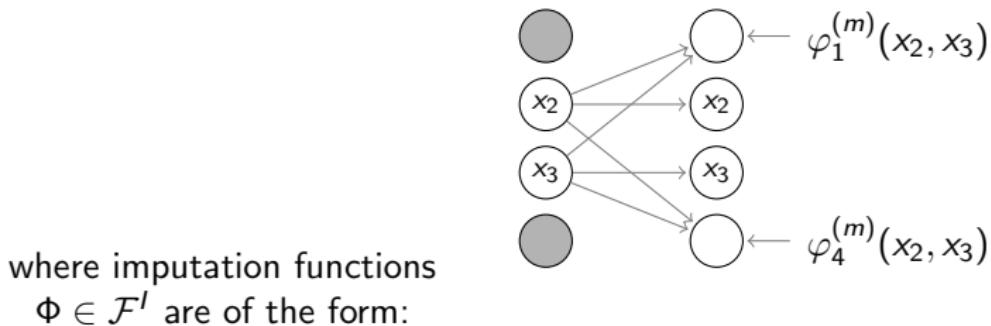
⁴³Twala et al. (2008). Methods for coping with missing data in decision trees. *Pattern Recog.*

1. Missing values mechanism
2. Single Imputation
3. Multiple Imputation
4. Imputation quality
5. Supervised Learning with Missing values
 - Decision trees as PbP predictors
 - Impute-then-regress procedures with consistent predictors
6. Linear models
 - Linear regression: A pattern-by-pattern approach
 - Linear regression: Impute-then-regress procedures via zero-imputation
 - Classification with missing values
7. Conclusion

- ▷ Impute-then-Regress procedures consist in
 1. Impute missing values
 2. train a supervised learning algorithm on the imputed data set.

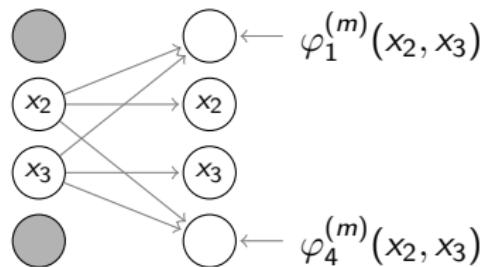
- ▷ Impute-then-Regress procedures consist in
 1. Impute missing values
 2. train a supervised learning algorithm on the imputed data set.
- ▷ More formally, define **Impute-then-Regress procedures** as functions of the form:

$$g \circ \Phi, \text{ where } \Phi \in \mathcal{F}^I, g : \mathbb{R}^d \mapsto \mathbb{R}.$$



- ▷ Impute-then-Regress procedures consist in
 1. Impute missing values
 2. train a supervised learning algorithm on the imputed data set.
- ▷ More formally, define **Impute-then-Regress procedures** as functions of the form:

$$g \circ \Phi, \text{ where } \Phi \in \mathcal{F}^I, g : \mathbb{R}^d \mapsto \mathbb{R}.$$



where imputation functions

$\Phi \in \mathcal{F}^I$ are of the form:

Can Impute-then-Regress procedures be Bayes optimal?

Impute-then-Regress procedures are Bayes optimal

55 / 99

Given an imputation function Φ , we define g_Φ^* the minimizer of the population risk on imputed data as

$$g_\Phi^* \in \operatorname{argmin}_{g: \mathbb{R}^d \mapsto \mathbb{R}} \mathbb{E} \left[\left(Y - g \circ \Phi(\tilde{X}) \right)^2 \right].$$

Given an imputation function Φ , we define g_Φ^* the minimizer of the population risk on imputed data as

$$g_\Phi^* \in \operatorname{argmin}_{g: \mathbb{R}^d \mapsto \mathbb{R}} \mathbb{E} \left[\left(Y - g \circ \Phi(\tilde{X}) \right)^2 \right].$$

Theorem (Le Morvan et al., 2021)

Assume that X admits a density, the response Y is generated as $Y = f^*(X) + \varepsilon$ and $\Phi \in \mathcal{F}'_\infty$ (C^∞ imputation functions). Then,

- for all missing data mechanisms,
- and for almost all imputation functions,

$g_\Phi^* \circ \Phi$ is Bayes optimal.

Given an imputation function Φ , we define g_Φ^* the minimizer of the population risk on imputed data as

$$g_\Phi^* \in \operatorname{argmin}_{g: \mathbb{R}^d \mapsto \mathbb{R}} \mathbb{E} \left[\left(Y - g \circ \Phi(\tilde{X}) \right)^2 \right].$$

Theorem (Le Morvan et al., 2021)

Assume that X admits a density, the response Y is generated as $Y = f^*(X) + \varepsilon$ and $\Phi \in \mathcal{F}'_\infty$ (C^∞ imputation functions). Then,

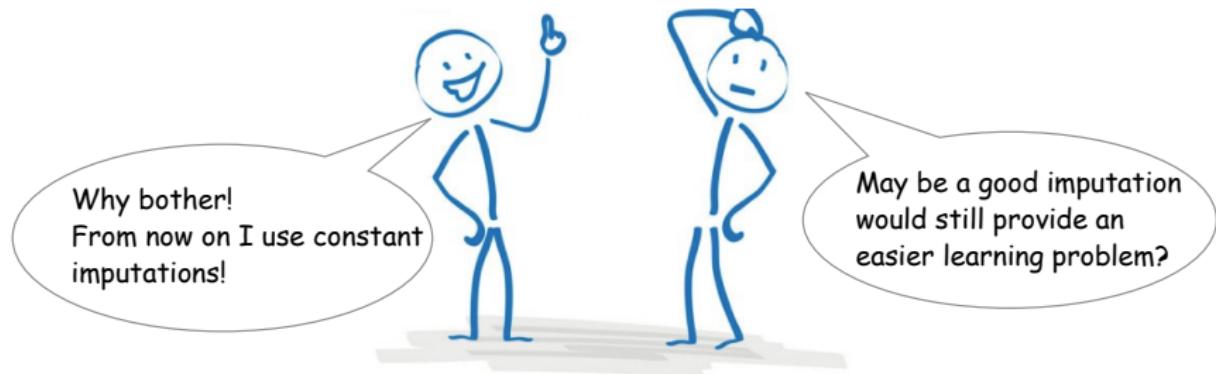
- for all missing data mechanisms,
- and for almost all imputation functions,

$g_\Phi^* \circ \Phi$ is Bayes optimal.

For almost all imputation functions, and all missing data mechanisms, a universally consistent algorithm trained on the imputed data is a consistent procedure.

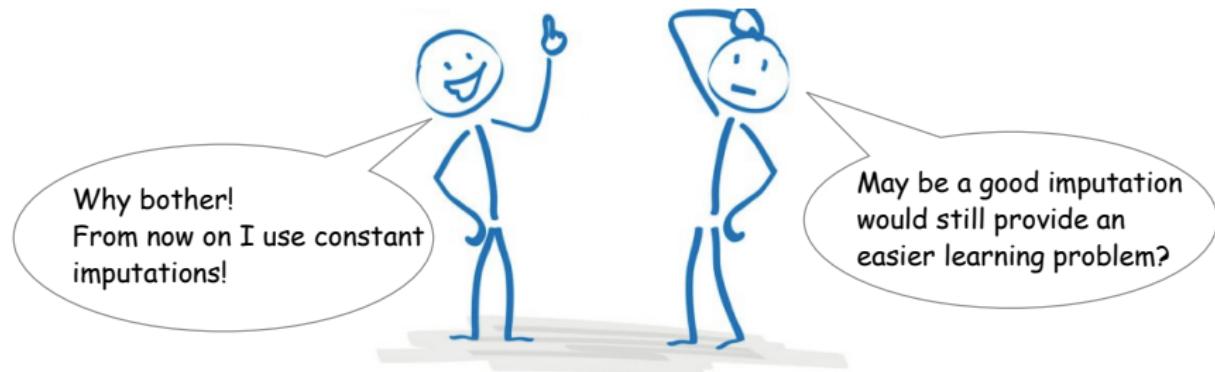
Which imputation function should one choose?

56 / 99



Which imputation function should one choose?

56 / 99



Question

Are there *continuous* Impute-then-Regress decompositions of Bayes predictors?

From now on, we suppose f^* (Byes predictor with complete data) is smooth and consider the conditional expectation Φ^{CI} .

Question

What can we say about the optimal predictor on the conditionally imputed data: $g_{\Phi^{CI}}^ \circ \Phi^{CI}$?*

Question

What can we say about the optimal predictor on the conditionally imputed data: $g_{\Phi^{CI}}^ \circ \Phi^{CI}$?*

Theorem (Le Morvan et al., 2021)

Suppose that $f^ \circ \Phi^{CI}$ is not Bayes optimal, and that the probability of observing all variables is strictly positive, i.e., $P(M = \mathbf{0}, X = x) > 0$, for all x . Then there is no continuous function g such that $g \circ \Phi^{CI}$ is Bayes optimal.*

- ▷ In the above setting, $g_{\Phi^{CI}}^*$ is not continuous. Thus, imputing via conditional expectation leads to a difficult learning problem.
- ▷ Almost all imputations lead to consistent estimators but some ease the training of the supervised learning algorithm.

Imputation-then-regress: does imputation matter?

58 / 99

Adding the mask to the input (one mask per feature):

$$\begin{array}{cc} X_1 & X_2 \\ \left(\begin{array}{cc} 1 & 2 \\ 3 & \text{NA} \\ \text{NA} & 4 \end{array} \right) & \rightarrow \end{array} \begin{array}{cccc} X_1 & X_2 & M_1 & M_2 \\ \left(\begin{array}{cccc} 1 & 2 & 0 & 0 \\ 3 & \text{NA} & 0 & 1 \\ \text{NA} & 4 & 1 & 0 \end{array} \right) \end{array}$$

From an empirical study over 19 datasets⁴⁴:

⁴⁴M. Le Morvan, G. Varoquaux, Imp. for pred.: beware of diminish. returns. (ICLR2025)

⁴⁵Mike et al. (2023). The Missing Indicator Method: From Low to High Dimensions. SIGKDD.

Imputation-then-regress: does imputation matter?

58 / 99

Adding the mask to the input (one mask per feature):

$$\begin{pmatrix} X_1 & X_2 \\ 1 & 2 \\ 3 & \text{NA} \\ \text{NA} & 4 \end{pmatrix} \rightarrow \begin{pmatrix} X_1 & X_2 & M_1 & M_2 \\ 1 & 2 & 0 & 0 \\ 3 & \text{NA} & 0 & 1 \\ \text{NA} & 4 & 1 & 0 \end{pmatrix}$$

From an empirical study over 19 datasets⁴⁴:

- ▷ Imputation accuracy matters less when using expressive models or when incorporating the mask as complementary inputs⁴⁵
- ▷ Imputation accuracy matters much more for generated linear outcomes than for real-data outcome
- ▷ Adding the mask as input is beneficial for prediction performances even for MCAR settings, where missingness is uninformative.

⁴⁴M. Le Morvan, G. Varoquaux, Imp. for pred.: beware of diminish. returns. (ICLR2025)

⁴⁵Mike et al. (2023). The Missing Indicator Method: From Low to High Dimensions. SIGKDD.

Adding the mask to the input (one mask per feature):

$$\begin{pmatrix} X_1 & X_2 \\ 1 & 2 \\ 3 & \text{NA} \\ \text{NA} & 4 \end{pmatrix} \rightarrow \begin{pmatrix} X_1 & X_2 & M_1 & M_2 \\ 1 & 2 & 0 & 0 \\ 3 & \text{NA} & 0 & 1 \\ \text{NA} & 4 & 1 & 0 \end{pmatrix}$$

From an empirical study over 19 datasets⁴⁴:

- ▷ Imputation accuracy matters less when using expressive models or when incorporating the mask as complementary inputs⁴⁵
- ▷ Imputation accuracy matters much more for generated linear outcomes than for real-data outcome
- ▷ Adding the mask as input is beneficial for prediction performances even for MCAR settings, where missingness is uninformative.

Investing in more flexible models is more efficient than investing in more complex imputations.

⁴⁴M. Le Morvan, G. Varoquaux, Imp. for pred.: beware of diminish. returns. (ICLR2025)

⁴⁵Mike et al. (2023). The Missing Indicator Method: From Low to High Dimensions. SIGKDD.

Bayes predictor

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

Two common strategies:

- ▷ Impute-then-regress strategies - impute the data then learn on the imputed data set
 - ◊ Computationally efficient but possibly inconsistent
 - ◊ Consistent if used with a non-parametric learning algorithm (forests, tree boosting, nearest neighbor...)
- ▷ Pattern-by-pattern strategies - use a different predictor for each missing pattern
 - ◊ Consistent by design but intractable in most situations

1. Missing values mechanism
2. Single Imputation
3. Multiple Imputation
4. Imputation quality
5. Supervised Learning with Missing values
 - Decision trees as PbP predictors
 - Impute-then-regress procedures with consistent predictors
6. Linear models
 - Linear regression: A pattern-by-pattern approach
 - Linear regression: Impute-then-regress procedures via zero-imputation
 - Classification with missing values
7. Conclusion

Our aim

Predict on new data, which may contain missing entries.

MCAR

(missing completely at random)

$$\mathbb{P}(M|X) = \mathbb{P}(M)$$

MAR (missing at random)

$$\mathbb{P}(M|X) = \mathbb{P}(M|X^{(obs)})$$

MNAR (missing not at random)

Linear model

$$Y = X^T \beta^* + \text{noise}$$

- ▷ $Y \in \mathbb{R}$ (regression) outcome is always observed
- ▷ $X \in \mathbb{R}^d$ contains missing values!
- ▷ β^* model parameter

Linear models do not remain linear

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.

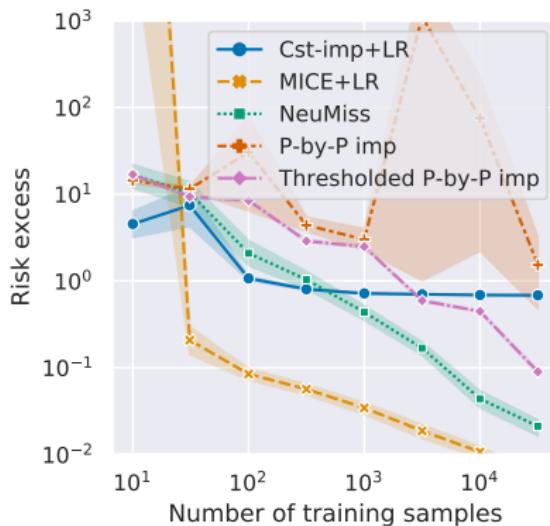
Here, the submodel for which only X_1 is observed is not linear.

- ⇒ There exists a large variety of submodels for a same linear model.
- ⇒ Submodel natures depend on the structure of X and on the missing-value mechanism.

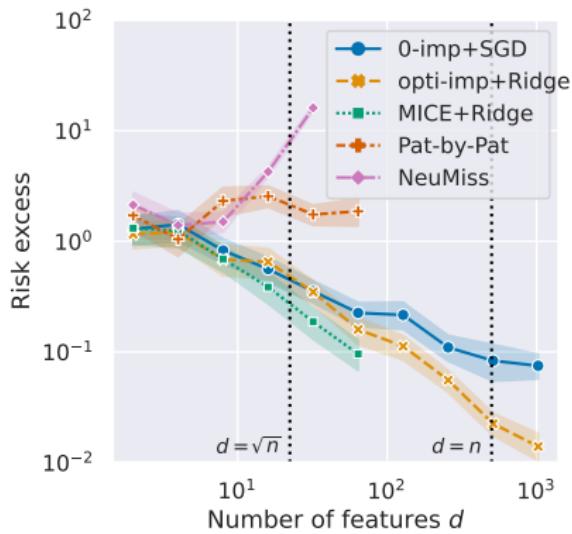
Handling missing values in linear models for prediction^{63 / 99}

2 possible approaches

- ▷ Pattern-by-pattern methods
- ▷ Impute-then-regress procedures



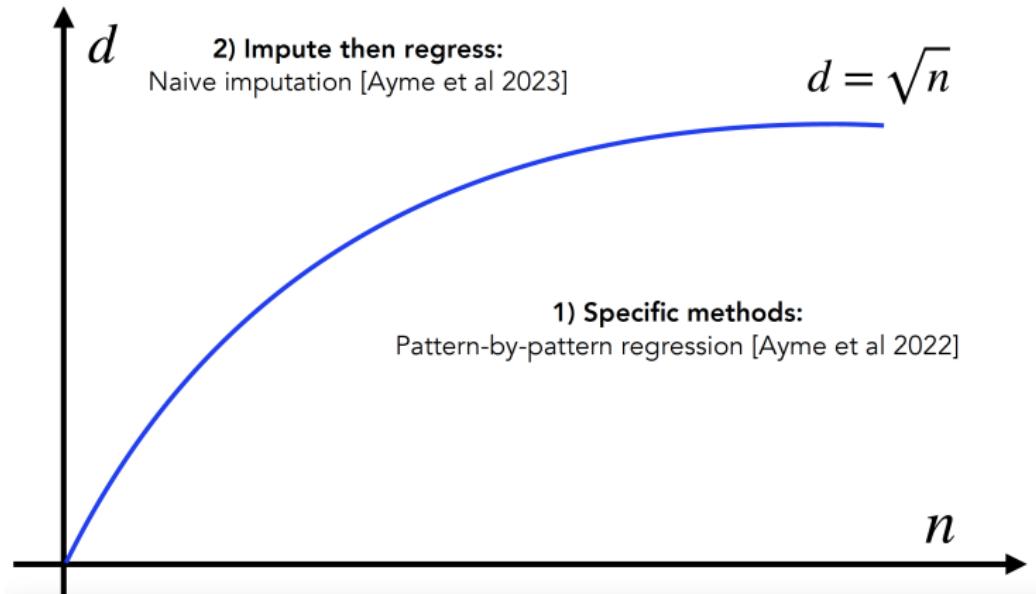
Fixed dimension



Fixed sample size

Different strategies for prediction

64 / 99



1. Missing values mechanism
2. Single Imputation
3. Multiple Imputation
4. Imputation quality
5. Supervised Learning with Missing values
 - Decision trees as PbP predictors
 - Impute-then-regress procedures with consistent predictors
6. Linear models
 - Linear regression: A pattern-by-pattern approach
 - Linear regression: Impute-then-regress procedures via zero-imputation
 - Classification with missing values
7. Conclusion

- ▷ Dataset $\mathcal{D}_n = \{(Z_i, Y_i), i \in [n]\}$ where

$$Z_i = (X_{\text{obs}(M_i)}, M_i).$$

- ▷ New test point $Z = (X_{\text{obs}(M)}, M)$ (with **unknown** target Y).

Goal in prediction

Find a **linear function** \hat{f} that minimizes the risk

$$R_{\text{miss}}(\hat{f}) = \mathbb{E} \left[(\hat{f}(Z) - Y)^2 \right].$$

Pattern-by-pattern Bayes predictor

Consider either

▷ $X \sim \mathcal{N}(\mu, \Sigma)$ Gaussian (G)

or,

▷ $X|(M = m) \sim \mathcal{N}(\mu^m, \Sigma^m)$ Gaussian pattern mixture model (GPMM)

Decompose the Bayes predictor

$$f^*(Z) = \sum_{m \in \mathcal{M}} f_m^*(X_{obs(m)}) \mathbb{1}_{M=m},$$

with f_m^* the Bayes predictor conditionally on the event $(M = m)$.

Proposition

[Le Morvan et al 2020]

If [(MCAR or MAR) and G] or GPMM then, for all $m \in \mathcal{M}$,

f_m^* is linear.

A missing-distribution-free upper bound

Predictor $\hat{f}(Z) = \sum_{m \in \mathcal{M}} \hat{f}_m(X_{obs(m)}) \mathbb{1}_{M=m}$ (pattern-by-pattern OLS)
 where \hat{f}_m is a modified least-square regression rule trained on

$$\mathcal{D}_m = \{(X_{i,obs(m)}, Y_i), M_i = m\}.$$

Theorem (simplified) [Le Morvan et al. 2020] [Ayme, Boyer, Dieuleveut, S. 2022]

If [(MCAR or MAR) and G] or GPMM then

$$\mathbb{E} \left[(\hat{f}(Z) - f^*(Z))^2 \right] \lesssim \log(n) 2^d \frac{d}{n}$$

where the constant depends on the level of noise.

A missing-distribution-free upper bound

Predictor $\hat{f}(Z) = \sum_{m \in \mathcal{M}} \hat{f}_m(X_{obs(m)}) \mathbb{1}_{M=m}$ (pattern-by-pattern OLS)
 where \hat{f}_m is a modified least-square regression rule trained on

$$\mathcal{D}_m = \{(X_{i,obs(m)}, Y_i), M_i = m\}.$$

Theorem (simplified) [Le Morvan et al. 2020] [Ayme, Boyer, Dieuleveut, S. 2022]

If [(MCAR or MAR) and G] or GPMM then

$$\mathbb{E} \left[(\hat{f}(Z) - f^*(Z))^2 \right] \lesssim \log(n) 2^d \frac{d}{n}$$

where the constant depends on the level of noise.

- ▷ This result does not depend on the distribution of missing patterns.
- ▷ Number of parameters is $p := d2^d$. This result suffers from the curse of dimensionality even with small d .

A missing pattern distribution adaptive bound

Idea: Regression only on **high frequency** missing patterns

$$\hat{f}(Z) = \sum_{m \in \mathcal{M}} \hat{f}_m(X_{obs(m)}) \mathbb{1}_{M=m} \mathbb{1}_{|\mathcal{D}_m| \geq d}.$$

A missing pattern distribution adaptive bound

69 / 99

Idea: Regression only on **high frequency** missing patterns

$$\hat{f}(Z) = \sum_{m \in \mathcal{M}} \hat{f}_m(X_{obs(m)}) \mathbb{1}_{M=m} \mathbb{1}_{|\mathcal{D}_m| \geq d}.$$

Theorem [Ayme, Boyer, Dieuleveut, S. 2022]

$$\mathbb{E} \left[\left(\hat{f}(Z) - f^*(Z) \right)^2 \right] \lesssim \log(n) \mathcal{E}_p(d/n),$$

with $\mathcal{E}_p(d/n) := \sum_m \min(p_m, d/n)$.

- ▷ Valid for MCAR, MAR and MNAR settings.
- ▷ Adaptive to missing data distribution via $\mathcal{E}_p(d/n) \leq \text{Card}(\mathcal{M})(d/n)$.

Examples

1. Uniform distribution: $\mathcal{E}_p\left(\frac{d}{n}\right) = 2^d d/n$
2. Bernoulli distribution: $M_j \sim \mathcal{B}(\varepsilon)$ with $\varepsilon \leq d/n$: $\mathcal{E}_p\left(\frac{d}{n}\right) = d^2/n$

A lower bound

70 / 99

Let \mathcal{P}_p be a class of data distributions $\left\{ \begin{array}{l} X|(M=m) \sim \mathcal{N}(\mu^m, \Sigma^m) \\ \text{Linear model} \\ \mathbb{P}[M=m] = p_m \end{array} \right.$

$$\text{Minimax error } (p) = \underbrace{\min_{\tilde{f}}}_{\text{Best algo}} \quad \underbrace{\max_{\mathbb{P} \in \mathcal{P}_p}}_{\substack{\text{Worst case on a class} \\ \mathcal{P}_p \text{ of problems}}} \quad \mathbb{E}_{\mathbb{P}} \left[(\tilde{f}(Z) - f^*(Z))^2 \right]$$

A lower bound

Let \mathcal{P}_p be a class of data distributions $\left\{ \begin{array}{l} X|(M=m) \sim \mathcal{N}(\mu^m, \Sigma^m) \\ \text{Linear model} \\ \mathbb{P}[M=m] = p_m \end{array} \right.$

$$\underset{\text{Minimax error}}{\text{(p)}} = \underbrace{\min_{\tilde{f}}}^{\text{Best algo}} \underbrace{\max_{\substack{\mathbb{P} \in \mathcal{P}_p \\ \mathcal{P}_p \text{ of problems}}} \mathbb{E}_{\mathbb{P}} \left[(\tilde{f}(Z) - f^*(Z))^2 \right]}$$

Theorem

[Ayme, Boyer, Dieuleveut, S. 2022]

$$\sigma^2 \mathcal{E}_p \left(\frac{1}{n} \right) \lesssim \underset{\text{Minimax error}}{\text{(p)}} \leq \mathbb{E} \left[\left(\hat{f}(Z) - f^*(Z) \right)^2 \right] \lesssim \log(n) \mathcal{E}_p \left(\frac{d}{n} \right)$$

A lower bound

Let \mathcal{P}_p be a class of data distributions $\left\{ \begin{array}{l} X|(M=m) \sim \mathcal{N}(\mu^m, \Sigma^m) \\ \text{Linear model} \\ \mathbb{P}[M=m] = p_m \end{array} \right.$

$$\text{Minimax error } (p) = \underbrace{\min_{\tilde{f}}}_{\text{Best algo}} \quad \underbrace{\max_{\substack{\mathbb{P} \in \mathcal{P}_p \\ \mathcal{P}_p \text{ of problems}}} \mathbb{E}_{\mathbb{P}} \left[(\tilde{f}(Z) - f^*(Z))^2 \right]}_{\text{Worst case on a class}}$$

Theorem

[Ayme, Boyer, Dieuleveut, S. 2022]

$$\sigma^2 \mathcal{E}_p \left(\frac{1}{n} \right) \lesssim \text{Minimax error } (p) \leq \mathbb{E} \left[\left(\hat{f}(Z) - f^*(Z) \right)^2 \right] \lesssim \log(n) \mathcal{E}_p \left(\frac{d}{n} \right)$$

Examples

- ▷ Uniform distribution
- ▷ Bernoulli distribution $M_j \sim \mathcal{B}(\varepsilon)$
with $\varepsilon \leq d/n$

$$\begin{aligned} \mathcal{E}_p \left(\frac{1}{n} \right) &= 2^d / n & \mathcal{E}_p \left(\frac{d}{n} \right) &= 2^d d / n \\ \mathcal{E}_p \left(\frac{1}{n} \right) &= d / n & \mathcal{E}_p \left(\frac{d}{n} \right) &= d^2 / n \end{aligned}$$

- ☞ For data regimes where n is large, several problems can be learned, even for MNAR.
- ☞ The procedure can be modified to adapt to the distribution of missing patterns.
- ☞ **The dimension is an issue**, even under the classical assumptions (MAR)

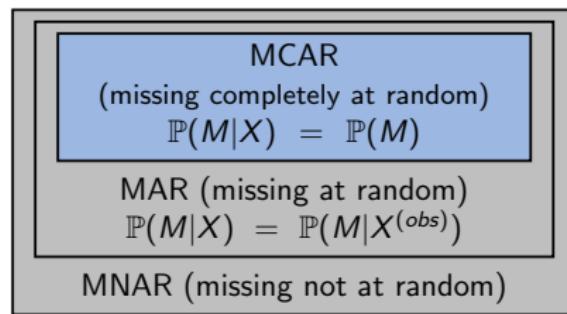
1. Missing values mechanism
2. Single Imputation
3. Multiple Imputation
4. Imputation quality
5. Supervised Learning with Missing values
 - Decision trees as PbP predictors
 - Impute-then-regress procedures with consistent predictors
6. Linear models
 - Linear regression: A pattern-by-pattern approach
 - Linear regression: Impute-then-regress procedures via zero-imputation**
 - Classification with missing values
7. Conclusion

- ▷ Impute-then-regress method
 - 1. Impute the missing values by 0 to get X_{imp} (e.g., via `df.fillna(0)`)
 - 2. Perform a SGD regression

Impute-then-regress?

73 / 99

- ▷ Impute-then-regress method
 1. Impute the missing values by 0 to get X_{imp} (e.g., via `df.fillna(0)`)
 2. Perform a SGD regression
- ▷ Focus on MCAR values: $M_1, \dots, M_d \sim \mathcal{B}(\rho)$
 ρ = probability to be observed



impute by 0 = doesn't exploit observed values?

Risk decomposition

- ▷ R^* = optimal risk without missing data
- ▷ R_{miss}^* = optimal risk with missing data

$$\Delta_{\text{miss}} := R_{\text{miss}}^* - R^* \quad (\text{missing data error})$$

- ▷ $R_{\text{imp}}(\theta) =$ the risk of $f_\theta(X_{\text{obs}}, M) = \theta^\top X_{\text{imp}}$
- ▷ $R_{\text{imp}}(\theta_{\text{imp}}^*) =$ optimal risk of linear prediction after imputation by 0

$$\Delta_{\text{imp}/\text{miss}} := R_{\text{imp}}(\theta_{\text{imp}}^*) - R_{\text{miss}}^* \quad (\text{imputation error})$$

- ▷ Risk decomposition:

$$R_{\text{miss}}(f_\theta) = R^* + \underbrace{\Delta_{\text{miss}} + \Delta_{\text{imp}/\text{miss}}}_{\text{missing data and imputation error}} + \underbrace{R_{\text{miss}}(f_\theta) - R_{\text{imp}}(\theta_{\text{imp}}^*)}_{\text{estimation/optimization error}}$$

Toy example: how imputed inputs disturb learning

75 / 99

- ▷ Complete model
 - ◊ $Y = X_1$
 - ◊ $X = (X_1, \dots, X_1)$
 - ◊ $R^* = 0$
 - ◊ $M_1, \dots, M_d \sim \mathcal{B}(1/2)$

Toy example: how imputed inputs disturb learning

75 / 99

- ▷ Complete model
 - ◊ $Y = X_1$
 - ◊ $X = (X_1, \dots, X_1)$
 - ◊ $R^* = 0$
 - ◊ $M_1, \dots, M_d \sim \mathcal{B}(1/2)$
- ▷ With **imputed** inputs and $\theta_1 = (1, 0, \dots, 0)^\top$
 - ◊ $X_{\text{imp}}^\top \theta_1 = X_1 M_1$
 - ◊ $R_{\text{imp}}(\theta_1) = \frac{1}{2} \mathbb{E}[Y^2]$
- ▷ With **imputed** inputs and $\theta_2 = 2(1/d, 1/d, \dots, 1/d)^\top$
 - ◊ $X_{\text{imp}}^\top \theta_2 = \frac{2}{d} X_1 \sum_j M_j$
 - ◊ $R_{\text{imp}}(\theta_2) = \frac{1}{d} \mathbb{E}[X_1^2]$
 - ◊ $\Delta_{\text{miss}} + \Delta_{\text{imp/miss}} \leq R_{\text{imp}}(\theta_2) - R^* \leq \frac{1}{d} \mathbb{E}[Y^2]$

Toy example: how imputed inputs disturb learning

75 / 99

- ▷ Complete model
 - ◊ $Y = X_1$
 - ◊ $X = (X_1, \dots, X_1)$
 - ◊ $R^* = 0$
 - ◊ $M_1, \dots, M_d \sim \mathcal{B}(1/2)$
- ▷ With **imputed** inputs and $\theta_1 = (1, 0, \dots, 0)^\top$
 - ◊ $X_{\text{imp}}^\top \theta_1 = X_1 M_1$
 - ◊ $R_{\text{imp}}(\theta_1) = \frac{1}{2} \mathbb{E}[Y^2]$
- ▷ With **imputed** inputs and $\theta_2 = 2(1/d, 1/d, \dots, 1/d)^\top$
 - ◊ $X_{\text{imp}}^\top \theta_2 = \frac{2}{d} X_1 \sum_j M_j$
 - ◊ $R_{\text{imp}}(\theta_2) = \frac{1}{d} \mathbb{E}[X_1^2]$
 - ◊ $\Delta_{\text{miss}} + \Delta_{\text{imp/miss}} \leq R_{\text{imp}}(\theta_2) - R^* \leq \frac{1}{d} \mathbb{E}[Y^2]$

correlation \Rightarrow low imputation/missing values error ?

- ▷ Ridge-regularized risk with complete data

$$R_\lambda(\theta) = R(\theta) + \lambda \|\theta\|_2^2$$

- ▷ Standard in high-dimension settings

Theorem

[Ayme, Boyer, Dieuleveut, S. 2023]

Under the MCAR Bernoulli model of probability ρ of observation and $\text{Var}(X_j) = 1 \forall j$,

$$R_{\text{imp}}(\theta) = R(\rho\theta) + \rho(1 - \rho)\|\theta\|_2^2$$

Consequences

1. $\Delta_{\text{miss}} + \Delta_{\text{imp}/\text{miss}} = \text{ridge bias for } \lambda = \frac{1-\rho}{\rho}$
2. θ_{imp}^* on a small ball around 0 (implicit regularization)

- ☞ Imputed MCAR missing values seem to be at the same price of ridge regularization

Learning with low-rank and imputed-by-0 data

77 / 99

- ▷ **Low-rank data:** covariance matrix $\Sigma = [XX^\top]$ is

$$\Sigma = \sum_{j=1}^r \lambda_j v_j v_j^\top,$$

with $\lambda_1 = \dots = \lambda_r$ and $r \ll d$.

- ▷ Bias on low-rank data:

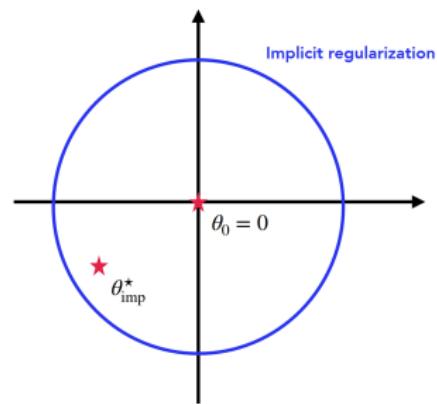
$$\Delta_{\text{miss}} + \Delta_{\text{imp/miss}} \lesssim \frac{1 - \rho}{\rho} \frac{r}{d} \mathbb{E}[Y^2]$$

correlation \Rightarrow low imputation/missing values error !

- ▷ Averaged SGD iterates:

$$\begin{cases} \theta_{\text{imp},t} &= [I - \gamma X_{\text{imp},t} X_{\text{imp},t}^\top] \theta_{\text{imp},t-1} + \gamma Y_t X_{\text{imp},t} \\ \bar{\theta}_{\text{imp},n} &= \frac{1}{n+1} \sum_{t=1}^n \theta_{\text{imp},t} \end{cases}$$

- ▷ Why use SGD ?
 1. Streaming online (one pass only)
 2. Minimizes directly the generalization risk R
 3. Friendly assumptions
 4. Leverage the implicit regularization of naive imputations choosing $\theta_{\text{imp},0} = 0$ and $\gamma = 1/d\sqrt{n}$.



Theorem

[Ayme, Boyer, Dieuleveut, S. 2023]

Under classical assumptions for SGD,

$$\mathbb{E} [R_{\text{imp}}(\bar{\theta}_{\text{imp},n})] - R^* \leq \Delta_{\text{miss}} + \Delta_{\text{imp}/\text{miss}} + \frac{d}{\sqrt{n}} \|\theta_{\text{imp}}^*\|_2^2 + \frac{\text{noise variance}}{\sqrt{n}}$$

Theorem

[Ayme, Boyer, Dieuleveut, S. 2023]

Under classical assumptions for SGD,

$$\mathbb{E} [R_{\text{imp}}(\bar{\theta}_{\text{imp},n})] - R^* \leq \Delta_{\text{miss}} + \Delta_{\text{imp}/\text{miss}} + \frac{d}{\sqrt{n}} \|\theta_{\text{imp}}^*\|_2^2 + \frac{\text{noise variance}}{\sqrt{n}}$$

▷ Example: low-rank setting

$$\mathbb{E} [R_{\text{imp}}(\bar{\theta}_{\text{imp},n})] - R^* \lesssim \left(\frac{1}{\rho\sqrt{n}} + \frac{1-\rho}{d} \right) \frac{r}{d} \mathbb{E} Y^2 + \frac{\text{noise variance}}{\sqrt{n}}$$

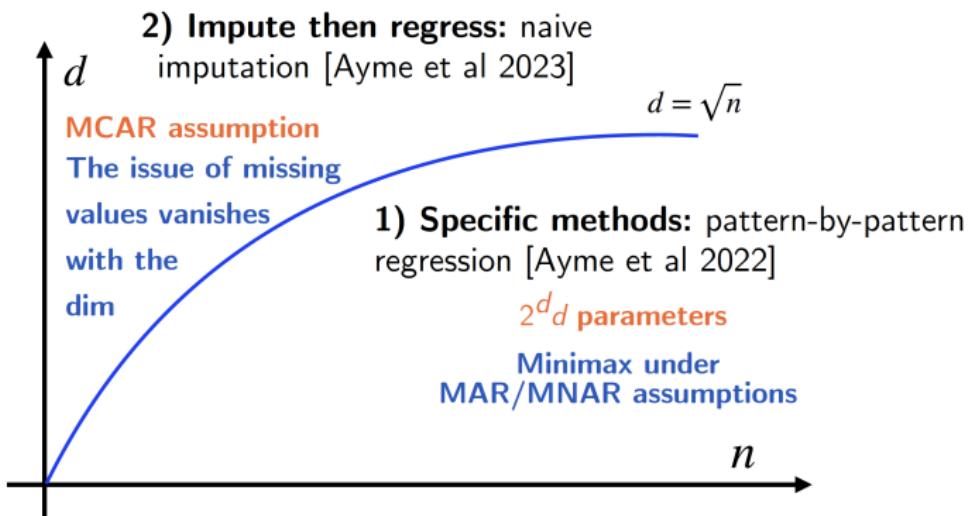
► Imputation bias vanishes for $d \gg \sqrt{n}$

Naive imputation implicitly regularizes HD linear models

- ▷ MCAR inputs
(observation rate = ρ)
- ▷ All in all

Performing standard linear regression on imputed-by-0 data

Adding a ridge regularization w/ parameter
 $\lambda = \frac{1-\text{observation rate}}{\text{observation rate}}$



1. Missing values mechanism
2. Single Imputation
3. Multiple Imputation
4. Imputation quality
5. Supervised Learning with Missing values
 - Decision trees as PbP predictors
 - Impute-then-regress procedures with consistent predictors
6. Linear models
 - Linear regression: A pattern-by-pattern approach
 - Linear regression: Impute-then-regress procedures via zero-imputation
 - Classification with missing values
7. Conclusion

LDA

Let $\mathbb{P}(Y = 1) = 0.5$ and $\forall k \in \{-1, 1\}$, $X|Y = k \sim \mathcal{N}(\mu_k, \Sigma)$.

Bayes predictor for the complete case:

$$h_{\text{comp}}^*(x) := \text{sign} \left((\mu_1 - \mu_{-1})^\top \Sigma^{-1} \left(x - \frac{\mu_1 + \mu_{-1}}{2} \right) \right).$$

⁴⁶A primer on linear classification with missing data A.D.R. Lobo, A. Ayme, C. Boyer, E. Scornet

Bayes predictor for the complete case:

$$h_{\text{comp}}^*(x) := \text{sign} \left((\mu_1 - \mu_{-1})^\top \Sigma^{-1} \left(x - \frac{\mu_1 + \mu_{-1}}{2} \right) \right).$$

Proposition: Bayes predictor for LDA+MCAR

Assume LDA + MCAR. Then the PbP Bayes classifier satisfies

$$\begin{aligned} h_m^*(x_{\text{obs}(m)}) &= \text{sign} \left((\mu_{1,\text{obs}(m)} - \mu_{-1,\text{obs}(m)})^\top \Sigma_{\text{obs}(m)}^{-1} \right. \\ &\quad \times \left. \left(x_{\text{obs}(m)} - \frac{\mu_{1,\text{obs}(m)} + \mu_{-1,\text{obs}(m)}}{2} \right) \right). \end{aligned}$$

⁴⁶A primer on linear classification with missing data A.D.R. Lobo, A. Ayme, C. Boyer, E. Scornet

Bayes predictor for the complete case:

$$h_{\text{comp}}^*(x) := \text{sign} \left((\mu_1 - \mu_{-1})^\top \Sigma^{-1} \left(x - \frac{\mu_1 + \mu_{-1}}{2} \right) \right).$$

Proposition: Bayes predictor for LDA+MCAR

Assume LDA + MCAR. Then the PbP Bayes classifier satisfies

$$\begin{aligned} h_m^*(x_{\text{obs}(m)}) &= \text{sign} \left((\mu_{1,\text{obs}(m)} - \mu_{-1,\text{obs}(m)})^\top \Sigma_{\text{obs}(m)}^{-1} \right. \\ &\quad \times \left. \left(x_{\text{obs}(m)} - \frac{\mu_{1,\text{obs}(m)} + \mu_{-1,\text{obs}(m)}}{2} \right) \right). \end{aligned}$$

- ▷ PbP strategy is Bayes optimal
- ▷ Constant imputation is not optimal (if Σ is nondiagonal)
- ▷ Extension to MNAR scenarios (GPMM)

⁴⁶A primer on linear classification with missing data A.D.R. Lobo, A. Ayme, C. Boyer, E. Scornet

Logistic model

$$\mathbb{P}[Y = 1|X] = \sigma(\beta_0^* + \sum_j \beta_j^* X_j) \text{ with } \sigma(t) = 1/(1 + e^{-t}).$$

Bayes classifier: $g^*(\tilde{x}) = \mathbb{1}_{\eta^*(\tilde{x}) > 0.5}$ with $\eta^*(\tilde{x}) = \mathbb{E}[Y|\tilde{X} = \tilde{x}]$.

III-specified PbP logistic regression

Assume MCAR data in a logistic model for complete data with X_1, \dots, X_d independent Gaussian random variables. Let $m \in \{0, 1\}^d$ and assume that there exists a vector $\beta_m^* \in \mathbb{R}^{|\text{obs}(m)|+1}$ such that

$$\mathbb{P}(Y = 1|X_{\text{obs}(m)}, M = m) = \sigma\left(\beta_{0,m}^* + \sum_{j \in \text{obs}(m)} \beta_{j,m}^* X_j\right).$$

Then, for all $j \in \text{mis}(m)$, $\beta_j^* = 0$

Ill-specified PbP logistic regression

Assume MCAR data in a logistic model for complete data with X_1, \dots, X_d independent Gaussian random variables. Let $m \in \{0, 1\}^d$ and assume that there exists a vector $\beta_m^* \in \mathbb{R}^{|\text{obs}(m)|+1}$ such that

$$\mathbb{P}(Y = 1 | X_{\text{obs}(m)}, M = m) = \sigma\left(\beta_{0,m}^* + \sum_{j \in \text{obs}(m)} \beta_{j,m}^* X_j\right).$$

Then, for all $j \in \text{mis}(m)$, $\beta_j^* = 0$

- ▷ Logistic model cannot hold on complete data AND on data with a given missing pattern
- ▷ Constant imputation Impute-then-Logistic-Regression is ill specified

$$\begin{aligned}\mathbb{E}[Y | X_{\text{obs}(M)}, M = m] &= \mathbb{E}\left[\mathbb{E}[Y | X] | X_{\text{obs}(M)}\right] = \mathbb{E}\left[\sigma(X) | X_{\text{obs}(M)}\right] \\ &\neq \sigma\left(\mathbb{E}[\beta_0^* + \sum_{j=1}^d \beta_j^* X_j | X_{\text{obs}(M)}]\right).\end{aligned}$$

Denote $\Phi(t)$ the probit function: $\Phi(t) = (2\pi)^{-1/2} \int_{-\infty}^t e^{-t^2/2} dt$,

Theorem

Assume a logistic model on complete data and a GPMM:

$X|M = m \sim \mathcal{N}(\mu_m, \Sigma_m)$. Then, for all m , the Bayes probability on pattern m , η_m^* , satisfies for all $x \in \mathbb{R}^{|\text{obs}(m)|}$,

$$\left| \eta_m^*(x) - \sigma \left(\frac{\alpha_{0,m} + \alpha_m^\top x}{\sqrt{1 + (\pi/8)\tilde{\sigma}_m^2}} \right) \right| \leq 2\|\varepsilon\|_\infty \approx 0.036,$$

where $\varepsilon(t) = \Phi(t) - \sigma(t\sqrt{8/\pi})$, and $\alpha_{0,m}, \alpha_m, \tilde{\sigma}_m^2$.

Theoretical ground for understanding why PbP logistic regression performs well in practice while being ill-specified.

See also⁴⁷

⁴⁷K.A. Verchand, A. Montanari, High-dimensional logistic regression with missing data: Imputation, regularization, and universality

⁴⁸C. Muller, E. Scornet, J. Josse, When Pattern-by-Pattern Works: Theoretical and Empirical Insights for Logistic Models with Missing Values

- ▷ Pattern-by-pattern (PbP): Logistic regression on each pattern

- ▷ Pattern-by-pattern (PbP)
- ▷ Mean imputation (Mean.IMP): Mean per covariate

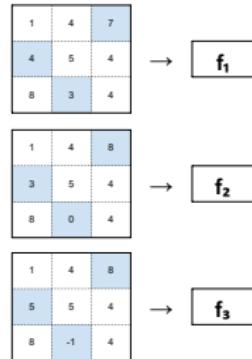
- ▷ Pattern-by-pattern (PbP)
- ▷ Mean imputation (Mean.IMP)
- ▷ Fully specified (SAEM): Fully parametrized model,
assuming normal covariates + logistic regression,
optimized by Iterative EM

- ▷ Pattern-by-pattern (PbP)
- ▷ Mean imputation (Mean.IMP)
- ▷ Fully specified (SAEM)
- ▷ Imputation by MICE (MICE.IMP): Iterative imputation by iterative MICE algorithm

Methods evaluated

86 / 99

- ▷ Pattern-by-pattern (PbP)
- ▷ Mean imputation (Mean.IMP)
- ▷ Fully specified (SAEM)
- ▷ Imputation by MICE (MICE.IMP)
 - Allow multiple imputations (MICE.K.IMP): Fit logistic on each dataset, average predictions



- ▷ Pattern-by-pattern (PbP)
- ▷ Mean imputation (Mean.IMP)
- ▷ Fully specified (SAEM)
- ▷ Imputation by MICE (MICE.IMP)
 - Allow multiple imputations (MICE.K.IMP)
 - Add M during imputation process (MICE.M.IMP)

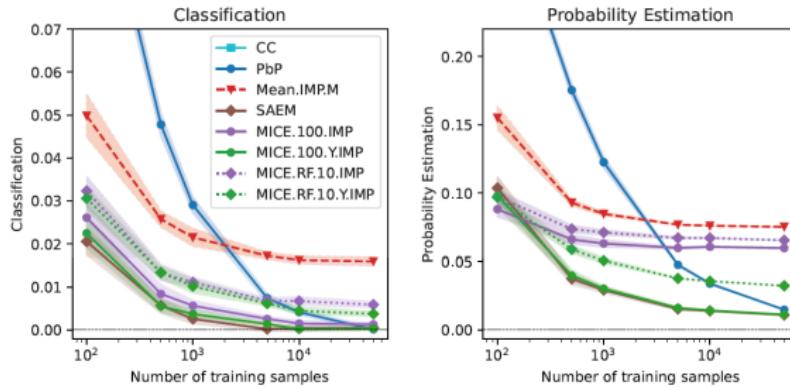
X'					
1	4	NA	0	0	1
NA	5	4	1	0	0
8	NA	4	0	1	0
NA	5	4	1	0	0
1	3	1	0	0	0
NA	1	4	1	0	0
7	8	4	0	0	0
4	5	NA	0	0	1

- ▷ Pattern-by-pattern (PbP)
- ▷ Mean imputation (Mean.IMP)
- ▷ Fully specified (SAEM)
- ▷ Imputation by MICE (MICE.IMP)
 - Allow multiple imputations (MICE.K.IMP)
 - Add M during imputation process (MICE.M.IMP)
 - Add Y during training of imputation process (MICE.Y.IMP)

X'

1	4	NA	0
NA	5	4	0
8	NA	4	1
NA	5	4	1
1	3	1	1
NA	1	4	0
7	8	4	0
4	5	NA	1

Gaussian features (MCAR)

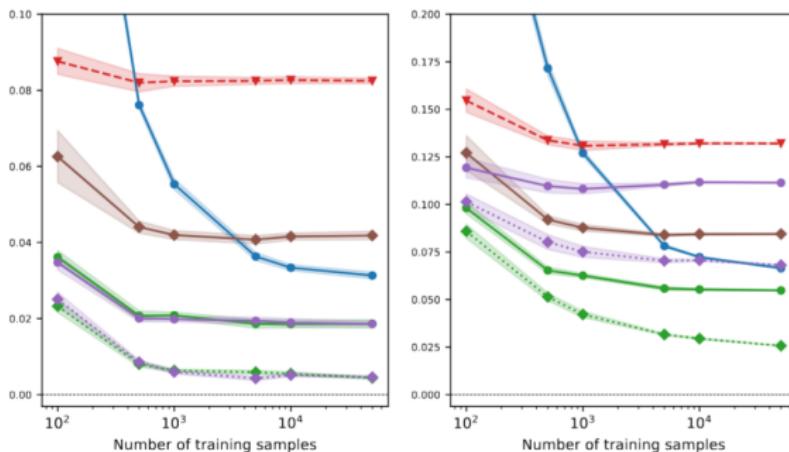


- ▷ $X \sim \mathcal{N}(0, \Sigma)$
- ▷ 5 dimensions
- ▷ 10 replicates
- ▷ Toeplitz correlation matrix (0.65 corr.)
- ▷ MCAR with prob. 0.25

-
- PbP approaching the Bayes prob. (large training set)
 - Necessary to use multiple imputations with MICE
 - Necessary to add Y to MICE imputation
 - SAEM and MICE.100.Y.IMP best overall

Non-linear features (MCAR)

88 / 99



- ▷ X non-linear transformation of $\mathcal{N}(0, \Sigma)$
- ▷ 5 dimensions
- ▷ 10 replicates
- ▷ Σ Toeplitz matrix (0.65)
- ▷ MCAR with prob 0.25

-
- No method can estimate Bayes probabilities
 - SAEM suffers from misspecification
 - PbP not approaching Bayes, coherent with our Theorem

Non-linear features (MCAR): details per pattern

89 / 99

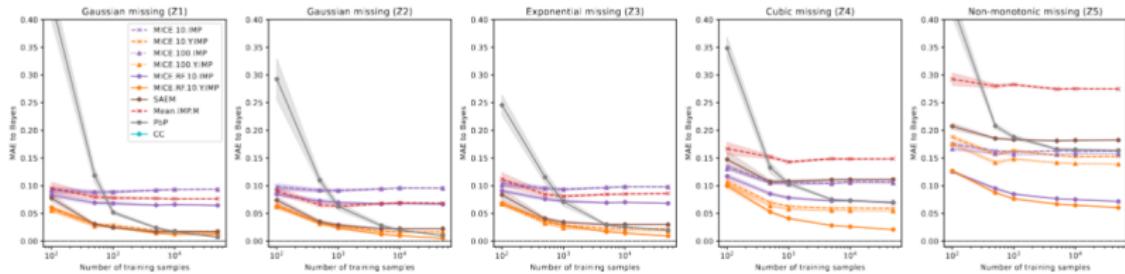
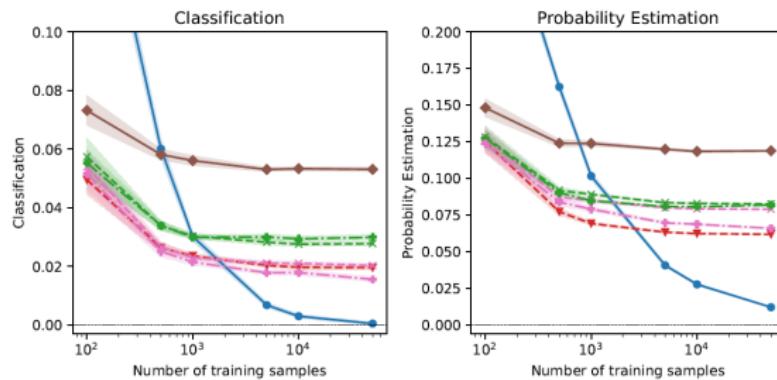


Figure C.7: Performances of selected procedures in terms of MAE from Bayes probabilities. The results are displayed by missing pattern in the test set (with one missing index: [1,0,0,0,0], ..., [0,0,0,0,1]). Means and standard errors over 10 replicates of non-linear features with MCAR missingness are displayed (see Section 5.1). The curves from MICE.10.IMP and MICE.100.IMP overlap in the first 4 plots.

Mixture of Gaussian (MNAR)

90 / 99



- ▷ $X|M = m \sim \mathcal{N}(\mu_m, \Sigma_m)$
- ▷ 5 dimensions
- ▷ 10 replicates
- ▷ Σ Toeplitz matrix
- ▷ MCAR with prob 0.25

-
- Only the PbP strategy performs well
 - Coherent with theory

- ▷ Theoretically, although misspecified, Pattern-by-pattern performs well under gaussian covariates (MCAR or Pattern Mixture Model)
- ▷ Confirmed experimentally: in GPMM-MNAR, PbP is one of the most competitive methods.

Empirically,

- ▷ MICE imputation consistently performing well in MCAR setting
 1. With the use of multiple imputations
 2. With the inclusion of Y in covariates
 3. But needs non-linear inner regressor for non-linear covariates
- ▷ M(N)AR settings are more tricky

1. Missing values mechanism
2. Single Imputation
3. Multiple Imputation
4. Imputation quality
5. Supervised Learning with Missing values
 - Decision trees as PbP predictors
 - Impute-then-regress procedures with consistent predictors
6. Linear models
 - Linear regression: A pattern-by-pattern approach
 - Linear regression: Impute-then-regress procedures via zero-imputation
 - Classification with missing values
7. Conclusion

Missing mechanisms

- ▷ Different missing data scenario (MCAR, MAR, MNAR).
- ▷ Both % of NA & structure matter (5% of NA can be an issue)
- ▷ MAR was designed for likelihood inference (e.g. EM algorithm) but can hide many complex distributions (distribution shift in MAR).
- ▷ Few implementations of EM strategies.

Missing mechanisms

- ▷ Different missing data scenario (MCAR, MAR, MNAR).
- ▷ Both % of NA & structure matter (5% of NA can be an issue)
- ▷ MAR was designed for likelihood inference (e.g. EM algorithm) but can hide many complex distributions (distribution shift in MAR).
- ▷ Few implementations of EM strategies.

Imputation

- ▷ Results in a complete data set, on which any method can be applied.
- ▷ *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.*

Missing mechanisms

- ▷ Different missing data scenario (MCAR, MAR, MNAR).
- ▷ Both % of NA & structure matter (5% of NA can be an issue)
- ▷ MAR was designed for likelihood inference (e.g. EM algorithm) but can hide many complex distributions (distribution shift in MAR).
- ▷ Few implementations of EM strategies.

Imputation

- ▷ Results in a complete data set, on which any method can be applied.
- ▷ *Imputation is both seductive & dangerous* (Dempster & Rubin, 1983).

Single imputation

- ▷ From simple (mean imputation) to more complex strategies (MissForest)
- ▷ Useful for point estimates
- ▷ Distort the marginal and joint distributions
- ▷ Lead to confidence interval with poor coverage

Single imputation

- ▷ From simple (mean imputation) to more complex strategies (MissForest)
- ▷ Useful for point estimates
- ▷ Distort the marginal and joint distributions
- ▷ Lead to confidence interval with poor coverage

Single imputation

- ▷ From simple (mean imputation) to more complex strategies (MissForest)
- ▷ Useful for point estimates
- ▷ Distort the marginal and joint distributions
- ▷ Lead to confidence interval with poor coverage

Multiple imputation

- ▷ Look for an imputation that preserve the joint distribution of the data
- ▷ MI aims at estimating the parameters and their variability taking into account the uncertainty of the missing values
- ▷ Useful for confidence intervals
- ▷ Compare imputations with distributional metrics like energy distance
- ▷ mice-DRF promising (code available) - mice-Engression^a

^aShen & Meinshausen (2024). Engression: extrapolation through the lens of distributional regression. *JRSS B*.

Aim

Estimating the Bayes predictor in presence of missing values

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

Two common strategies:

- ▷ **Impute-then-regress strategies** - impute the data then learn on the imputed data set
 - ◊ Computationally efficient but possibly inconsistent
- ▷ **Pattern-by-pattern strategies** - use a different predictor for each missing pattern
 - ◊ Consistent by design but intractable in most situations

Decision trees

- ▷ Decision trees are among the few methods able to natively handle missing values (MIA)
- ▷ Amounts to PbP strategies with a data-driven selection of relevant patterns

Impute-then-regress

- ▷ Consistent for any imputation method when the predictor is universally consistent
- ▷ Use the same imputation for train and test sets
- ▷ In finite sample, some imputation may ease the training of the predictor (e.g., Conditional Imputation is not well-suited in general)
- ▷ Rethinking imputation: a good imputation is the one that makes the prediction easy

Pattern-by-Pattern

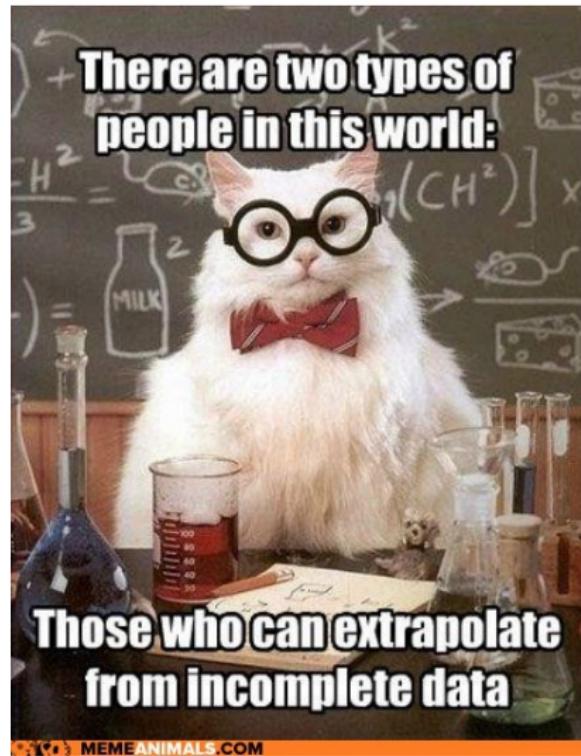
- ▷ Rate in $2^d/n$ in the worst case
- ▷ Improved by performing regressions on the most frequent patterns only
- ▷ Rate in d^2/n for MCAR Bernoulli, with a probability of missingness small enough
- ▷ MNAR/MAR is not suited for prediction (GPMM)

Impute-then-Regress

- ▷ Inconsistent in fixed dimension
- ▷ Consistent in high dimensions with a slow rate $n^{-1/2}$
- ▷ Imputation by zero amounts to a ridge regularization with a strength depending on the missing probability

Logistic regression model

- ▷ PbP and constant imputation result in inconsistent predictor
- ▷ But in presence of Gaussian features, Bayes probabilities are correctly estimated by PbP
- ▷ PbP competitive in GPMM-MNAR scenario but deteriorates when input distribution is not Gaussian



⁴⁹ More ressources: <https://rmisstastic.netlify.app/>