



HAL
open science

What's a good imputation to predict with missing values?

Marine Le Morvan, Julie Josse, Erwan Scornet, Gaël Varoquaux

► **To cite this version:**

Marine Le Morvan, Julie Josse, Erwan Scornet, Gaël Varoquaux. What's a good imputation to predict with missing values?. 2021. hal-03243931v1

HAL Id: hal-03243931

<https://hal.science/hal-03243931v1>

Preprint submitted on 31 May 2021 (v1), last revised 29 Nov 2021 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

What’s a good imputation to predict with missing values?

Marine Le Morvan^{1,2} Julie Josse⁴ Erwan Scornet³ Gaël Varoquaux¹

¹ Université Paris-Saclay, Inria, CEA, Palaiseau, 91120, France

² Université Paris-Saclay, CNRS/IN2P3, IJCLab, 91405 Orsay, France

³ CMAP, UMR7641, Ecole Polytechnique, IP Paris, 91128 Palaiseau, France

⁴ Inria Sophia-Antipolis, Montpellier, France

{marine.le-morvan, julie.josse, gael.varoquaux}@inria.fr
erwan.scornet@polytechnique.edu

Abstract

How to learn a good predictor on data with missing values? Most efforts focus on first imputing as well as possible and second learning on the completed data to predict the outcome. Yet, this widespread practice has no theoretical grounding. Here we show that for almost all imputation functions, an impute-then-regress procedure with a powerful learner is Bayes optimal. This result holds for all missing-values mechanisms, in contrast with the classic statistical results that require missing-at-random settings to use imputation in probabilistic modeling. Moreover, it implies that perfect conditional imputation may not be needed for good prediction asymptotically. In fact, we show that on perfectly imputed data the best regression function will generally be discontinuous, which makes it hard to learn. Crafting instead the imputation so as to leave the regression function unchanged simply shifts the problem to learning discontinuous imputations. Rather, we suggest that it is easier to learn imputation and regression jointly. We propose such a procedure, adapting NeuMiss, a neural network capturing the conditional links across observed and unobserved variables whatever the missing-value pattern. Experiments confirm that joint imputation and regression through NeuMiss is better than various two step procedures in our experiments with finite number of samples.

1 Introduction

Data with missing values are ubiquitous in many applications, as in health or business: some observations come with missing features. There is a rich statistical literature on imputation as well as inference with missing values [Rubin, 1976, Little and Rubin, 1987, 2002, 2019]. Most of the theory and practices build upon the *Missing At Random* (MAR) assumption that allows to maximize the likelihood of observed data while ignoring the missing-values mechanism, for instance using expectation maximization [Dempster et al., 1977]. On the opposite, Missing Not At Random settings, where missingness depends on the unobserved values, may not be identifiable and require dedicated methods with a model of the missing-values mechanism.

Learning predictive models with missing values poses distinct challenges compared to inference tasks [Josse et al., 2019]. Indeed, when the input is an arbitrary subset of variables in dimension d , there are 2^d potential missing data patterns and as many sub-models to learn. Consequently, even simple data-generating mechanisms lead to complex decision rules [Le Morvan et al., 2020b]. To date, there are few supervised-learning models natively suited for partially-observed data. A notable

exception is found with tree-based models [Twala et al., 2008, Chen and Guestrin, 2016], widely used in data-science practice.

The most common practice however remains by far to use off-the-shelf methods first for imputation of missing values and second for supervised-learning on the resulting completed data. Such a procedure may benefit from progress in missing-value imputation with machine learning [van Buuren 2018, Yoon et al. 2018]. However, there is a lack of learning theory to support such Impute-then-Regress procedures: Under what conditions are they Bayes consistent? Which aspects of the imputation are important?

There is empirical realization that the choice of imputation matters for predictive performance. The NADIA R package [Borowski and Fic] can select an imputation method to minimize a prediction error on a test set. Auto-ML is used to optimize full pipelines, including imputation [eg Jarrett et al., 2021]. Ipsen et al. [2020] optimize a constant imputation for supervised learning. However, the imputation is only weakly guided by the target in these approaches, it is set either from a family of black-box methods using gradient-free model selection, or from trivial imputation functions. In addition, there is a lack of insight on what drives a good imputation for prediction.

We contribute a systematic analysis of Impute-the-Regress procedures in a general setting: non-linear response function and any missingness mechanism (no MAR assumptions). We show that:

- Impute-then-Regress procedures are Bayes optimal for *all missing data mechanisms* and for *almost all imputation functions*, whatever the number of variables that may be missing. This very general result gives theoretical grounding to such widespread procedures.
- We study “natural” choices of imputation and regression functions: the oracle imputation by the conditional expectation and oracle regression function on the complete data. We show that chaining these oracles is not Bayes optimal in general and quantify its excess risk. We show that in both cases, choosing an oracle for one step, imputation or regression, imposes discontinuities on the other step, thus making it harder to learn.
- As these results suggest that imputation and regression should be adapted to one another, we contribute a method that jointly optimizes imputation and regression, using NeuMiss networks [Le Morvan et al., 2020a] as a differentiable imputation procedure.
- We compare empirically a number of Impute-then-Regress procedures on simulated non-linear regression tasks. Joint optimization of both steps provides the best performance.

2 Problem setting

Notations We consider a dataset of i.i.d. realizations of the random variable $(X, M, Y) \in \mathbb{R}^d \times \{0, 1\}^d \times \mathbb{R}$ where X are the complete covariates, M a missingness indicator, and Y a response of interest. For each realization (x, m, y) , $m_j = 1$ indicates that x_j is missing, and $m_j = 0$ that it is observed. We denote by $mis(m) \subset \llbracket 1, d \rrbracket$ the indices corresponding to the missing covariates (and similarly $obs(m)$ the observed indices), so that $x_{obs(m)}$ corresponds to the entries actually observed. We define the incomplete covariate vector $\tilde{X} \in (\mathbb{R} \cup \{\text{NA}\})^d$ as $\tilde{X}_j = X_j$ if $M_j = 0$ and $\tilde{X}_j = \text{NA}$ otherwise, where NA represents a missing value.

Assumptions We assume that X admits a density on \mathbb{R}^d and that, for all $j \in \llbracket 1, d \rrbracket$, each component X_j has finite expectation and variance, that is $\mathbb{E}[X_j^2] < \infty$. Moreover, we assume that the response Y is generated according to:

$$Y = f^*(X) + \epsilon, \quad \text{with } \mathbb{E}[\epsilon | X_{obs(M)}, M] = 0 \quad \text{and } \mathbb{E}[Y^2] < \infty. \quad (1)$$

where $f^* : \mathbb{R}^d \rightarrow \mathbb{R}$ is a function of the complete input data X , $\epsilon \in \mathbb{R}$ is a random noise variable.

2.1 Supervised learning with missing values

Optimization problem In practice, in the presence of missing values, we do not have access to the complete data (X, M, Y) but only to the subset of it that is observed, i.e. $(X_{obs(M)}, M, Y)$. Thus instead of learning a mapping from \mathbb{R}^d to \mathbb{R} , we need to learn a mapping from $(\mathbb{R} \cup \{\text{NA}\})^d$ to \mathbb{R} , where the set of observed covariates can be any subset of $\llbracket 1, d \rrbracket$. It is this unusual input space, partly discrete, that makes supervised learning with missing values challenging and different from classical

supervised learning problems. Formally, the optimization problem we wish to solve is:

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

Bayes predictor The function which minimizes (2), called the *Bayes predictor*, is given by:

$$\tilde{f}^*(\tilde{X}) = \mathbb{E} [Y | X_{obs(M)}, M] = \mathbb{E} [f^*(X) | X_{obs(M)}, M]. \quad (3)$$

As \tilde{X} is a function of X_{obs} and M , we will sometimes slightly abuse notations and write $\tilde{f}^*(\tilde{X}) = \tilde{f}^*(X_{obs}, M)$. The risk of the Bayes predictor is called the *Bayes risk*, which we denote as \mathcal{R}^* . It is the lowest achievable risk for a given supervised learning problem.

Definition 1 (Bayes optimality). A Bayes optimal function f achieves the Bayes rate, i.e. $\mathcal{R}(f) = \mathcal{R}^*$.

As can be seen from (3), the Bayes predictor is a function of M , a discrete random variable that can take one of 2^d values since $M \in \{0, 1\}^d$. The function \tilde{f}^* can thus be viewed as 2^d different functions, one for each possible subset of variables. This view raises questions that are central to this paper: How should we parametrize functions on such input domains? And which function families should we consider to approximate \tilde{f}^* ? These questions have been studied in the case where f^* is assumed to be a linear function, and X follows a Gaussian distribution. Indeed, under these assumptions, Le Morvan et al. [2020b,a] have derived analytical expressions for the Bayes predictor and deduced appropriate parametric estimators. However, aside from specific cases, it is impossible to derive an analytical expression for the Bayes predictor, and novel arguments are needed to understand which classes of functions should be considered in general.

3 Asymptotic analysis of Impute-then-regress procedures

3.1 Impute-then-regress procedures

Let $|mis(m)|$ (resp. $|obs(m)|$) be the number of missing entries (resp. observed) for any missing data pattern m . For each $m \in \{0, 1\}^d$, we define an *imputation function* $\phi^{(m)} : \mathbb{R}^{|obs(m)|} \rightarrow \mathbb{R}^{|mis(m)|}$ which outputs values for the missing entries based on the observed ones. We denote by $\phi_j^{(m)} : \mathbb{R}^{|obs(m)|} \rightarrow \mathbb{R}$ the component function of $\phi^{(m)}$ that imputes the j -th component in X if it is missing. Classical choices of imputation functions include constant functions or linear functions. Finally, we introduce the family of functions \mathcal{F}^I that transform an incomplete vector into a complete one, precisely:

$$\mathcal{F}^I = \left\{ \Phi : (\mathbb{R} \cup \{\text{NA}\})^d \rightarrow \mathbb{R}^d : \forall j \in \llbracket 1, d \rrbracket, \Phi_j(\tilde{X}) = \begin{cases} X_j & \text{if } M_j = 0 \\ \phi_j^{(M)}(X_{obs(M)}) & \text{if } M_j = 1 \end{cases} \right\}. \quad (4)$$

Let us define \mathcal{F}_∞^I in the exact same way but for imputation functions $\phi^{(m)} \in \mathcal{C}^\infty$, for all $m \in \{0, 1\}^d$. Here we study *Impute-then-regress procedures*, which we define as two-step procedures where the data is first imputed using a function $\Phi \in \mathcal{F}^I$, and then a regression is performed on the imputed data. Such a procedure is quite natural to deal with arbitrary subsets of inputs variables. It embeds the data into \mathbb{R}^d to reduce the problem to a classical one. In practice, impute-then-regress procedures are widely used. However, the choice of function class is so far mostly ad-hoc and raises multiple questions: How close to the Bayes rate can functions obtained via such procedures be? Should we prefer some choices of imputation functions over others? What happens when the missing data mechanism is missing not at random? In this section, we will give answers to these questions.

Below, we write *obs* (resp. *mis*) instead of $obs(M)$ (resp. $mis(M)$) to lighten notations.

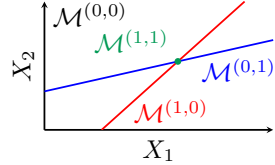
3.2 Impute-then-regress procedures are Bayes optimal

Definition 2 (Universal consistency). An estimator f_n is Bayes consistent if $\lim_{n \rightarrow \infty} \mathcal{R}(f_n) = \mathcal{R}^*$. It is said to be universally consistent if the previous statement holds for all distributions of (X, Y) .

The following theorem shows that Impute-then-regress procedures are Bayes optimal for almost all imputation functions. In other words, it means that a universal learner trained on imputed data provides optimal performances asymptotically for almost all imputation functions. Let us now define,

for all imputation function $\Phi \in \mathcal{F}^I$, the function $g_\Phi^* \in \underset{g: \mathbb{R}^d \rightarrow \mathbb{R}}{\operatorname{argmin}} \mathbb{E} \left[\left(Y - g \circ \Phi(\tilde{X}) \right)^2 \right]$.

Figure 1: **Example - Imputation manifolds in two dimensions** — Manifolds represented for linear imputation functions. $\mathcal{M}^{(0,0)}$ is the whole plane. Note that $\mathcal{M}^{(1,1)}$ need not be at the intersection of the two lines, it depends on the imputation function chosen. Here, $\mathcal{M}^{(0,1)}$ and $\mathcal{M}^{(1,0)}$ are transverse if and only if the two lines are not coincident.



Theorem 3.1 (Bayes consistency of Impute-then-regress procedures). *Assume the data is generated according to (1). Then, for almost all imputation function $\Phi \in \mathcal{F}_\infty^I$, the function $g_\Phi^* \circ \Phi$ is Bayes optimal. In other words, for almost all imputation functions $\Phi \in \mathcal{F}_\infty^I$, a universally consistent algorithm trained on the imputed data $\Phi(\tilde{X})$ is Bayes consistent.*

Appendix A.3 gives the proof. Theorem 3.1 states a very general result: Impute-then-regress procedures are Bayes consistent for all missing data mechanisms, almost all imputation functions, regardless of the distribution of (X, Y) and the number of missing covariates. Since Theorem 3.1 holds for almost all imputation functions, it implies that good imputations are not required to obtain good predictive performances, at least asymptotically. Note that here, the notion of *almost all* is to be understood in its topological sense, and not in its measure theory sense. Moreover, this theorem does not make any assumption on the missing data mechanism, and is therefore valid for Missing Not At Random (MNAR) data. This contrasts with most methods for inference and imputation with missing values, valid only for MAR data. Finally, the theorem remains valid for any configuration of variables that may contain missing values, including the case in which all variables may contain missing values. Bayes consistency of Impute-the-Regress procedures has already been studied, but in much more restricted settings. Josse et al. [2019] proved that such procedures are Bayes consistent under the MAR assumption, for constant imputations functions and for only one potentially missing variable. Bertsimas et al. [2021] refined this result to almost surely continuous imputation functions. While these two prior works build on very similar proofs, we use here very different arguments summarized in the next paragraph.

The first key idea of the proof is that, after imputation, all data points with a given missing data pattern m are mapped to a manifold $\mathcal{M}^{(m)}$ of dimension $|\text{obs}(m)|$. For example in 2D, data points are mapped to \mathbb{R}^2 when completely observed, to 1D manifolds when they have one value missing, and to one point when all values are missing (see fig. 1). Thus, Impute-then-Regress procedures first map data points to various manifolds depending on their missing data patterns and then apply a prediction function defined on the whole space including manifolds. The second key idea of the proof is to ensure that the missing data patterns of imputed points can almost surely be identified. For this, the proof requires that all manifolds of the same dimension are pairwise transverse. This assumption is sufficient, though not necessary, to ensure that the intersection of two manifolds of dimension $|\text{obs}(m)|$ cannot itself be of dimension $|\text{obs}(m)|$. Transversality is a weak assumption. In fact, Thom’s transversality theorem, (which we rely on in our proof) says that it is a generic property: it holds for “typical examples”, i.e *almost all* imputation functions will lead to transverse manifolds.

The proof is constructive and exhibits a function g_Φ^* which achieves the Bayes rate for a given set of imputation functions. For each manifold $\mathcal{M}^{(m)}$, ordered from smallest dimension to largest, we require that g_Φ^* on $\mathcal{M}^{(m)}$ equals the Bayes predictor for missing data pattern m except on points for which g_Φ^* has already been defined, i.e, the points where $\mathcal{M}^{(m)}$ intersects with the manifolds ranked before it. Thus, we obtain a function g_Φ^* that does not depend on m , and which for each manifold, equals the Bayes predictor except on subsets of measure zero under the assumption that manifolds of the same dimension are pairwise transverse. Refer to appendix A.3 for more details.

While this theorem is a very general result, it does not say what the optimal function associated to a given imputation looks like. In fact, depending on the imputation function it may be non-continuous, vary widely, and require a very large number of samples to be approximated.

4 Imputation versus regression: choosing one may break the other

Theorem 3.1 gives a theoretical grounding to Impute-then-regress procedures. As it holds for almost any imputation function, one could very well choose simple and cheap imputations such as imputing by a constant. However, the difficulty of the ensuing learning problem will depend on the choice of

imputation function. Indeed, the function g_{Φ}^* that achieves the Bayes rate depends on the imputation function Φ . In general, it may not be continuous or smooth. Thus g_{Φ}^* can be more or less difficult to approximate by machine learning algorithms depending on the chosen imputation function.

Le Morvan et al. [2020b] showed that even if Y is a linear function of X , imputing by a constant leads to a complicated Bayes predictor: piecewise affine but with 2^d regions. This result highlights how imputations neglecting the structure of covariates can result in additional complexity for the regression function g_{Φ}^* . Rather, another common practice is to impute by the conditional expectation: it minimizes the mean squared error between the imputed matrix and the complete one and is the target of most imputation methods. One hope may be that if f^* has desirable properties, such as smoothness, conditional imputation will lead to a function g_{Φ}^* which inherits from these properties.

In this section we first show that replacing missing values by their conditional expectation in the oracle regression function f^* gives a small but non-zero risk. Characterizing the optimal function on the conditionally-imputed data, we find that it suffers from discontinuities and thus forms a difficult estimation problem. Rather, we study whether the imputation can be corrected for f^* to form the Bayes predictor on partially-observed data.

4.1 Applying f^* on conditional imputations: chaining oracles isn't without risks.

The conditional imputation function $\Phi^{CI} : (\mathbb{R} \cup \{\text{NA}\})^d \rightarrow \mathbb{R}^d$ is defined as follows:

$$\forall j \in \llbracket 1, d \rrbracket, \Phi_j^{CI}(\tilde{X}) = \begin{cases} X_j & \text{if } M_j = 0 \\ \mathbb{E}[X_j | X_{\text{obs}}, M] & \text{if } M_j = 1 \end{cases}$$

Note that $\Phi^{CI} \in \mathcal{F}^I$. To lighten notations, we will write $X^{CI} := \Phi^{CI}(\tilde{X})$ to denote the conditionally imputed data.

Lemma 1 (First order approximation). *Assume that the data is generated according to (1). Moreover assume that (i) $f^* \in \mathcal{C}^2(\mathcal{S}, \mathbb{R})$ where $\mathcal{S} \subset \mathbb{R}^d$ is the support of the data, and that (ii) there exists positive semidefnite matrices $\bar{H}^+ \in P_d^+$ and $\bar{H}^- \in P_d^+$ such that for all X in \mathcal{S} , $\bar{H}^- \preceq H(X) \preceq \bar{H}^+$ with $H(X)$ the Hessian of f^* at X . Then for all X in \mathcal{S} and for all missing data patterns:*

$$\frac{1}{2} \text{tr}(\bar{H}_{\text{mis}, \text{mis}}^- \Sigma_{\text{mis} | \text{obs}, M}) \leq \tilde{f}^*(\tilde{X}) - f^*(X^{CI}) \leq \frac{1}{2} \text{tr}(\bar{H}_{\text{mis}, \text{mis}}^+ \Sigma_{\text{mis} | \text{obs}, M}) \quad (5)$$

where $\Sigma_{\text{mis} | \text{obs}, M}$ is the covariance matrix of the distribution of X_{mis} given X_{obs} and M .

Appendix A.5 gives the proof. The assumption that $\bar{H}^- \preceq H(X) \preceq \bar{H}^+$ for any X means that the minimum and maximum curvatures of f^* in any direction are uniformly bounded over the entire space. Lemma 1 shows that applying f^* to the conditionally imputed (CI) data is a good approximation of the Bayes predictor when there is no direction in which both the curvature of f^* and the conditional variance of the missing data given the observed one are high. Intuitively, if a low quality imputation is compensated by a flat function, or conversely, if a fast varying function is compensated by a high quality imputation, then f^* applied to the CI data approximates well the Bayes predictor.

Proposition 4.1 ((Non-)Consistency of chaining oracles). *The function $f^* \circ \Phi^{CI}$ is Bayes optimal if and only if the function f^* and the imputed data X^{CI} satisfy:*

$$\forall M \text{ s.t. } P(M) > 0, \quad \mathbb{E}[f^*(X) | X_{\text{obs}}, M] = f^*(X^{CI}) \quad \text{almost everywhere.} \quad (6)$$

Besides, under the assumptions of Lemma 1, the excess risk of chaining oracles compared to the Bayes risk \mathcal{R}^* is upper-bounded by:

$$\mathcal{R}(f^* \circ \Phi^{CI}) - \mathcal{R}^* \leq \frac{1}{4} \mathbb{E}_M \left[\max \left(\text{tr}(\bar{H}_{\text{mis}, \text{mis}}^- \Sigma_{\text{mis} | \text{obs}, M})^2, \text{tr}(\bar{H}_{\text{mis}, \text{mis}}^+ \Sigma_{\text{mis} | \text{obs}, M})^2 \right) \right] \quad (7)$$

Appendix A.6 gives the proof. Condition (6) for Bayes optimality is clearly stringent. Indeed, if one variable is missing, condition (6) says that the expectation of the regression function should be equal to the regression function applied at the expected entry. Although such functions are difficult to characterize precisely, it is clear that condition (6) is difficult to fulfill for generic regression functions (linear functions are among the few examples that do satisfy it). Therefore, for most functions f^* , $f^* \circ \Phi^{CI}$ is not Bayes optimal. Proposition 4.1 also gives an upper bound for the excess risk of the predictor $f^*(X^{CI})$ compared to the Bayes rate, showing here again that if there is no direction in

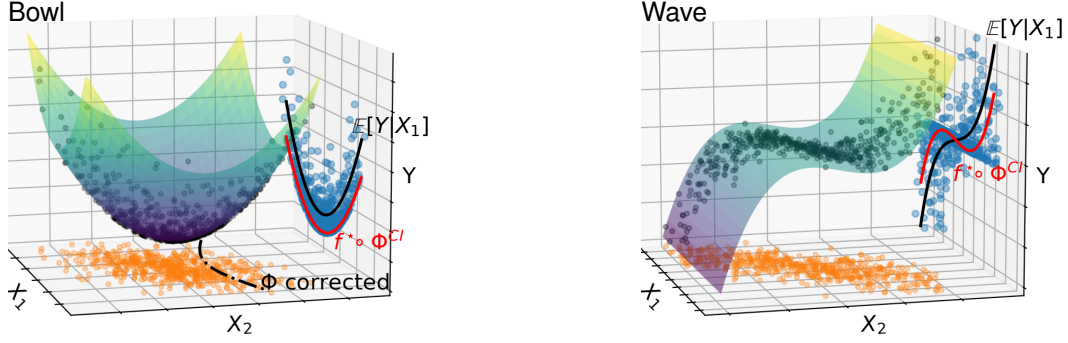


Figure 2: **Left: corrected imputation** The regression function is $f^*(x_1, x_2) \mapsto x_1^2 + x_2^2$. When x_2 is missing, chaining perfect conditional imputation with the regression function ($f^* \circ \Phi^{CI}$) gives a biased predictor, shown in red, as the unexplained variance in x_2 is turned into bias. However, using as an imputation $\Phi(x_1) = \sqrt{\rho^2 x_1^2 + (1 - \rho^2)}$ corrects this bias, with ρ the correlation between x_1 and x_2 . **Right: no continuous corrected imputation exists.** The function is defined as $f^*(x_1, x_2) \mapsto x_2^2 - 3x_2$. No continuous corrected imputation is possible because the Bayes predictor on the partially-observed data $\mathbb{E}[Y|X_1]$ is monotonous, while the regression function f^* is not.

which both the curvature and the variance of the missing data given the observed one are high, the excess risk is small.

The special case of linear regression: When f^* is a linear function, the curvature is 0, hence eq. (7) implies no excess risk. This is also visible from the expression of the Bayes predictor (3), where the expectation on unobserved data can be pushed inside f^* as it is linear. The Bayes predictor can thus be exactly written as f^* applied to conditionally-imputed data.

4.2 Regressing on conditional imputations, a good idea?

Proposition 4.2 (Regression function discontinuities). *Suppose that $f^* \circ \Phi^{CI}$ is not Bayes optimal, and that the probability of observing all variables is strictly positive, i.e., for all x , $P(M = (0, \dots, 0), X = x) > 0$. Then there is no continuous function g such that $g \circ \Phi^{CI}$ is Bayes optimal.*

In other words, when conditional imputation is used, the optimal regression function experiences discontinuities unless it is f^* . The proof is given in appendix A.7. From a finite-sample learning standpoint, discontinuous functions are in general harder to learn: in the general case, non-parametric regression requires more samples to achieve a given error on functions without specific regularities as opposed to functions with a form of smoothness [see e.g., Györfi et al., 2006, chap 3]. Hence, while regression on conditional imputation may be consistent (Theorem 3.1), it can require an inordinate number of samples.

4.3 Fasten your seat belt: corrected imputations may experience discontinuities.

Another possible route is to find *corrected imputations* which we define as imputation functions Φ such that, if f^* is used as regression function, the impute-then-regress procedure $f^* \circ \Phi$ is Bayes optimal. Intuitively, given a fixed regression function f^* , the question is: can we "correct" an imputation function and thus the manifold that it describes so that f^* restricted to this manifold is equal to the Bayes predictor? Assuming f^* is continuous, the intermediate value theorem gives a first answer to this question by ensuring the existence of imputations functions satisfying

$$f^* \circ \Phi(X_{obs(M)}, M) = \mathbb{E}[f^*(X)|X_{obs(M)}, M].$$

For the same reasons as above, determining that such imputations not only exist but are *continuous* is important from a practical perspective. Indeed, assuming f^* is continuous, the Bayes predictor with missing values could then be tackled as the composition of two continuous functions, with an Impute-then-Regress strategy. Intuitively in 2D, the existence of a continuous corrected imputation can be seen as the existence of a continuous path in the 2D plane whose value by f^* equals the Bayes predictor. Figure 2 (left) gives a simple example in 2D for which a continuous corrected imputation exists. However, as illustrated in Figure 2 (right), *continuous* corrected imputations do not always exist. Indeed, on this example the Bayes predictor is non-decreasing but there is no continuous path

in the 2D plane on which f^* is non-decreasing and maps at some point to both the 'purple' and 'yellow values' (proof in Appendix A.8). It is thus of interest to clarify when continuous corrected imputations exist. Proposition 4.3 establishes such conditions.

Proposition 4.3 (Existence of continuous corrected imputations). *Assume that f^* is uniformly continuous, twice continuously differentiable and that, for all missing patterns m and all x_{obs} , the support of $X_{mis}|X_{obs} = x_{obs}, M = m$ is connected. Additionally, assume that for all missing patterns m , and all (x_{obs}, x_{mis}) , the gradient of f^* with respect to the missing coordinates is nonzero:*

$$\nabla_{x_{mis}} f^*(x_{obs}, x_{mis}) \neq 0. \quad (8)$$

Then, for all m , there exist continuous imputation functions $\phi^{(m)} : \mathbb{R}^{|obs(m)|} \rightarrow \mathbb{R}^{|mis(m)|}$ such that $f^ \circ \Phi$ is Bayes optimal.*

Appendix A.9 gives a proof based on a global implicit function theorem. Assumption 8 is restrictive: it is for instance not met for our example in Figure 2 (left), which still admits continuous corrected imputations. This highlights the fact that continuous corrected imputations also exist under weaker conditions, but it is difficult to conclude on "how often" it is the case.

5 Jointly optimizing an impute-n-regress procedure: NeuMiss+MLP

The above suggests that it is beneficial to adapt the regression function to the imputation procedure and vice versa. Hence, we introduce a method for the joint optimization of these two steps by chaining a NeuMiss network, which learns an imputation, with an MLP (multi-layer perceptron).

NeuMiss [Le Morvan et al., 2020a] is a neural-network architecture originally designed to approximate the Bayes predictor for a linear model with missing values. It contains a theoretically-grounded Neumann block that can efficiently approximate the conditional expectation of the missing values given the observed ones. Here, we reuse the Neumann block as it outputs a learned imputation: each observed coordinate is mapped to its observed value and each missing coordinate is mapped to a function of the observed ones. The whole architecture can be seen as an Impute-then-Regress architecture, but that can be jointly optimized.

We performed a few minor improvements on the NeuMiss architecture. First, though the theory behind NeuMiss points to using shared weights in the Neumann block as well as residual connections going from the input to each hidden layer of the Neumann block, Le Morvan et al. [2020a] used neither. We found empirically that shared weights as well as residual connections improved performance. Moreover, while Le Morvan et al. [2020a] initialized the weights of the Neumann block randomly, we chose to initialize them with sample estimates of quantities that should be targeted to perform well according to Le Morvan et al. [2020a]. For clarity, the (non-linear) NeuMiss architecture and its initialization strategy are described in detail in Appendix A.11.

6 Empirical study of impute-n-regress procedures

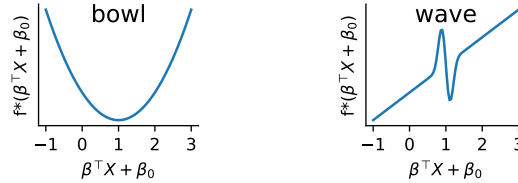
6.1 Experimental setup

Data generation The data $X \in \mathbb{R}^{n \times d}$ are generated according to a multivariate Gaussian distribution $\mathcal{N}(\mu, \Sigma)$ where the mean is drawn from a standard Gaussian $\mathcal{N}(0, Id)$ and the covariance is generated as $\Sigma = BB^\top + D$. $B \in \mathbb{R}^{d \times q}$ is a matrix with entries drawn from a standard normal Gaussian distribution, and D is a diagonal matrix with small entries that ensures that the covariance matrix is full rank. We study two correlation settings called *high* and *low* corresponding respectively to $q = \text{int}(0.3 * d)$ and $q = \text{int}(0.7 * d)$. The experiments are run with $d = 50$.

Choice of f^* The response Y is generated according to $Y = f^*(X) + \epsilon$ with two choices of f^* named *bowl* and *wave*, depicted in Figure 3 (exact expression in appendix A.10). β is a vector of ones normalized such that the quantity $z = \beta^\top X + \beta_0$ follows a Gaussian distribution centered on 1 with variance 1. Note that f_{bowl}^* and f_{wave}^* were designed so that the desired variations occur over the support of the data. The noise ϵ is chosen so as to have a signal-to-noise ratio of 10.

Missing values 50% of the entries of X were deleted according to one of two missing data mechanisms: Missing Completely At Random (MCAR) or Gaussian self-masking [GSM, see Le Morvan et al., 2020a]. Gaussian self-masking is a Missing Not At Random (MNAR) mechanism, where the probability that a variable j is missing depends on X_j via a Gaussian function.

Figure 3: **Bowl and wave functions** used for f^* in the empirical study.



Baseline methods benchmarked For each level of correlation (*low* or *high*), for each function f^* (*bowl* or *wave*), and each missing data mechanism (MCAR or GSM), we compare a number of methods. First, for reference, we compute various oracle predictors:

- **Bayes predictor:** This is the function that achieves the lowest achievable risk. In general cases, its expression cannot be derived analytically. However, we show that it can be derived for ridge functions, i.e. functions of the form $x \mapsto g(\beta^\top x)$, for some choices of g including polynomials and the Gaussian cdf. We thus built f_{bowl}^* and f_{wave}^* as combination of these base functions which allows us to compute their corresponding Bayes predictors. Appendix A.10 gives their expressions.
- **Chained oracles:** $f^* \circ \Phi^{CI}$ consists in imputing by the conditional expectation and then applying f^* . The analytical expression of Φ^{CI} can be derived analytically for both MCAR and GSM, and we thus use this analytical expression to impute.
- **Oracle + MLP:** The data is imputed using the analytical expression of the conditional expectation, and then a MLP is fitted to the completed data.

These three predictors all use ground truth information (parameters μ, Σ of the data distribution, expression of f^* or of the missing data mechanism) which are unavailable in practice. They are mainly useful as reference points. We then compare the NeurMiss+MLP architecture and a number of classic Impute-then-Regress methods as well as gradient boosted regression trees:

- **Mean + MLP** The data is imputed by the mean, and a multilayer perceptron (MLP) is fitted on the completed data.
- **MICE + MLP** The data is imputed using Scikit-learn’s [Pedregosa et al., 2012, BSD licensed] conditional imputer `IterativeImputer` that adapts the popular Multivariate Imputation by Chained Equations [MICE, van Buuren, 2018] to be able to impute a test set. A multilayer perceptron (MLP) is then fitted on the completed data.
- **GBRT:** Gradient boosted regression trees (Scikit-learn’s `HistGradientBoostingRegressor` with default parameters). This predictor readily supports missing values: during training, missing values on the decision variable for a given split are sent to the left or right child depending on which provides the largest gain, the Missing Incorporated Attribute strategy [Twala et al., 2008].

Finally, we also run **Mean + mask + MLP** as well as **MICE + mask + MLP** in which the mask is concatenated to the imputed data before fitting a MLP. Concatenating the mask is a widespread practice to account for MNAR data.

All the MLPs used have the same hyperparameters: 0, 1 or 2 hidden layers (chosen on a validation set), ReLU activation functions, and a width of 100 hidden neurons. Adam is used with an adaptive learning rate: the learning rate is divided by 5 each time 2 consecutive epochs fail to decrease the training loss by at least $1e-4$. Early stopping is triggered when the validation score does not improve by at least $1e-4$ for 10 consecutive epochs. Finally for NeuMiss the depth of the Neumann block is 5 or 15 (chosen on a validation set). The experiments use training sets of $n = 100\,000$ and $n = 500\,000$ and validation and test sets of size $n = 10\,000$. For $n = 100,000$, running all methods in one setting of correlation, missing data mechanisms, and choice of f^* , takes 3 hours on one core.

6.2 Experimental results

The results are presented in Figure 4.

Chaining oracles fail when both curvature is high and correlation is low. The chained oracle has a performance close to that of the Bayes predictor in all cases except when the wave function is applied to low correlation data. This observation illustrates well Proposition 4.1. Intuitively, the Bayes predictor for each missing data pattern is a smoothed version of f^* , and it is all the more

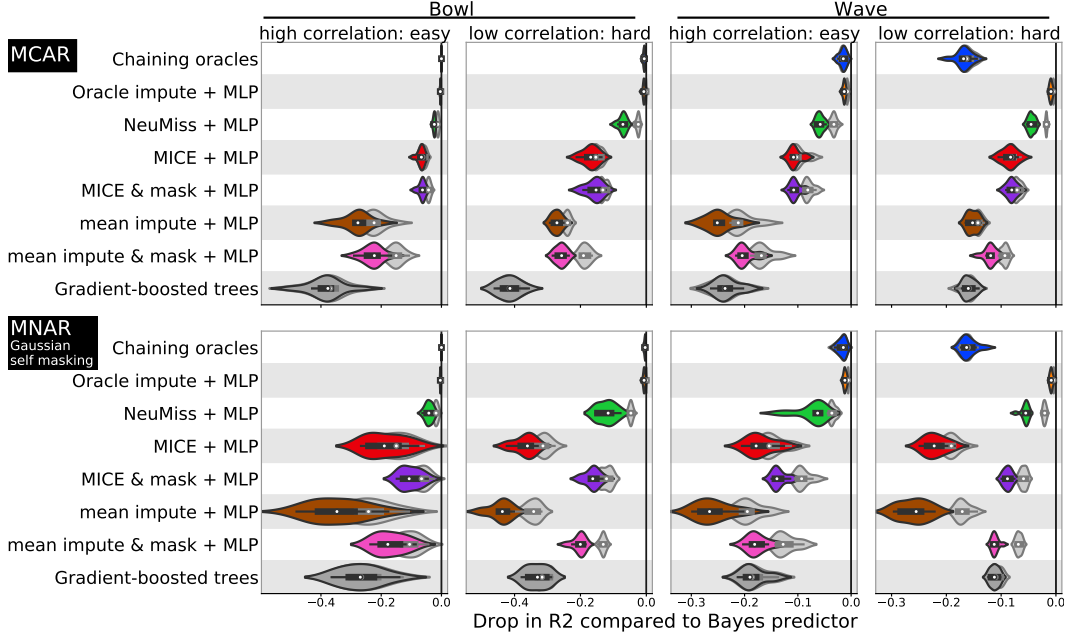


Figure 4: Performances (R2 score on a test set) compared to that of the Bayes predictor across 10 repeated experiments. The dark plots are for $n = 100\,000$ and the light ones for $n = 500\,000$.

smoothed that there is uncertainty around the likely values of the missing data. In the low correlation setting, the uncertainty is such that f^* is not a good proxy anymore for the Bayes predictor.

Regressing on oracle conditional imputation provide excellent performances. Contrary to the chained oracles, *Oracle + MLP* is close to the Bayes rate in all cases. This result should be put in perspective to Proposition 4.2, which states that there is no *continuous* regression function g such that $g \circ \Phi^{CI}$ is Bayes optimal if f^* is not. Indeed, as the MLP can only learn continuous functions, it shows that there are continuous functions g such that $g \circ \Phi$, albeit non-consistent, performs very well.

Adding the mask is critical in MNAR settings with *mean* and *MICE* imputations In MNAR settings, missingness carries information that can be useful for prediction. However, both the mean and iterative conditional imputation discard this information and output an imputed dataset in which the missingness information is lost or more difficult to retrieve. For this reason, it is common practice to concatenate the mask with the imputed data to expose the missingness information to the predictor. Our experiments show that under self-masking (MNAR), adding the mask to the mean or iteratively imputed data markedly improve performances. Note that NeuMiss does not require adding the mask as an input since the missingness information is already incorporated via the non-linearities.

NeuMiss+MLP performs best among Impute-then-Regress predictors. In *all* settings, NeuMiss performs best. GBRT perform poorly here possibly because they are not well adapted to approximate smooth functions. Finally, note that when the difficulty of the problem increases, for example with a lower correlation, then (i) the performance of the Bayes predictor decreases and (ii) the differences between the performances of methods is reduced, as in the lower right panel.

7 Conclusion

Impute-then-regress procedures assemble standard statistical routines to build predictors suited for data with missing values. However, we have shown that seeking the best prediction of the outcome leads to different tradeoffs compared to inferential purposes. Given a powerful learner, *almost all imputations* lead asymptotically to the optimal prediction, *whatever the missingness mechanism*. A good choice of imputation can however reduce the complexity of the function to learn. Though conditional expectation can lead to discontinuous optimal regression functions, our experiments show that it still leads to easier learning problems compared to simpler imputations. In order to adapt the

imputation to the regression function, we proposed to jointly learn these two steps by chaining a trainable imputation via the NeuMiss networks and a classical MLP. An empirical study of non-linear regression shows that it outperforms impute-then-regress procedures built on standard imputation methods as well as gradient-boosted trees with incorporated handling of missing values. In further work, it would be useful to theoretically characterize the learning behaviors of Impute-then-Regress methods in finite sample regimes.

References

- AV Arutyunov and SE Zhukovskiy. Application of methods of ordinary differential equations to global inverse function theorems. *Differential Equations*, 55(4):437–448, 2019.
- Dimitris Bertsimas, Arthur Delarue, and Jean Pauphilet. Prediction with Missing Data, 2021.
- Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2006. ISBN 0387310738.
- Jan Borowski and Piotr Fic. NADIA r package. <https://cran.r-project.org/web/packages/NADIA/index.html>. Accessed: 2021-05-26.
- Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.
- Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39:1, 1977.
- G B Folland. *Advanced Calculus*. Featured Titles for Advanced Calculus Series. Prentice Hall, 2002. ISBN 9780130652652.
- Martin Golubitsky. *Stable mappings and their singularities / M. Golubitsky, V. Guillemin*. Graduate texts in mathematics. Springer-Verlag, 1973. ISBN 0-387-90072-1.
- Victor W. Guillemin and Alan Pollack. *Differential topology*. Prentice-Hall Englewood Cliffs, N.J, 1974. ISBN 0132126052.
- László Györfi, Michael Kohler, Adam Krzyzak, and Harro Walk. *A distribution-free theory of nonparametric regression*. Springer Science & Business Media, 2006.
- Niels Ipsen, Pierre-Alexandre Mattei, and Jes Frellsen. How to deal with missing data in supervised deep learning? In *ICML Workshop on the Art of Learning with Missing Values (Artemiss)*, 2020.
- Daniel Jarrett, Jinsung Yoon, Ioana Bica, Zhaozhi Qian, Ari Ercole, and Mihaela van der Schaar. Clairvoyance: A pipeline toolkit for medical time series. In *International Conference on Learning Representations*, 2021.
- Julie Josse, Nicolas Prost, Erwan Scornet, and Gaël Varoquaux. On the consistency of supervised learning with missing values. *arXiv preprint arXiv:1902.06931*, 2019.
- Marine Le Morvan, Julie Josse, Thomas Moreau, Erwan Scornet, and Gaël Varoquaux. NeuMiss networks: differentiable programming for supervised learning with missing values. *Advances in Neural Information Processing Systems*, 33, 2020a.
- Marine Le Morvan, Nicolas Prost, Julie Josse, Erwan Scornet, and Gael Varoquaux. Linear predictor on linearly-generated data with missing values: non consistency and solutions. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 3165–3174. PMLR, 2020b.
- Roderick JA Little and Donald B Rubin. *Statistical analysis with missing data*. John Wiley & Sons, 1987, 2002, 2019.

- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems*, 32: 8026–8037, 2019.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12(oct): 2825–2830, 2012.
- Donald B Rubin. Inference and missing data. *Biometrika*, 63(3):581–592, 1976.
- B. E. T. H. Twala, M. C. Jones, and D. J. Hand. Good methods for coping with missing data in decision trees. *Pattern Recogn. Lett.*, 29:950–956, 2008. ISSN 0167-8655. doi: 10.1016/j.patrec.2008.01.010.
- Stef van Buuren. *Flexible Imputation of Missing Data, Second Edition*. 2018. doi: 10.1201/9780429492259.
- Jinsung Yoon, James Jordon, and Mihaela Schaar. Gain: Missing data imputation using generative adversarial nets. In *International Conference on Machine Learning*, page 5689. PMLR, 2018.

A Appendix

A.1 Proof of Lemma 2

Lemma 2. Let $\phi^{(m)} \in \mathcal{C}^\infty(\mathbb{R}^{|\text{obs}(m)|}, \mathbb{R}^{|\text{mis}(m)|})$ be the imputation function for missing data pattern m , and let $\mathcal{M}^{(m)} = \{x \in \mathbb{R}^d : x_{\text{mis}} = \phi^{(m)}(x_{\text{obs}})\}$. For all m , $\mathcal{M}^{(m)}$ is an $|\text{obs}|$ -dimensional manifold.

Proof. Let:

$$h^{(m)} : \mathbb{R}^d \rightarrow \mathbb{R}^{|\text{mis}|}$$

$$x \mapsto x_{\text{mis}} - \phi^{(m)}(x_{\text{obs}})$$

Regular value: We will show that $\mathbf{0}_{\text{mis}}$ is a regular value of $h^{(m)}$. By definition [see p21 in Guillemin and Pollack, 1974], a point $y \in \mathbb{R}^{|\text{mis}|}$ is a regular value of $h^{(m)}$ if $dh_x^{(m)}$ is surjective at every point x such that $h^{(m)}(x) = y$. The mapping $dh_x^{(m)}$ is linear and can be represented by the Jacobian of $h^{(m)}$ at x :

$$J_{h^{(m)}}(x) = \left(\begin{array}{c|c} A & Id \end{array} \right), \quad A \in \mathbb{R}^{|\text{mis}| \times |\text{obs}|}, \quad Id \in \mathbb{R}^{|\text{mis}| \times |\text{mis}|}.$$

Given the structure of $J_{h^{(m)}}(x)$, it is obviously of rank $|\text{mis}|$ at every point x . Thus $dh_x^{(m)}$ is surjective at every point x , and it is true in particular for the points x such that $h^{(m)}(x) = \mathbf{0}$. We conclude that by definition, $\mathbf{0}_{\text{mis}}$ is a regular value of $h^{(m)}$.

Preimage theorem: By the Preimage theorem ([Guillemin and Pollack, 1974], p.21), since $\mathbf{0} \in \mathbb{R}_{\text{mis}}$ is a regular value of $h^{(m)} : \mathbb{R}^d \rightarrow \mathbb{R}^{|\text{mis}|}$, then the the preimage $(h^{(m)})^{-1}(\mathbf{0})$ is a submanifold of \mathbb{R}^d of dimension $d - |\text{mis}| = |\text{obs}|$.

Since by definition, $(h^{(m)})^{-1}(\mathbf{0}) = \mathcal{M}^{(m)}$, we have that $\mathcal{M}^{(m)}$ is a $|\text{obs}|$ -dimensional manifold. \square

A.2 Proof of Lemma 3

Lemma 3. Let m and m' be two distinct missing data patterns with the same number of missing values $|\text{mis}|$. Let $\phi^{(m)} \in \mathcal{C}^\infty(\mathbb{R}^{|\text{obs}(m)|}, \mathbb{R}^{|\text{mis}(m)|})$ be the imputation function for missing data pattern m , and let $\mathcal{M}^{(m)} = \{x \in \mathbb{R}^d : x_{\text{mis}} = \phi^{(m)}(x_{\text{obs}})\}$. We define similarly $\phi^{(m')}$ and $\mathcal{M}^{(m')}$. For almost all imputation functions $\phi^{(m)}$ and $\phi^{(m')}$,

$$\dim(\mathcal{M}^{(m)} \cap \mathcal{M}^{(m')}) = \begin{cases} 0 & \text{if } |\text{mis}| > \frac{d}{2} \\ d - 2|\text{mis}| & \text{otherwise.} \end{cases} \quad (9)$$

Proof. According to Thom Transversality theorem ([Golubitsky, 1973], p.54) with:

- $W = \mathcal{M}^{(m')}$,
- $f = \phi^{(m)}$,
- $k = 0$ (note that as stated p.37, $J^0(X, Y) = X \times Y$ and $j^0 f(x) = \text{graph}(f)$),

we have that $\{\phi^{(m)} \in \mathcal{C}^\infty(\mathbb{R}^{|\text{obs}|}, \mathbb{R}^{|\text{mis}|}) \mid \text{graph}(\phi^{(m)}) \pitchfork \mathcal{M}^{(m')}\}$ is a residual subset of $\mathcal{C}^\infty(\mathbb{R}^{|\text{obs}|}, \mathbb{R}^{|\text{mis}|})$ in the \mathcal{C}^∞ topology. In other words, the fact that $\text{graph}(\phi^{(m)})$ is transverse to $\mathcal{M}^{(m')}$ is a generic property. Put differently, almost all functions $\phi^{(m)}$ have their graph transverse to $\mathcal{M}^{(m')}$. Note that here the notion of *almost all* has to be understood in its topological sense, and not in its measure theory sense.

Suppose that $|\text{obs}| < \frac{d}{2}$. According to Lemma 2, $\mathcal{M}^{(m')}$ is a $|\text{obs}|$ -dimensional manifold. Moreover we just showed that for almost all $\phi^{(m)}$, $\text{graph}(\phi^{(m)}) \pitchfork \mathcal{M}^{(m')}$. Applying Proposition 4.2 of

[Golubitsky, 1973] (p.51) with $W = \mathcal{M}^{(m')}$ and $f = \text{graph}(\phi^{(m)})$, we obtain that $\mathcal{M}^{(m)}$ and $\mathcal{M}^{(m')}$ are disjoint, since, by definition, $\mathcal{M}^{(m)}$ is the image of $\text{graph}(\phi^{(m)})$. Consequently, the dimension of their intersection is 0.

Suppose that $|\text{obs}| \geq \frac{d}{2}$. According to the theorem p.30 of [Guillemin and Pollack, 1974], since $\mathcal{M}^{(m)}$ and $\mathcal{M}^{(m')}$ are transverse submanifolds of \mathbb{R}^d , their intersection is again a manifold with $\text{codim}(\mathcal{M}^{(m)} \cap \mathcal{M}^{(m')}) = \text{codim}(\mathcal{M}^{(m)}) + \text{codim}(\mathcal{M}^{(m')})$. This implies that $\dim(\mathcal{M}^{(m)} \cap \mathcal{M}^{(m')}) = 2|\text{obs}| - d$. \square

A.3 Proof of Theorem 3.1

Theorem 3.1 (Bayes consistency of Impute-then-regress procedures). *Assume the data is generated according to (1). Then, for almost all imputation function $\Phi \in \mathcal{F}_{\infty}^I$, the function $g_{\Phi}^* \circ \Phi$ is Bayes optimal. In other words, for almost all imputation functions $\Phi \in \mathcal{F}_{\infty}^I$, a universally consistent algorithm trained on the imputed data $\Phi(\tilde{X})$ is Bayes consistent.*

Proof. Let $\phi^{(m)} \in \mathcal{C}^{\infty}(\mathbb{R}^{|\text{obs}(m)|}, \mathbb{R}^{|\text{mis}(m)|})$ be the imputation function for missing data pattern m , and let $\mathcal{M}^{(m)} = \{x \in \mathbb{R}^d : x_{\text{mis}} = \phi^{(m)}(x_{\text{obs}})\}$. According to Lemma 2, for all m , $\mathcal{M}^{(m)}$ is an $|\text{obs}|$ -dimensional manifold. $\mathcal{M}^{(m)}$ corresponds to the subspace where all points with missing data pattern m are mapped after imputation.

Let us order missing data patterns according to their number of missing values, with the pattern of all missing entries ranked first and the pattern of all observed entries ranked last. Two patterns with the same number of missing values are ordered arbitrarily. We use $m(i)$ to refer to the missing data pattern ranked in i^{th} position.

Let g^* be the function defined as follows: for all i ,

$$\forall Z = \Phi(\tilde{X}) \in \mathcal{M}^{(m(i))} \setminus \bigcup_{m(k) < m(i)} \mathcal{M}^{(m(k))}, \quad g^*(Z) = \tilde{f}^*(\tilde{X}).$$

For a given missing data pattern $m(i)$, by distributivity of intersections across unions, we have:

$$\mathcal{M}^{(m(i))} \cap \left(\bigcup_{m(k) < m(i)} \mathcal{M}^{(m(k))} \right) = \bigcup_{m(k) < m(i)} \left(\mathcal{M}^{(m(i))} \cap \mathcal{M}^{(m(k))} \right)$$

If $m(k)$ has strictly more missing values than $m(i)$, then by Lemma 2 $\dim(\mathcal{M}^{(m(k))}) < \dim(\mathcal{M}^{(m(i))})$, and thus $\dim(\mathcal{M}^{(m(k))} \cap \mathcal{M}^{(m(i))}) < \dim(\mathcal{M}^{(m(i))})$. Moreover, If $m(k)$ has the same number of missing values as $m(i)$, then by Lemma 3, for almost all imputation functions $\phi^{(m(k))}$ and $\phi^{(m(i))}$, $\dim(\mathcal{M}^{(m(k))} \cap \mathcal{M}^{(m(i))}) < \dim(\mathcal{M}^{(m(i))})$. We conclude that for all $m(k) < m(i)$, $\mathcal{M}^{(m(k))} \cap \mathcal{M}^{(m(i))}$ is a subset of measure zero in $\mathcal{M}^{(m(i))}$. Finally, since a countable union of sets of measure zero has measure zero, we obtain that $\bigcup_{m(k) < m(i)} (\mathcal{M}^{(m(i))} \cap \mathcal{M}^{(m(k))})$ has measure zero in $\mathcal{M}^{(m(i))}$.

Let's now compute the risk of $g^* \circ \Phi$:

$$\mathcal{R}(g^* \circ \Phi) = \sum_{M=m} P(M=m) \int_{X_{\text{obs}}} P(X_{\text{obs}}|M=m) \left(\tilde{f}^*(\tilde{X}) - g^* \circ \Phi(\tilde{X}) \right)^2 \quad (10)$$

For a given missing data pattern m , $\Phi(\tilde{X}) \in \mathcal{M}^{(m)}$. Moreover, we constructed g^* such that $g^* \circ \Phi(\tilde{X}) = \tilde{f}^*(\tilde{X})$ for all $\Phi(\tilde{X}) \in \mathcal{M}^{(m)}$ except on a set that we just showed to be of measure zero for almost all imputation functions. As a result, the function $X_{\text{obs}} \mapsto \tilde{f}^*(\tilde{X}) - g^* \circ \Phi(\tilde{X})$ is zero almost everywhere for a given m , and the function $X_{\text{obs}} \mapsto P(X_{\text{obs}}|M=m) \left(\tilde{f}^*(\tilde{X}) - g^* \circ \Phi(\tilde{X}) \right)^2$ is also zero almost everywhere. Since the integral of a function that vanishes almost everywhere is equal to 0, we conclude that $\mathcal{R}(g^* \circ \Phi) = 0$. Since the risk cannot be negative, $g^* \circ \Phi$ is a minimizer of the risk and thus it is Bayes optimal. \square

A.4 Proof of Lemma 4

Lemma 4.

$$\forall X \in \mathbb{R}^p, \forall mis \subseteq \llbracket 1, p \rrbracket, H(X) \preceq \bar{H}^+ \implies H_{mis,mis}(X) \preceq \bar{H}_{mis,mis}^+$$

Proof. Let $X \in \mathbb{R}^p$, and let m be a missing data pattern with observed (resp. missing) indices obs (resp. mis). $H(X) \preceq \bar{H}^+$ is equivalent to:

$$\forall u \in \mathbb{R}^p, u^\top (\bar{H}^+ - H(X)) u \geq 0. \quad (11)$$

Let $\mathcal{V} \subseteq \mathbb{R}^p$ be a subspace such that for any v in \mathcal{V} , $v_{obs} = 0$. Since $\mathcal{V} \subseteq \mathbb{R}^p$, (11) implies:

$$\begin{aligned} \forall v \in \mathcal{V}, v^\top (\bar{H}^+ - H(X)) v &\geq 0 \\ \iff \forall v_{mis} \in \mathbb{R}^{|mis|}, v_{mis}^\top (\bar{H}_{mis,mis}^+ - H_{mis,mis}(X)) v_{mis} &\geq 0 \\ \iff H_{mis,mis}(X) &\preceq \bar{H}_{mis,mis}^+ \end{aligned}$$

□

A.5 Proof of Lemma 1

Lemma 1 (First order approximation). *Assume that the data is generated according to (1). Moreover assume that (i) $f^* \in \mathcal{C}^2(\mathcal{S}, \mathbb{R})$ where $\mathcal{S} \subset \mathbb{R}^d$ is the support of the data, and that (ii) there exists positive semidefnite matrices $\bar{H}^+ \in P_d^+$ and $\bar{H}^- \in P_d^+$ such that for all X in \mathcal{S} , $\bar{H}^- \preceq H(X) \preceq \bar{H}^+$ with $H(X)$ the Hessian of f^* at X . Then for all X in \mathcal{S} and for all missing data patterns:*

$$\frac{1}{2} \text{tr} (\bar{H}_{mis,mis}^- \Sigma_{mis|obs,M}) \leq \tilde{f}^*(\tilde{X}) - f^*(X^{CI}) \leq \frac{1}{2} \text{tr} (\bar{H}_{mis,mis}^+ \Sigma_{mis|obs,M}) \quad (5)$$

where $\Sigma_{mis|obs,M}$ is the covariance matrix of the distribution of X_{mis} given X_{obs} and M .

Proof. Without loss of generality, suppose that we reorder variables such that we can write $X = (X_{obs}, X_{mis})$. Consider the function

$$\begin{aligned} f_{mis}^* : \mathbb{R}^{|mis|} &\rightarrow \mathbb{R} \\ X_{mis} &\mapsto f^*(X_{obs}, X_{mis}) \end{aligned}$$

Since $f^* \in \mathcal{C}^2(\mathbb{R}^d, \mathbb{R})$, we have $f_{mis}^* \in \mathcal{C}^2(\mathbb{R}^{|mis|}, \mathbb{R})$. Therefore, we can write the first order Taylor expansion (see Theorem 2.68 in Folland [2002]) of f_{mis}^* around $E[X_{mis}|X_{obs}, M]$:

$$\begin{aligned} f_{mis}^*(X_{mis}) &= f^*(X_{obs}, \mathbb{E}[X_{mis}|X_{obs}, M]) \\ &\quad + \nabla f_{mis}^*(X_{obs}, \mathbb{E}[X_{mis}|X_{obs}, M])^\top (X_{mis} - \mathbb{E}[X_{mis}|X_{obs}, M]) \\ &\quad + R(X_{mis} - \mathbb{E}[X_{mis}|X_{obs}, M]), \end{aligned} \quad (12)$$

where R is the Lagrange remainder satisfying

$$\begin{aligned} R(X_{mis} - \mathbb{E}[X_{mis}|X_{obs}, M]) &= \\ &\quad \frac{1}{2} (X_{mis} - \mathbb{E}[X_{mis}|X_{obs}, M])^\top H_{mis,mis}(c) (X_{mis} - \mathbb{E}[X_{mis}|X_{obs}, M]), \end{aligned}$$

for some c in the ball $\mathcal{B}(\mathbb{E}[X_{mis}|X_{obs}, M], \|X_{mis} - \mathbb{E}[X_{mis}|X_{obs}, M]\|_2)$. By assumption, for all X , $H(X) \preceq \bar{H}^+$. Therefore, according to Lemma 4, we have $H_{mis,mis}(X) \preceq \bar{H}_{mis,mis}^+$ for any missing data pattern, which leads to:

$$\begin{aligned} R(X_{mis} - \mathbb{E}[X_{mis}|X_{obs}, M]) &\leq \\ &\quad \frac{1}{2} (X_{mis} - \mathbb{E}[X_{mis}|X_{obs}, M])^\top \bar{H}_{mis,mis}^+ (X_{mis} - \mathbb{E}[X_{mis}|X_{obs}, M]). \end{aligned}$$

Using equality (12), we get:

$$\begin{aligned} f^*(X_{obs}, X_{mis}) - f^*(X_{obs}, \mathbb{E}[X_{mis}|X_{obs}, M]) &= \\ &\quad - \nabla f_{mis}^*(X_{obs}, \mathbb{E}[X_{mis}|X_{obs}, M])^\top (X_{mis} - \mathbb{E}[X_{mis}|X_{obs}, M]) \\ &\leq \frac{1}{2} (X_{mis} - \mathbb{E}[X_{mis}|X_{obs}, M])^\top \bar{H}_{mis,mis}^+ (X_{mis} - \mathbb{E}[X_{mis}|X_{obs}, M]) \end{aligned}$$

Finally, taking the expectation with regards to $P(X_{mis}|X_{obs}, M)$ on both sides, we obtain

$$\mathbb{E}[f^*(X_{obs}, X_{mis})|X_{obs}, M] - f^*(X_{obs}, \mathbb{E}[X_{mis}|X_{obs}, M]) \leq \frac{1}{2} \text{tr}(H_{mis, mis}^{+\top} \Sigma_{mis|obs, M}), \quad (13)$$

where we have used the fact that, for any vector $X \in \mathbb{R}^d$ and for any $H \in P_d^+$,

$$X^\top H X = \text{tr}(X^\top H X) = \text{tr}(H X X^\top).$$

Following a similar reasoning, we can show that:

$$\mathbb{E}[f^*(X_{obs}, X_{mis})|X_{obs}, M] - f^*(X_{obs}, \mathbb{E}[X_{mis}|X_{obs}, M]) \geq \frac{1}{2} \text{tr}(H_{mis, mis}^{-\top} \Sigma_{mis|obs, M}) \quad (14)$$

Together, inequalities (13) and (14) conclude the proof. \square

A.6 Proof of Proposition 4.1

Proposition 4.1 ((Non-)Consistency of chaining oracles). *The function $f^* \circ \Phi^{CI}$ is Bayes optimal if and only if the function f^* and the imputed data X^{CI} satisfy:*

$$\forall M \text{ s.t. } P(M) > 0, \quad \mathbb{E}[f^*(X)|X_{obs}, M] = f^*(X^{CI}) \quad \text{almost everywhere.} \quad (6)$$

Besides, under the assumptions of Lemma 1, the excess risk of chaining oracles compared to the Bayes risk \mathcal{R}^* is upper-bounded by:

$$\mathcal{R}(f^* \circ \Phi^{CI}) - \mathcal{R}^* \leq \frac{1}{4} \mathbb{E}_M \left[\max \left(\text{tr}(\bar{H}_{mis, mis}^- \Sigma_{mis|obs, M})^2, \text{tr}(\bar{H}_{mis, mis}^+ \Sigma_{mis|obs, M})^2 \right) \right] \quad (7)$$

Proof.

$$Y - f^*(X^{CI}) = (Y - \tilde{f}^*(\tilde{X})) + (\tilde{f}^*(\tilde{X}) - f^*(X^{CI})) \quad (15)$$

$$(Y - f^*(X^{CI}))^2 = (Y - \tilde{f}^*(\tilde{X}))^2 + (\tilde{f}^*(\tilde{X}) - f^*(X^{CI}))^2 \quad (16)$$

$$+ 2(Y - \tilde{f}^*(\tilde{X}))(\tilde{f}^*(\tilde{X}) - f^*(X^{CI})) \quad (17)$$

$$= (Y - \tilde{f}^*(\tilde{X}))^2 + (\tilde{f}^*(\tilde{X}) - f^*(X^{CI}))^2 \quad (18)$$

$$+ 2(f^*(X) - \tilde{f}^*(\tilde{X}))(\tilde{f}^*(\tilde{X}) - f^*(X^{CI})) \quad (19)$$

$$+ 2\epsilon(\tilde{f}^*(\tilde{X}) - f^*(X^{CI})) \quad (20)$$

$$\mathbb{E} \left[(Y - f^*(X^{CI}))^2 \right] = \mathcal{R}^* + \mathbb{E} \left[(\tilde{f}^*(\tilde{X}) - f^*(X^{CI}))^2 \right] \quad (21)$$

where we used the definition of the Bayes rate. Moreover, term (20) vanishes when taking the expectation w.r.t ϵ because $\mathbb{E}[\epsilon|X_{obs}, M] = 0$ and ϵ is uncorrelated with X or M , and term (19) vanishes when taking the expectation w.r.t $X_{mis}|X_{obs}, M$ because by definition $\mathbb{E}_{X_{mis}|X_{obs}, M}[f^*(X_{obs}, X_{mis})] = \tilde{f}^*(\tilde{X})$.

Clearly, $f^* \circ \Phi^{CI}$ is Bayes optimal if and only if:

$$\mathbb{E} \left[(\tilde{f}^*(\tilde{X}) - f^*(X^{CI}))^2 \right] = 0 \quad (22)$$

$$\iff \sum_M \int P(X_{obs}, M) (\tilde{f}^*(\tilde{X}) - f^*(X^{CI}))^2 dX_{obs} = 0 \quad (23)$$

$$\iff \forall M, X_{obs} : P(X_{obs}, M) > 0, \quad \tilde{f}^*(\tilde{X}) = f^*(X^{CI}) \quad \text{almost everywhere.} \quad (24)$$

where equality 24 is true since all terms are positive.

Besides, by Lemma 1, we have:

$$\frac{1}{2} \text{tr}(\bar{H}_{mis, mis}^- \Sigma_{mis|obs, M}) \leq \tilde{f}^*(\tilde{X}) - f^*(X^{CI}) \leq \frac{1}{2} \text{tr}(\bar{H}_{mis, mis}^+ \Sigma_{mis|obs, M}). \quad (25)$$

By convexity of the square function, it follows that:

$$\left(\tilde{f}^*(\tilde{X}) - f^*(X^{CI})\right)^2 \leq \frac{1}{2} \max\left(\text{tr}\left(\bar{H}_{mis,mis}^- \Sigma_{mis|obs,M}\right)^2, \text{tr}\left(\bar{H}_{mis,mis}^+ \Sigma_{mis|obs,M}\right)^2\right). \quad (26)$$

Finally, by taking the expectation on both sides:

$$\begin{aligned} \mathbb{E}\left[\left(\tilde{f}^*(\tilde{X}) - f^*(X^{CI})\right)^2\right] &\leq \\ &\frac{1}{2}\mathbb{E}_M\left[\max\left(\text{tr}\left(\bar{H}_{mis,mis}^- \Sigma_{mis|obs,M}\right)^2, \text{tr}\left(\bar{H}_{mis,mis}^+ \Sigma_{mis|obs,M}\right)^2\right)\right]. \end{aligned} \quad (27)$$

Combining equation (21) with inequality (27) concludes the proof. \square

A.7 Proof of Proposition 4.2

Proposition 4.2 (Regression function discontinuities). *Suppose that $f^* \circ \Phi^{CI}$ is not Bayes optimal, and that the probability of observing all variables is strictly positive, i.e., for all x , $P(M = (0, \dots, 0), X = x) > 0$. Then there is no continuous function g such that $g \circ \Phi^{CI}$ is Bayes optimal.*

Proof. We will prove this result by contradiction. Suppose that (i) $f^* \circ \Phi^{CI}$ is not Bayes optimal, (ii) the probability of observing all variables is strictly positive, (iii) there exists a function g continuous such that $g \circ \Phi^{CI}$ is Bayes optimal.

Following a reasoning similar to the one in the proof of proposition 4.1, we can show that $g \circ \Phi^{CI}$ is Bayes optimal if and only if:

$$\forall M, X_{obs} : P(X_{obs}, M) > 0, \quad \mathbb{E}[f^*(X)|X_{obs}, M] = g(X^{CI}) \quad \text{almost everywhere.}$$

In particular since for all x , the joint probability $P(M = (0, \dots, 0), X = x)$ of observing all variables is strictly positive, g should satisfy this equality for $M = (0, \dots, 0)$, i.e.:

$$f^*(X) = g(X) \quad \text{almost everywhere.}$$

Since g is continuous, it implies $g = f^*$. Since by assumption, f^* is not Bayes optimal, then g is not either, which is a contradiction. \square

A.8 Example of a case where no continuous corrected imputation exists.

Let:

$$\begin{aligned} f^* : \mathbb{R}^2 &\rightarrow \mathbb{R} \\ (X_1, X_2) &\mapsto X_2^3 - 3X_2 \end{aligned}$$

and let:

$$\begin{aligned} X_2 = X_1 + \epsilon \quad \text{with} \quad \mathbb{E}[\epsilon|X_1, M = (0, 1)] &= 0 \\ \mathbb{E}[\epsilon^2|X_1, M = (0, 1)] &= \sigma^2, \sigma^2 > 1 \\ \mathbb{E}[\epsilon^3|X_1, M = (0, 1)] &= 0 \end{aligned}$$

Suppose that X_2 is missing. Then the Bayes predictor is given by:

$$\begin{aligned} \tilde{f}^*(X_1, M = (0, 1)) &= \mathbb{E}[f^*(X)|X_1, M = (0, 1)] \\ &= \mathbb{E}[X_2^3 - 3X_2|X_1, M = (0, 1)] \\ &= \mathbb{E}\left[(X_1 + \epsilon)^3 - 3(X_1 + \epsilon)|X_1, M = (0, 1)\right] \\ &= \mathbb{E}\left[X_1^3 + \epsilon^3 + 3X_1\epsilon^2 + 3X_1^2\epsilon - 3X_1 - 3\epsilon|X_1, M = (0, 1)\right] \\ &= X_1^3 + 3X_1(\sigma^2 - 1) \end{aligned}$$

Clearly, the Bayes predictor for $M = (0, 1)$ is:

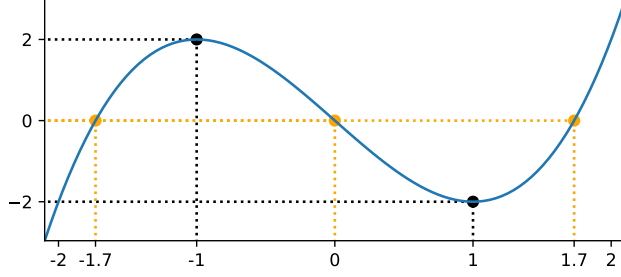


Figure 5: Graph of $X_2 \mapsto f^*(X_1, X_2)$

- continuous,
- non-decreasing since $\sigma^2 > 1$,
- $\lim_{X_1 \rightarrow +\infty} \tilde{f}^*(X_1, M = (0, 1)) = +\infty$ and $\lim_{X_1 \rightarrow -\infty} \tilde{f}^*(X_1, M = (0, 1)) = -\infty$.

Proof by contradiction: Suppose that there exists a function $\Phi : \mathbb{R} \rightarrow \mathbb{R}$ (i) continuous and (ii) such that for all X_1 , $f^*(X_1, \Phi(X_1)) = \tilde{f}^*(X_1, M = (0, 1))$.

Let $x_1^+ \in \mathbb{R}$ such that $\tilde{f}^*(X_1 = x_1^+, M = (0, 1)) > 2$. x_1^+ exists since $\lim_{X_1 \rightarrow +\infty} \tilde{f}^*(X_1, M = (0, 1)) = +\infty$. Clearly,

$$f^*(x_1^+, X_2) = \tilde{f}^*(x_1^+, M = (0, 1)) \iff X_2 = x_2^+ \quad \text{with} \quad x_2^+ > 2.$$

Similarly, let $x_1^- \in \mathbb{R}$ such that $\tilde{f}^*(X_1 = x_1^-, M = (0, 1)) < -2$. x_1^- exists since $\lim_{X_1 \rightarrow -\infty} \tilde{f}^*(X_1, M = (0, 1)) = -\infty$. Clearly,

$$f^*(x_1^-, X_2) = \tilde{f}^*(x_1^-, M = (0, 1)) \iff X_2 = x_2^- \quad \text{with} \quad x_2^- < -2.$$

So Φ must satisfy:

$$\Phi(x_1^-) = x_2^- < -2$$

$$\Phi(x_1^+) = x_2^+ > 2$$

Note that since the Bayes predictor is non-decreasing, we have $x_1^- < x_1^+$. Since Φ is continuous, there exists $\tilde{x}_1 \in [x_1^-, x_1^+]$ and $\hat{x}_1 \in [x_1^-, x_1^+]$ such that $\tilde{x}_1 < \hat{x}_1$ and $\Phi(\tilde{x}_1) = -1$ and $\Phi(\hat{x}_1) = 1$. It implies that:

$$f^*(\tilde{x}_1, \Phi(\tilde{x}_1)) = f^*(\tilde{x}_1, -1) = 2 > -2 = f^*(\hat{x}_1, 1) = f^*(\hat{x}_1, \Phi(\hat{x}_1)).$$

This implies that the function $X_1 \mapsto f^*(X_1, \Phi(X_1))$ cannot be non-decreasing. Since the Bayes predictor is non-decreasing, the two cannot be equal. CONTRADICTION.

A.9 Proof of Proposition 4.3

We start by proving the result for a given missing pattern $m \in \{0, 1\}^d$. Take $r \in \{1, \dots, d-1\}$ and consider a missing pattern m such that $|\text{obs}(m)| = r$. We let $F : \mathbb{R}^r \times \mathbb{R}^{d-r} \rightarrow \mathbb{R}$ defined, for all $(x_{\text{obs}}, x_{\text{mis}})$ as

$$F(x_{\text{obs}}, x_{\text{mis}}) = f^*(x_{\text{obs}}, x_{\text{mis}}) - \tilde{f}^*(x_{\text{obs}}, m). \quad (28)$$

Our aim is to find, for all x_{obs} , a value x_{mis} (depending continuously on x_{obs}) satisfying

$$F(x_{\text{obs}}, x_{\text{mis}}) = 0. \quad (29)$$

To this aim, we check the assumptions of Theorem 6 in Arutyunov and Zhukovskiy [2019] for the function F . The desired conclusion will follow.

Since f^* is uniformly continuous and twice continuously differentiable, condition 1 – 3 of Theorem 6 in Arutyunov and Zhukovskiy [2019] are satisfied. To verify the next condition, we have to prove that

there exists $(x_{obs,0}, x_{mis,0})$ such that $F(x_{obs,0}, x_{mis,0}) = 0$. Note that this is equivalent to finding $(x_{obs,0}, x_{mis,0})$ satisfying

$$f^*(x_{obs,0}, x_{mis,0}) = \tilde{f}^*(x_{obs,0}, m) = \mathbb{E}[f^*(X)|X_{obs} = x_{obs,0}, M = m], \quad (30)$$

by definition of the regression function \tilde{f}^* . By assumption, the support of $X_{mis}|X_{obs} = x_{obs,0}, M = m$ is connected. Therefore, the intermediate value theorem can be applied and proves the existence of a pair $(x_{obs,0}, x_{mis,0})$ satisfying equation (30). Finally, by assumption, the regularity condition (GR1) in Arutyunov and Zhukovskiy [2019] is satisfied. This proves that there exists a continuous mapping $\phi^{(m)} : \mathbb{R}^r \rightarrow \mathbb{R}^{d-r}$ such that

$$F(x_{obs}, \phi^{(m)}(x_{obs})) = 0. \quad (31)$$

The previous reasoning holds for all missing patterns m , such that $|mis(m)| \geq 1$. Besides the result is clear for $r = 0$ since the imputation function is reduced to a constant in this case (no components of X are observed). On the contrary, in the case where all covariates are observed ($r = d$), no imputation function is needed. Therefore, the result holds for all $0 \leq r \leq d$, which concludes the proof.

A.10 Expressions of f_{bowl}^* and f_{wave}^* and their corresponding Bayes predictors.

Expressions of f_{bowl}^* and f_{wave}^* . The functions f^* used in the experimental study are defined as:

$$\begin{aligned} f_{bowl}^*(X) &= (\beta^\top X + \beta_0 - 1)^2 \\ f_{wave}^*(X) &= (\beta^\top X + \beta_0 - 1) + \sum_{(a_i, b_i) \in S} a_i \Phi(\gamma(\beta^\top X + \beta_0 + b_i)) \end{aligned}$$

where Φ the standard Gaussian cdf, $\gamma = 20\sqrt{\frac{\pi}{8}}$ and $S = \{(2, -0.8), (-4, -1), (2, -1.2)\}$. β is chosen as a vector of ones rescaled so that $\text{var}(\beta^\top X) = 1$. These functions are depicted in Figure 3.

Expressions of the Bayes predictors. The expressions of the Bayes predictors corresponding to f_{bowl}^* and f_{wave}^* are given by:

$$\tilde{f}_{bowl}^*(\tilde{X}) = \mathbb{E}[f_{bowl}^*(X)|X_{obs}, M] \quad (32)$$

$$= (\beta_{obs}^\top X_{obs} + \beta_{mis}^\top \mu_{mis|obs, M} + \beta_0 - 1)^2 + \beta_{mis}^\top \Sigma_{mis|obs, M} \beta_{mis} \quad (33)$$

and:

$$\tilde{f}_{wave}^*(\tilde{X}) = \mathbb{E}[f_{wave}^*(X)|X_{obs}, M] \quad (34)$$

$$= \beta_{obs}^\top X_{obs} + \beta_{mis}^\top \mu_{mis|obs, M} + \beta_0 - 1 \quad (35)$$

$$+ \sum_{(a_i, b_i) \in S} a_i \Phi\left(\frac{\beta_{obs}^\top X_{obs} + \beta_{mis}^\top \mu_{mis|obs, M} + \beta_0 + b_i}{\sqrt{1/\gamma^2 + \beta_{mis}^\top \Sigma_{mis|obs, M} \beta_{mis}}}\right) \quad (36)$$

with $\mu_{mis|obs, M}$ and $\Sigma_{mis|obs, M}$ the mean and covariance matrix of the conditional distribution $P(X_{mis}|X_{obs}, M)$. Below, we give the expression of these parameters for the MCAR and Gaussian self-masking missing data mechanisms. Let $\mu_{mis|obs}$ and $\Sigma_{mis|obs}$ the mean and covariance matrix of the conditional distribution $P(X_{mis}|X_{obs})$. Since the data is generated according to a multivariate Gaussian distribution $\mathcal{N}(\mu, \Sigma)$, we have:

$$\mu_{mis|obs} = \mu_{mis} + \Sigma_{mis|obs} \Sigma_{obs}^{-1} (X_{obs} - \mu_{obs})$$

$$\Sigma_{mis|obs} = \Sigma_{mis, mis} - \Sigma_{mis, obs} \Sigma_{obs}^{-1} \Sigma_{obs, mis}$$

In the MCAR case, we simply have $\Sigma_{mis|obs, M} = \Sigma_{mis|obs}$ and $\mu_{mis|obs, M} = \mu_{mis|obs}$. In the Gaussian self-masking case, it has been shown in Le Morvan et al. [2020a] that $P(X_{mis}|X_{obs}, M)$ is again Gaussian but with parameters:

$$\Sigma_{mis|obs, M} = \left(D_{mis, mis}^{-1} + \Sigma_{mis|obs}^{-1} \right)^{-1}$$

$$\mu_{mis|obs, M} = \Sigma_{mis|obs, M} \left(D_{mis, mis}^{-1} \tilde{\mu}_{mis} + \Sigma_{mis|obs}^{-1} \mu_{mis|obs} \right)$$

where $\tilde{\mu}$ and D are parameters of the Gaussian self-masking missing data mechanism. Finally, we give below the derivations to obtain the expression of the Bayes predictors.

Derivation of the Bayes predictor for f_{bowl}^* .

$$f_{bowl}^*(X) = (\beta^\top X + \beta_0 - 1)^2 \quad (37)$$

$$= (\beta_{obs}^\top X_{obs} + \beta_{mis}^\top X_{mis} + \beta_0 - 1)^2 \quad (38)$$

$$= (\beta_{obs}^\top X_{obs} + \beta_{mis}^\top (X_{mis} - \mu_{mis|obs,M}) + \beta_{mis}^\top \mu_{mis|obs,M} + \beta_0 - 1)^2 \quad (39)$$

$$= (\beta_{obs}^\top X_{obs} + \beta_{mis}^\top \mu_{mis|obs,M} + \beta_0 - 1)^2 + (\beta_{mis}^\top (X_{mis} - \mu_{mis|obs,M}))^2 \quad (40)$$

$$+ 2\beta_{mis}^\top (X_{mis} - \mu_{mis|obs,M}) (\beta_{obs}^\top X_{obs} + \beta_{mis}^\top \mu_{mis|obs,M} + \beta_0 - 1) \quad (41)$$

Now taking the expectation with regards to $P(X_{mis}|X_{obs}, M)$, the last term vanishes and we get:

$$\mathbb{E}[f_{bowl}^*(X)|X_{obs}, M] = (\beta_{obs}^\top X_{obs} + \beta_{mis}^\top \mu_{mis|obs,M} + \beta_0 - 1)^2 + \beta_{mis}^\top \Sigma_{mis|obs,M} \beta_{mis} \quad (42)$$

Derivation of the Bayes predictor for f_{wave}^* .

$$f_{wave}^*(X) = (\beta^\top X + \beta_0 - 1) + \sum_{(a_i, b_i) \in S} a_i \Phi(\gamma(\beta^\top X + \beta_0 + b_i)) \quad (43)$$

$$= (\beta_{obs}^\top X_{obs} + \beta_{mis}^\top X_{mis} + \beta_0 - 1) \quad (44)$$

$$+ \sum_{(a_i, b_i) \in S} a_i \Phi(\gamma(\beta_{obs}^\top X_{obs} + \beta_{mis}^\top X_{mis} + \beta_0 + b_i)) \quad (45)$$

Define $T^{(m)} = \beta_{mis}^\top X_{mis}$. Since $P(X_{mis}|X_{obs}, M)$ is Gaussian in both the MCAR and Gaussian self-masking cases, $P(T^{(m)}|X_{obs}, M)$ is also Gaussian with mean and variance given by:

$$\mu_{T^{(m)}|X_{obs}, M} = \beta_{mis}^\top \mu_{mis|obs,M} \quad (46)$$

$$\sigma_{T^{(m)}|X_{obs}, M}^2 = \beta_{mis}^\top \Sigma_{mis|obs,M} \beta_{mis} \quad (47)$$

To compute the Bayes predictor, we now need to compute the quantity:

$$\mathbb{E}_{T^{(m)}|X_{obs}, M} \left[\Phi \left(\gamma \left(\beta_{obs}^\top X_{obs} + T^{(m)} + \beta_0 + b_i \right) \right) \right]$$

This expectation can then be computed following [Bishop, 2006] (section 4.5.2) which gives the result.

A.11 NeuMiss architecture

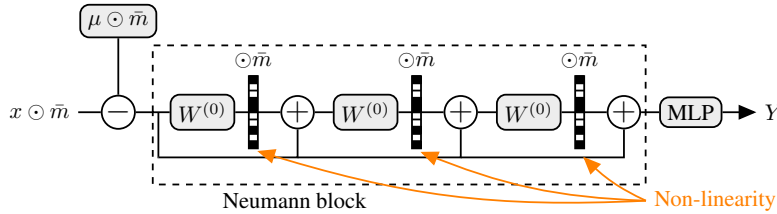


Figure 6: **(Non-linear) NeuMiss network architecture with a Neumann block of depth 3** — $\bar{m} = 1 - m$. MLP stands for a standard multi-layer perceptron with ReLU activations.

Initialization strategy for the Neumann block [Le Morvan et al., 2020a] designed the Neumann block so that the network may be able to approximate well the conditional expectation $\mathbb{E}[X_{mis}|X_{obs}, M]$. In particular, if the weights of the Neumann block are chosen equal to certain quantities depending on μ and Σ (the mean and covariance matrix of the data distribution), then the network computes an approximation of the conditional expectation whose error tends to zero exponentially fast with the depth of the Neumann block. Since μ and Σ are unknown in practice, the weights of the network cannot be set to the desired quantities and are instead learned by gradient descent.

In our experiments, we computed estimates $\hat{\mu}$ and $\hat{\Sigma}$ as the sample mean and sample covariance matrix over the observed entries only, i.e, if variable j is observed in n_j samples, then $\hat{\mu}_j$ is computed

as the mean over these n_j values only. Similarly, the covariance between variables i and j is estimated using only the samples for which both are observed. Based on these estimates, we initialized the weights of the network to their theoretical values given in Le Morvan et al. [2020a]:

$$\begin{aligned}\mu &= \hat{\mu} \\ W^{(0)} &= Id - \frac{2}{\hat{L}} \hat{\Sigma}\end{aligned}$$

where \hat{L} is the largest eigenvalue of Σ .

We use a pytorch-based [Paszke et al., 2019] implementation of NeuMiss as well as the MLP with which we combine it.