

# **Conformal Prediction with Missing Values**

Margaux Zaffran, Aymeric Dieuleveut, Julie Josse, Yaniv Romano

### ▶ To cite this version:

Margaux Zaffran, Aymeric Dieuleveut, Julie Josse, Yaniv Romano. Conformal Prediction with Missing Values. 2023. hal-03896384v2

# HAL Id: hal-03896384 https://hal.science/hal-03896384v2

Preprint submitted on 7 Mar 2023 (v2), last revised 9 Nov 2023 (v4)

**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.

## Conformal Prediction with Missing Values

Margaux Zaffran<sup>\*,1,2,3</sup>, Aymeric Dieuleveut<sup>3</sup>, Julie Josse<sup>2</sup>, and Yaniv Romano<sup>4</sup>

<sup>1</sup>Electricité De France R&D, Palaiseau, France

 <sup>2</sup>PreMeDICaL project team, INRIA Sophia-Antipolis, Montpellier, France
 <sup>3</sup>CMAP, CNRS, École polytechnique, Institut Polytechnique de Paris, Palaiseau, France
 <sup>4</sup>Departments of Electrical Engineering and of Computer Science, Technion - Israel Institute of Technology, Haifa, Israel

#### Abstract

Conformal prediction is a theoretically grounded framework for constructing predictive intervals. We study conformal prediction with missing values in the covariates - a setting that brings new challenges to uncertainty quantification. We first show that the marginal coverage guarantee of conformal prediction holds on imputed data for any missingness distribution and almost all imputation functions. However, we emphasize that the average coverage varies depending on the pattern of missing values: conformal methods tend to construct prediction intervals that under-cover the response conditionally to some missing patterns. This motivates our novel generalized conformalized quantile regression framework, missing data augmentation, which yields prediction intervals that are valid conditionally to the patterns of missing values, despite their exponential number. We then show that a universally consistent quantile regression algorithm trained on the imputed data is Bayes optimal for the pinball risk, thus achieving valid coverage conditionally to any given data point. Moreover, we examine the case of a linear model, which demonstrates the importance of our proposal in overcoming the heteroskedasticity induced by missing values. Using synthetic and data from critical care, we corroborate our theory and report improved performance of our methods.

#### 1 Introduction

By leveraging increasingly large data sets, statistical algorithms and machine learning methods can be used to support high-stakes decision-making problems such as autonomous driving, medical or civic applications, and more. To ensure the safe deployment of predictive models it is crucial to quantify the uncertainty of the resulting predictions, communicating the limits of predictive performance. Uncertainty quantification attracts a lot of attention in recent years, particularly methods that are based on Conformal Prediction (CP) (Vovk et al., 2005; Papadopoulos et al., 2002; Lei et al., 2018). CP provides controlled predictive regions for any underlying predictive algorithm (e.g., neural networks and random forests), in finite samples with no assumption on the data distribution except for the exchangeability of the train and test data. More precisely, for a *miscoverage rate*  $\alpha \in [0, 1]$ , CP outputs a prediction interval  $\hat{C}_{\alpha}$  for the test response Y given its corresponding covariates X such that:

$$\mathbb{P}(Y \in C_{\alpha}(X)) \ge 1 - \alpha.$$
(1)

Split CP (Papadopoulos et al., 2002; Lei et al., 2018) achieves Eq. (1) by keeping a hold-out set, the *calibration set*, used to evaluate the performance of a fixed predictive model.

At the same time, as the volume of data increases, the volume of missing values also increases. There is a vast literature on this topic (Little, 2019; Josse & Reiter, 2018), and a recent survey even identified more than 150 different implementations (Mayer et al., 2019). Missing values create additional challenges to the task of supervised learning, as traditional machine learning algorithms can not handle incomplete data (Josse et al., 2019; Le Morvan et al., 2020b,a, 2021; Ayme et al., 2022; Van Ness et al., 2022). One of the most popular strategies to deal with missing values suggests imputing the missing entries with plausible values to get completed data, on which any analysis can be performed. The drawback of this "impute-then-predict" approach is that single imputation can distort the joint and marginal distribution of the data. Yet, Josse et al. (2019); Le Morvan et al. (2020b, 2021) showed that such impute-then-predict strategies are Bayes consistent, under the assumption that a universally consistent learner is applied on an imputed data set. However, this line of work focuses on point prediction with missing values that aim to predict the most likely outcome. In contrast, our goal is quantifying predictive uncertainty, which surprisingly was not explored with missing values although its enormous importance.

#### Contributions

We study CP in the presence of missing values in the covariates. More precisely, we study downstream quantileregression (QR) based CP, such as CQR (Romano et al., 2019), on impute-then-predict strategies.

<sup>\*</sup>Corresponding author: margaux.zaffran@inria.fr

After recalling background and notations in Section 2, our first contribution is showing that CP on impute-then-predict is *marginally* valid regardless of the model, missingness distribution, and imputation function (Section 3).

Then, we consider valid coverage *conditionally on the missing data pattern*, referred to as a *mask*. In Section 4, we describe how different masks introduce additional heteroskedasticity: *the uncertainty on the output strongly depends on the set of predictive features observed*. We illustrate the need for a novel method to handle this phenomenon on synthetic data in Figure 1 (left panel): CQR (orange crosses) does not reach the target value of 90% on the mask with the lowest coverage (worst group).

This motivates our second contribution: we show in Section 5 how to form prediction intervals that are valid conditional on any given mask. This is highly challenging since there are exponentially many possible patterns to consider. Therefore, the naive solution to perform a calibration for each possible mask would fail as in finite samples, we often observe test samples with missing patterns that have low (or even null) frequency of appearance in the calibration set. To tackle this issue, we suggest two conformal methods that share the same core idea of missing data augmentation (MDA): we artificially mask the calibration data to match the mask of the point we consider at test time. The first method, CP-MDA with exact masking, relies on building an ideal calibration set for which the data points have the exact same mask as of the test point. We show its validity under exchangeability and Missing Completely At Random assumptions. Our second method, CP-MDA with nested masking, does not require such an ideal calibration set. Instead, we artificially construct a calibration set in which the data points have at least the same mask as the test point, i.e., this artificial masking results in calibration points having possibly more missing values than the test point. We show the latter methodology also achieves the desired coverage conditional on the mask, but at the cost of an additional assumption for validity: stochastic domination of the quantiles. Figure 1 illustrates those findings: both methods are valid even for small training set-size.

Our third contribution further supports our design choice to use QR. We show that QR on impute-then-predict strategy is Bayes-consistent – it can achieve the strongest form of coverage conditional on the observed test features (Section 6).

Lastly, we support our proposal using both (semi)-synthetic experiments and real medical data (Section 7). The code to reproduce our experiments is made available.

#### 2 Background

**Background on missing values.** Consider a data set with n exchangeable realizations of the random variable  $(X, M, Y) \in \mathcal{X} \times \{0, 1\}^d \times \mathbb{R}: \{(X^{(k)}, M^{(k)}, Y^{(k)})\}_{k=1}^n$ , where X represents the features, M the missing pattern, or mask, and Y an outcome to predict. For  $j \in [\![1, d]\!], M_j = 0$  when  $X_j$  is observed and  $M_j = 1$  when  $X_j$  is missing, i.e. NA (Not Available). We note  $\mathcal{M} = \{0, 1\}^d$  the set of masks. For a pattern  $m \in \mathcal{M}, X_{\text{obs}(m)}$  is the random vector of observed components, and  $X_{\min(m)}$  is the random vector of unobserved ones. For example, if we observe (NA, 6, 2) then m = (1, 0, 0) and  $X_{\text{obs}(m)} = (6, 2)$ . Our goal is to predict a new outcome  $Y^{(n+1)}$  given  $X_{\text{obs}(M^{(n+1)})}^{(n+1)}$  and  $M^{(n+1)}$ .

Assumption A1 (exchangeability). The random variables  $(X^{(k)}, M^{(k)}, Y^{(k)})_{k=1}^{n+1}$  are exchangeable.

Following Rubin (1976), we consider three well-known missingness mechanisms.

**Definition 2.1** (Missing Completely At Random (MCAR)). For any  $m \in \mathcal{M}$ ,  $\mathbb{P}(M = m | X) = \mathbb{P}(M = m)$ .

**Definition 2.2** (Missing At Random (MAR)). For any  $m \in \mathcal{M}$ ,  $\mathbb{P}(M = m|X) = \mathbb{P}(M = m|X_{obs(m)})$ .

**Definition 2.3** (Missing Non At Random (MNAR)). If the missing data is not MAR, it is MNAR. Thus, its probability



Figure 1: Coverage and length of the predictive intervals as a function of the training size, where the performance metrics correspond to the pattern of missing values giving the lowest coverage. Methods are Quantile Regression (QR), Conformalized Quantile Regression (CQR), and two missing data augmentation procedures, CP-MDA-Exact and CP-MDA-Nested, on top of CQR, using a Quantile Neural Network. Settings are given in Section 7. Results evaluated over 100 repetitions, and the (tiny) error bars correspond to standard error.

distribution depends on X, including the missing values.

As most predictive algorithms can not directly handle missing values, we impute the incomplete data using an imputation function  $\Phi$  which maps observed values to themselves and missing values to a function of the observed values. With notation from Le Morvan et al. (2021),  $\Phi$  belongs to  $\mathcal{F}^{I} := \{ \Phi : \mathcal{X} \times \mathcal{M} \to \mathcal{X} : \forall j \in [\![1,d]\!],$ 

 $\Phi_j(X,M) = X_j \mathbb{1}_{M_j=0} + \phi_j^M \left( X_{\text{obs}(M)} \right) \mathbb{1}_{M_j=1} \right\}.$ Additionally,  $\mathcal{F}_{\infty}^I$  is the restriction of  $\mathcal{F}^I$  to  $\mathcal{C}^{\infty}$  functions which include deterministic imputation, such as mean imputation or imputation by regression. The imputed data set is formed by the realizations of the *n* random variables  $(\Phi(X,M), M, Y)$ . In practice,  $\Phi$  is obtained as the result of an algorithm  $\mathcal{I}$  trained on  $\{(X^{(k)}, M^{(k)})\}_{k=1}^{n+1}$ .

Assumption A2 (Symmetrical imputation). The imputation function  $\Phi$  is the output of an algorithm  $\mathcal{I}$  treating its input data points symmetrically:  $\mathcal{I}((X^{(\sigma(k))}, M^{(\sigma(k))})_{k=1}^{n+1}) \stackrel{(d)}{=} \mathcal{I}((X^{(k)}, M^{(k)})_{k=1}^{n+1})$  conditionally on  $(X^{(k)}, M^{(k)})_{k=1}^{n+1}$  and for any permutation  $\sigma$  on [1, n+1].

Assumption A2 is very mild and satisfied by all existing imputation methods for exchangeable data. In particular, it is valid for iterative regression imputation which allows out-of-sample imputation.

Background on (split) conformal prediction. Split, or inductive, CP (SCP) (Papadopoulos et al., 2002; Lei et al., 2018) builds predictive regions by first splitting the n points of the training set into two disjoint sets  $Tr, Cal \subset [\![1, n]\!]$ , to create a proper training set, Tr, and a calibration set, Cal. On the proper training set, a model f (chosen by the user) is fitted, and then used to predict on the calibration set. Conformity scores  $S_{\text{Cal}} = \{(s(X^{(k)}, Y^{(k)}))_{k \in \text{Cal}}\}$  are computed to assess how well the fitted model  $\hat{f}$  predicts the response values of the calibration points. For example, Conformalized Quantile Regression (CQR, Romano et al., 2019) fits two quantile regressions  $\hat{q}_{low}$  and  $\hat{q}_{upp}$ , on the proper training set. The conformity scores are defined by  $s(x, y) = \max(\hat{q}_{low}(x) - y, y - \hat{q}_{upp}(x))$ . Finally, a corrected  $(1 - \tilde{\alpha})$ -th quantile of these scores  $\widehat{Q}_{1-\tilde{\alpha}}(S_{\text{Cal}})$  is computed (called correction term) to define the predictive region:  $\widehat{C}_{\alpha}(x) := \{y \text{ such that } s(y, \widehat{f}(x)) \leq \widehat{Q}_{1-\tilde{\alpha}}(S_{\operatorname{Cal}})\}.^1$ An illustration of CQR is provided in Appendix A.

This procedure satisfies Eq. (1) for any  $\hat{f}$ , any (finite) sample size n, as long as the data points are exchangeable.<sup>2</sup> Moreover, if the scores are almost surely distinct, the coverage holds almost exactly:  $\mathbb{P}(Y \in \hat{C}_{\alpha}(X)) \leq 1 - \alpha + \frac{1}{\#\text{Cal}+1}$ .

For more details on SCP, we refer to Angelopoulos & Bates (2021); Vovk et al. (2005), as well as to Manokhin (2022).

#### **3** Warm-up: marginal coverage with NAs

A first idea to get valid predictive intervals  $\widehat{C}_{\alpha}(X, M)$  in the presence of missing values M is to apply CP in combination with impute-then-predict, which we refer to as *impute-then-predict+conformalization*. More details on this approach are given in Appendix B.1 for both classification and regression tasks, although our main focus is regression. It turns out that such a simple approach is marginally (exactly) valid.

**Definition 3.1** (Marginal validity). A method outputting intervals  $\hat{C}_{\alpha}$  is marginally valid if the lower bound is satisfied, and exactly valid if the upper bound is also satisfied:

$$1 - \alpha \leq_{\text{validity}} \mathbb{P}\left(Y^{(n+1)} \in \widehat{C}_{\alpha}\left(X^{(n+1)}, M^{(n+1)}\right)\right)$$
$$\leq_{\text{exact validity}} 1 - \alpha + \frac{1}{\#\text{Cal} + 1}.$$

Indeed, symmetric imputation preserves exchangeability.

**Lemma 3.2** (Imputation preserves exchangeability). Let A1 hold. Then, for any missing mechanism, for any imputation function  $\Phi$  satisfying A2, the imputed random variables  $\left(\Phi\left(X^{(k)}, M^{(k)}\right), M^{(k)}, Y^{(k)}\right)_{k=1}^{n+1}$  are exchangeable.

Note that if we replace A1 by an i.i.d. assumption, the imputed data set remains only exchangeable without further assumptions on  $\mathcal{I}$ : indeed, imputing by the empirical mean for example breaks the independence.

**Proposition 3.3** ((Exact) validity of impute-then-predict+conformalization). *If A1 and A2 are satisfied, imputethen-predict+conformalization is marginally valid. If moreover the scores are almost surely distinct, it is exactly valid.* 

This is an important first positive result (proved in Appendix B.2) showing that CP applied on an imputed data set has the same validity properties as on complete data, regardless of the missing value mechanism (MCAR, MAR or MNAR) and of the symmetric imputation scheme. Note that similar propositions could be derived for full CP (Vovk et al., 2005) and Jackknife+ (Barber et al., 2021b).

#### 4 Challenge: NAs induce heteroskedasticity

To better understand the interplay between missing values and conditional coverage with respect to the mask, we consider an illustrative example of a Gaussian linear model.

**Model 4.1** (Gaussian linear model). The data is generated according to a linear model and the covariates are Gaussian conditionally to the pattern:

- $Y = \beta^T X + \varepsilon, \varepsilon \sim \mathcal{N}(0, \sigma_{\varepsilon}^2) \perp (X, M), \beta \in \mathbb{R}^d.$
- for all  $m \in \mathcal{M}$ , there exist  $\mu^m$  and  $\Sigma^m$  such that  $X|(M = m) \sim \mathcal{N}(\mu^m, \Sigma^m)$ .

In particular, Model 4.1 is verified when X is Gaussian and the missing data is MCAR. Model 4.1 is more general: it even includes MNAR examples (Ayme et al., 2022).

<sup>&</sup>lt;sup>1</sup>The correction  $\alpha \rightarrow \tilde{\alpha}$  is needed because of the inflation of quantiles in finite sample (see Lemma 2 in Romano et al. (2019) or Section 2 in Lei et al. (2018)).

<sup>&</sup>lt;sup>2</sup>Only the calibration and test data points need to be exchangeable.

**Proposition 4.2** (Oracle intervals). The oracle predictive interval is defined as the smallest valid interval knowing  $X_{\text{obs}(M)}$  and M. Under Model 4.1, its length only depends on the mask. For any  $m \in M$  this oracle length is:

$$\mathcal{L}^*_{\alpha}(m) = 2q_{1-\frac{\alpha}{2}}^{\mathcal{N}(0,1)} \sqrt{\beta_{\min(m)}^T \Sigma_{\min|obs}^m \beta_{\min(m)} + \sigma_{\varepsilon}^2}.$$
 (2)

See Appendix C for the definition of  $\mu_{\text{mis|obs}}^m$  and  $\Sigma_{\text{mis|obs}}^m$ and the quantiles of  $Y|(X_{\text{obs}(m)}, M = m)$ .

Eq. (2) stresses that even when the noise of the generative model is homoskedastic, *missing values induce heteroskedasticity*. Indeed, the covariance of the conditional distribution of  $Y|(X_{obs(m)}, M = m)$  depends on m. Furthermore, the uncertainty increases when missing values are associated with larger regression coefficients (i.e. the most predictive variables): if  $\beta_{mis(m)}$  is large, then  $\mathcal{L}^*_{\alpha}(m)$  is also large, as  $\Sigma^m_{mis|obs}$  is positive. In the extreme case where all the variables are missing, i.e.  $m = (1, \dots, 1)$ ,  $\mathcal{L}^*_{\alpha}(m) = 2q_{1-\frac{\alpha}{2}}^{\mathcal{N}(0,1)}\sqrt{\beta\Sigma^m\beta^T + \sigma_{\varepsilon}^2} = q_{1-\frac{\alpha}{2}}^{\mathcal{P}} - q_{\frac{\alpha}{2}}^{\mathcal{Y}}$ . On the contrary, if  $m = (0, \dots, 0)$  (that is all  $X_j$  are observed),  $\beta_{mis(m)}$  is empty and  $\mathcal{L}^*_{\alpha}(m) = 2q_{1-\frac{\alpha}{2}}^{\mathcal{N}(0,1)}\sigma_{\varepsilon} = q_{1-\frac{\alpha}{2}}^{\varepsilon} - q_{\frac{\alpha}{2}}^{\varepsilon}$ . We illustrate this induced heteroskedasticity and the impact of the predictive power in Appendix C, along with a discussion emphasizing that even with the Bayes predictor for the conditional mean, mean-based CP does not yield intervals that are valid conditionally on the mask.

The above analysis motivates the following two design choices we make in this work. First, we advocate working with OR models rather than classic regression ones, as the former can handle heteroskedastic data. Second, we recommend providing the mask information to the model in addition to the input covariates, as the mask may further encourage the model to construct an interval with a length adaptive to the given mask. Therefore, we focus on CQR (Romano et al., 2019), an adaptive version of SCP, and concatenate the mask to the features. However, the predictive intervals of this procedure may not necessarily provide valid coverage conditionally on the masks, especially in finite samples as shown in Figure 1 (orange crosses). This is due to the exponential number of patterns  $(2^d)$ , which strongly varies the quality of the finite sample predictions obtained according to the mask, whereas the correction term is calculated independently of the masks.

#### 5 Achieving mask-conditional-validity

We now aim at achieving *mask-conditional-validity* defined as follows using an ordering on the masks.

**Definition 5.1** (Included masks). Let  $(\mathring{m}, \breve{m}) \in \mathcal{M}^2, \mathring{m} \subset \breve{m}$  if for any  $j \in [\![1,d]\!]$  such that  $\mathring{m}_j = 1$  then  $\breve{m}_j = 1$ , i.e.  $\breve{m}$  includes at least the same missing values than  $\mathring{m}$ .

**Definition 5.2** (Mask-conditional-validity). A method is mask-conditionally-valid if for any  $m \in \mathcal{M}$  the lower bound is satisfied, and exactly mask-conditionally-valid if for any

 $m \in \mathcal{M}$  the upper bound is also satisfied:

$$\begin{aligned} 1 - \alpha &\leq \mathop{\mathbb{P}}_{\text{valid}} \left( Y^{(n+1)} \in \widehat{C}_{\alpha} \left( X^{(n+1)}, m \right) | M^{(n+1)} = m \right) \\ &\leq \\ &\leq \\ & \text{exactly valid} \quad 1 - \alpha + \frac{1}{\# \operatorname{Cal}^{\mathrm{m}} + 1}, \end{aligned}$$

where  $\operatorname{Cal}^m = \{ k \in \operatorname{Cal} \text{ such that } m^{(k)} \subset m \}.$ 

#### 5.1 Missing Data Augmentation (MDA)

To obtain a mask-conditionally-valid procedure, we suggest modifying the calibration set according to the mask of the test point, while the training step is unchanged. More precisely, the mask of the test point is applied to the calibration set, as illustrated in Figure 2. The rationale is to mimic the missing pattern of the test point by artificially augmenting the calibration set with that mask. It ensures that the correction term is computed using data with (at least) the same missing values as the test point. We refer to this strategy as *CP with Missing Data Augmentation* (CP-MDA), and derive two versions of it. Algorithms 1 and 2 are written using CQR as the base conformal procedure, but they work with any conformal method as we describe in Appendix D.1.

Algorithm 1 – CP-MDA-Exact. CP-MDA with *exact masking* consists of keeping the *artificially* masked calibration points (1. 7) that have exactly the same missing pattern as the test point (1. 5). Then Algorithm 1 performs as imputethen-predict+conformalization: impute the calibration set (1. 10), predict on it and get the calibration scores (1. 11), compute their quantile to obtain the correction term (1. 14), and finally impute and predict the test point with the fixed fitted model by adding and subtracting the correction term (1. 15) to the initial conditional quantile estimates. Note that Algorithm 1 is described for one test point for simplicity but extends easily to many test points. The computations are



Figure 2: CP-MDA illustration. Augmented calibration set according to one test point. For CP-MDA-Nested, the augmented masks of the calibration set are also applied temporarily to the test point.

#### Algorithm 1 CP-MDA-Exact (with COR)

**Input:** Imputation algorithm  $\mathcal{I}$ , quantile regression algorithm QR, significance level  $\alpha$ , training set  $\{(x^{(k)}, m^{(k)}, y^{(k)})\}_{k=1}^{n}$ , test point  $(x^{(\text{test})}, m^{(\text{test})})$ 

**Output:** Prediction interval  $\widehat{C}_{\alpha}(x^{(\text{test})}, m^{(\text{test})})$ 

- 1: Randomly split  $\{1, \ldots, n\}$  into 2 disjoint sets Tr & Cal Fit the imputation  $\Phi(\cdot) \leftarrow \mathcal{I}\left(\left\{\left(x^{(k)}, m^{(k)}\right), k \in \mathrm{Tr}\right\}\right)$ 2: Fit function:
- 3: Impute the training set:  $\forall k \in \text{Tr}, x_{\text{imp}}^{(k)} = \Phi(x^{(k)}, m^{(k)})$ 4: Fit QR:

$$\hat{q}_{\frac{\alpha}{2}}(\cdot) \leftarrow \mathcal{QR}\left(\left\{\left(x_{\text{imp}}^{(k)}, y^{(k)}\right), k \in \text{Tr}\right\}, \alpha/2\right)$$
$$\hat{q}_{1-\frac{\alpha}{2}}(\cdot) \leftarrow \mathcal{QR}\left(\left\{\left(x_{\text{imp}}^{(k)}, y^{(k)}\right), k \in \text{Tr}\right\}, 1-\alpha/2\right)$$

// Generate an augmented calibration set:

- 5:  $\operatorname{Cal}^{(\text{test})} = \{ k \in \operatorname{Cal} \text{ such that } m^{(k)} \subset m^{(\text{test})} \}$
- 6: for  $k \in \operatorname{Cal}^{(\operatorname{test})}$  do
- $\widetilde{m}^{(k)} = m^{(\text{test})}$  //Additional masking 7:
- 8: end for Augmented calibration set generated. //
- 9: for  $k \in \operatorname{Cal}^{(\text{test})}$  do

10: Impute the calibration set: 
$$x_{imp}^{(k)} = \Phi(x^{(k)}, \widetilde{m}^{(k)})$$

11: Set 
$$s^{(k)} = \max(\hat{q}_{\frac{\alpha}{2}}(x_{\text{imp}}^{(k)}) - y^{(k)}, y^{(k)} - \hat{q}_{1-\frac{\alpha}{2}}(x_{\text{imp}}^{(k)}))$$

- 12: end for
- 13: Set  $S = \{s^{(k)}, k \in \text{Cal}^{(\text{test})}\}$
- 14: Compute  $\widehat{Q}_{1-\tilde{\alpha}}(S)$ , the  $1-\tilde{\alpha}$ -th empirical quantile of S, with  $1 - \tilde{\alpha} := (1 - \alpha) (1 + 1/\#S)$

15: Set 
$$\widehat{C}_{\alpha}(x^{(\text{test})}, m^{(\text{test})}) = \left[ \widehat{q}_{\frac{\alpha}{2}} \circ \Phi(x^{(\text{test})}, m^{(\text{test})}) - \widehat{Q}_{1-\tilde{\alpha}}(S); \widehat{q}_{1-\frac{\alpha}{2}} \circ \Phi(x^{(\text{test})}, m^{(\text{test})}) + \widehat{Q}_{1-\tilde{\alpha}}(S) \right]$$

then shared: the training part (1. 1-4) is common to any test point and the correction term (1.5-14) can be reused for any new test point with the same mask.

In high dimensions, many calibration points may be discarded when applying CP-MDA-Exact since it is likely that their missing patterns would not be included in the one of the test point.<sup>3</sup> This limitation brings us to the second algorithm we propose, CP-MDA-Nested.

Algorithm 2 – CP-MDA-Nested. CP-MDA with nested masking avoids the removal of calibration points whose masks are not included in that of the test point. Instead, we apply the mask of the test point to the calibration points, and so we keep all the observations (1. 3). Next, we impute the masked calibration points (1.6) before computing their scores  $s^{(k)}$  (1. 7). Then, for each calibration point, the fitted quantile regressors are used to predict on the test point with a temporary mask, which matches the mask of the given augmented calibration point. These predictions are corrected with the score of the calibration point (l. 8-9) and stored in two bags  $Z_{\frac{\alpha}{2}}$  for the lower interval boundary, and  $Z_{1-\frac{\alpha}{2}}$ for the upper interval boundary (l. 11-12). The prediction

Algorithm 2 CP-MDA-Nested (with COR) **Input:** Same as Algorithm 1 **Output:** Same as Algorithm 1 1: Compute lines 1 to 4 of Algorithm 1 // Generate an augmented calibration set: 2: for  $k \in \text{Cal } \mathbf{do}$  Additional nested masking  $\widetilde{m}^{(k)} = \max(m^{(\text{test})}, m^{(k)})$ 3: 4: end for Augmented calibration set generated. // 5: for  $k \in \operatorname{Cal} \operatorname{do}$ Impute the calibration set:  $x_{imp}^{(k)} := \Phi\left(x^{(k)}, \widetilde{m}^{(k)}\right)$ 6: Set  $s^{(k)} = \max(\hat{q}_{\frac{\alpha}{2}}(x_{\text{imp}}^{(k)}) - y^{(k)}, y^{(k)} - \hat{q}_{1-\frac{\alpha}{2}}(x_{\text{imp}}^{(k)}))$ 7: Set  $z_{\frac{\alpha}{2}}^{(k)} = \hat{q}_{\frac{\alpha}{2}} \circ \Phi\left(x^{(\text{test})}, \widetilde{m}^{(k)}\right) - s^{(k)}$ 8: Set  $z_{1-\frac{\alpha}{2}}^{(k)} = \hat{q}_{1-\frac{\alpha}{2}} \circ \Phi\left(x^{(\text{test})}, \widetilde{m}^{(k)}\right) + s^{(k)}$ 9: 10: end for 11: Set  $Z_{\frac{\alpha}{2}} = \{z_{\frac{\alpha}{2}}^{(k)}, k \in \operatorname{Cal}\}$ 12: Set  $Z_{1-\frac{\alpha}{2}} = \{z_{1-\frac{\alpha}{2}}^{(k)}, k \in \text{Cal}\}$ 13: Compute  $\widehat{Q}_{\tilde{\alpha}}(Z_{\frac{\alpha}{2}})$ 14: Compute  $\widehat{Q}_{1-\tilde{\alpha}}(Z_{1-\frac{\alpha}{2}})$ 15: Set  $\widehat{C}_{\alpha}\left(x^{(\text{test})}, m^{(\text{test})}\right)^2 = [\widehat{Q}_{\tilde{\alpha}}\left(Z_{\frac{\alpha}{2}}\right); \widehat{Q}_{1-\tilde{\alpha}}\left(Z_{1-\frac{\alpha}{2}}\right)]$ 

is finally obtained by taking the  $\alpha$  quantiles of the bags Z (1. 13-15).

The rationale for predicting on temporary test points with the mask of a given augmented calibration point is that we want to treat the test and calibration points in the same way. We should note that this method may tend to achieve conservative coverage, since the augmented calibration set may have masks that overly include the missing pattern of the test point, i.e., the augmented points may have more missing values than the test point.

#### 5.2 Theoretical guarantees in finite sample

Before analyzing the conditional coverage of CP-MDA, we start by proving the marginal validity of CP-MDA-Nested.

Theorem 5.3 (Marginal validity of CP-MDA-Nested). If A1 and A2 hold, then CP-MDA-Nested is marginally valid at the level  $1 - 2\alpha$ , for any missing mechanism.

For conditional guarantees, let us consider the following assumptions.

Assumption A3 (Y is not explained by M).  $(Y \perp M)|X$ .

Assumption A4 (Stochastic domination of the quantiles). Let  $(\mathring{m}, \breve{m}) \in \mathcal{M}^2$ . If  $\mathring{m} \subset \breve{m}$  then for any  $\delta \in [0, 0.5]$ :

$$\begin{array}{l} \bullet \ q_{1-\delta/2}^{Y|(X_{\mathrm{obs}(\tilde{m})},M=\tilde{m})} \leq q_{1-\delta/2}^{Y|(X_{\mathrm{obs}(\tilde{m})},M=\tilde{m})}, \\ \bullet \ q_{\delta/2}^{Y|(X_{\mathrm{obs}(\tilde{m})},M=\tilde{m})} \geq q_{\delta/2}^{Y|(X_{\mathrm{obs}(\tilde{m})},M=\tilde{m})}. \end{array}$$

This assumption grasps the underlying intuition that the conditional distribution of  $Y|(X_{obs(m)}, M = m)$  tends to have larger deviations when the number of observed variables is smaller, in concordance with the intuition that observ-

<sup>&</sup>lt;sup>3</sup>Yet, these discarded points could be used for training but this comes at the cost of fitting a different model for each pattern; such a path is reasonable if the data is scarce.

ing predictive variables reduce the conditional randomness of  $Y|X_{\rm obs}$ .

The following theorems (proved in Appendix D) state the finite sample guarantees of CP-MDA.

**Theorem 5.4** (Mask-conditional-validity of CP-MDA). *Assume the missing mechanism is MCAR, and A1 to A3. Then:* 

1. CP-MDA-Exact is mask-conditionally-valid;

2. *if the scores are almost surely distinct, CP-MDA-Exact is exactly mask-conditionally-valid;* 

3. if A4 also holds, CP-MDA-Nested is mask-conditionallyvalid, up to a technical minor modification of the output.

The challenge in proving the validity of CP-MDA-Nested is that the augmented calibration and test points are not exchangeable conditional on the mask and thus may result in under-coverage. However, by imposing A4 we prove that this violation of exchangeability still leads to valid (and often conservative) coverage (see Lemma D.3). We conjecture that CP-MDA-Nested attains valid coverage (without any modification), as also demonstrated by our experiments. However, we could not prove its validity without making an independence assumption which we prefer to avoid as exchangeability is key to imputation methods. Instead, we prove in Theorem D.4 the validity of any variant outputting  $[\widehat{Q}_{\tilde{\alpha}}(Z_{\frac{\tilde{m}}{2}}^{\tilde{m}}); \widehat{Q}_{1-\tilde{\alpha}}(Z_{1-\frac{\tilde{m}}{2}}^{\tilde{m}})]$  for  $Z_{\frac{\tilde{m}}{2}}^{\tilde{m}}$  the subset of  $Z_{\frac{\tilde{m}}{2}}$  composed with points using mask  $\tilde{m}$  at 1. 6-9.

**Theorem 5.5** (Marginal validity of CP-MDA-Exact). Assume the missing mechanism is MCAR, and A1 to A3. Then CP-MDA-Exact is marginally valid.

# 6 Towards asymptotic individualized coverage

Achieving validity conditionally on the mask is an important step toward reliable predictions. Nevertheless, in practice one aims at the strongest coverage conditional on both Xand M. Lei & Wasserman (2014); Vovk (2012); Barber et al. (2021a) studied a related question (without considering missing patterns) and concluded that it is impossible to achieve informative intervals satisfying conditional coverage,  $\mathbb{P}(Y \in C_{\alpha}(x) | X = x) \geq 1 - \alpha$  for any  $x \in \mathcal{X}$  in the distribution-free and finite samples setting. Still, we can analyze the asymptotic regime, similarly to Theorem 1 of Sesia & Candès (2020), which proves the asymptotic conditional validity of CQR (without the presence of missing values) under consistency assumptions on the underlying quantile regressor. Here, by contrast, we study the asymptotic conditional validity of the impute-then-predict+conformalization procedure, by analyzing the consistency of impute-thenregress in quantile regression. That is, we aim at showing that we satisfy the required assumption of consistency to invoke Theorem 1 of Sesia & Candès (2020). The proofs of this section are given in Appendix E.

To analyze the consistency of impute-then-predict procedures for quantile regression, we extend the work of Le Morvan et al. (2021) on mean regression. Quantile regression with missing values, for a quantile level  $\beta$ , aims at solving

$$\min_{f:\mathcal{X}\times\mathcal{M}\to\mathbb{R}}\mathcal{R}_{\ell_{\beta}}(f) := \mathbb{E}\left[\ell_{\beta}\left(Y,f\left(X,M\right)\right)\right],\qquad(3)$$

with  $\ell_{\beta}$  the pinball loss  $\ell_{\beta}(y, \hat{y}) = \rho_{\beta}(y - \hat{y})$  and  $\rho_{\beta}(u) = \beta |u| \mathbb{1}_{\{u \ge 0\}} + (1 - \beta) |u| \mathbb{1}_{\{u \le 0\}}.$ 

An associated  $\ell_{\beta}$ -Bayes predictor minimizes Eq. (3). Its risk is called the  $\ell_{\beta}$ -Bayes risk, noted  $\mathcal{R}^*_{\ell_{\beta}}$ . Impute-then-predict procedure in quantile regression aims at solving

$$\min_{g:\mathcal{X}\to\mathbb{R}}\mathcal{R}_{\ell_{\beta},\Phi}(g) := \mathbb{E}\left[\ell_{\beta}\left(Y,g\circ\Phi\left(X,M\right)\right)\right],\quad(4)$$

for  $\Phi$  any imputation. Let  $g^*_{\ell_{\beta},\Phi} \in \operatorname{arg\,min}_g \mathcal{R}_{\ell_{\beta},\Phi}(g)$ . The following proposition states that  $\mathcal{R}_{\ell_{\beta},\Phi}(g^*_{\ell_{\beta},\Phi}) = \mathcal{R}^*_{\ell_{\beta}}$  and the consistency of a universal learner.

**Proposition 6.1** ( $\ell_{\beta}$ -consistency of an universal learner). Let  $\beta \in [0, 1]$ . If X admits a density on X, then, for almost all imputation function  $\Phi \in \mathcal{F}_{\infty}^{I}$ , (i)  $g_{\ell_{\beta},\Phi}^{*} \circ \Phi$  is  $\ell_{\beta}$ -Bayes-optimal (ii) any universally consistent algorithm for quantile regression trained on the data imputed by  $\Phi$  is  $\ell_{\beta}$ -Bayes-consistent (i.e., asymptotically in the training set size).

Note that this quantile regression case does not require  $\mathbb{E}\left[\varepsilon|X_{\text{obs}(M)}, M\right] = 0$ , contrary to the quadratic loss case (Le Morvan et al., 2021). We conclude our asymptotic analysis of conditional coverage with Corollary 6.2.

**Corollary 6.2.** For any missing mechanism, for almost all imputation function  $\Phi \in \mathcal{F}^I_{\infty}$ , if  $F_{Y|(X_{\text{obs}(M)},M)}$  is continuous, a universally consistent quantile regressor trained on the imputed data set yields asymptotic conditional coverage.

In words, the intervals obtained by taking Bayes predictors of levels  $\alpha/2$  and  $1 - \alpha/2$  are exactly valid conditionally to both the mask M and the observed variables  $X_{obs(M)}$ , if  $F_{Y|(X_{obs(M)},M)}$  is continuous. Importantly, while this result is asymptotic, it holds for *any* missing mechanism and it considers individualized conditional coverage.

#### 7 Empirical study

**Setup.** In all experiments, the data are imputed using iterative regression (iterative ridge implemented in Scikit-learn, Pedregosa et al. (2011)).<sup>4</sup> We compare the performance of our CQR-MDA-Exact and CQR-MDA-Nested (that is CP-MDA based on CQR) to CQR as well as to a vanilla QR (without any calibration). The predictive models are fitted on the imputed data concatenated with the mask. Without concatenating the mask to the features, the mask-conditional validity of QR is worsened, as demonstrated in Section 4. The prediction algorithm is a Neural Network (NN), fitted to minimize the pinball loss (Sesia & Romano, 2021, see Appendix F.1 for details). For the vanilla QR, we use both the training and calibration sets for training.

<sup>&</sup>lt;sup>4</sup>Theoretical results hold for any symmetric imputation. In practice, constant, mean and MICE imputations gave similar results.

Synthetic and semi-synthetic experiments. We designed the training and calibration data to have 20% of MCAR values. To evaluate the test marginal coverage  $\mathbb{P}(Y \in \hat{C}_{\alpha}(X, M))$ , missing values are introduced in the test set according to the same distribution as on the training and calibration sets. Then, to compute an estimator of  $\mathbb{P}(Y \in \hat{C}_{\alpha}(X, m)|M = m)$  for each  $m \in \mathcal{M}$ , we fix to a constant the number of observations per pattern, to ensure that the variability in coverage is not impacted by  $\mathbb{P}(M = m)$ . All experiments are repeated 100 times with different splits.

#### 7.1 Synthetic experiments: Gaussian linear data

**Data generation.** The data is generated according to Model 4.1, with  $X \sim \mathcal{N}(\mu, \Sigma)$ , with  $\mu = (1, \dots, 1)^T$  and  $\Sigma = \varphi(1, \dots, 1)^T(1, \dots, 1) + (1-\varphi)I_d$ ,  $\varphi = 0.8$ , Gaussian noise  $\varepsilon \sim \mathcal{N}(0, 1)$  and the following regression coefficients  $\beta = (1, 2, -1, 3, -0.5, -1, 0.3, 1.7, 0.4, -0.3)^T$ . Here, the oracle intervals are known (Proposition 4.2).

**Coverage and length for worst and best groups.** Figures 1 and 8 (Appendix F.2) show the coverage and the length of the intervals as a function of the training set size for, respectively, the "worst" and "best group". The calibration size is fixed to 1000 and the test set contains 2000 points with the "hardest" mask, i.e. the one leading to the lowest coverage (here it corresponds to cases where only  $X_4$  is observed) and 2000 points with the "easiest" mask, i.e. with the highest coverage (here it corresponds to all the variables observed). These figures highlight that:

- **CQR** and **QR** conditional coverage improve when the training size increases (Corollary 6.2);
- **Both versions of CQR-MDA** are mask-conditionallyvalid even for the worst pattern (Theorem 5.4);
- **CQR-MDA-Exact** is exactly mask-conditionally-valid as the coverage on the worst and best patterns are exactly 90% (Theorem 5.4).

**Coverage and length by mask size.** Figure 3 displays the average coverage and intervals' length as a function of the pattern size, i.e., the performance metrics are aggregated by the masks with the same number of missing variables; the first violin plot of each panel corresponds to the marginal coverage (see Appendix F.2 for QR results). Note that only the pattern sizes are presented and not the patterns themselves as there are  $2^d = 1024$  possible masks.<sup>5</sup> For each pattern size, 100 observations are drawn according to the distribution of M | size(M) in the test set. The training and calibration sizes are respectively 500 and 250 (Figure 10 contains the results for other sizes). Figure 3 shows that:

• CQR is marginally valid (Proposition 3.3);

• CQR and QR undercover with an increasing number of



Figure 3: Average coverage (top) and length (bottom) as a function of the number of missing values (NA). The first violin plot shows the marginal coverage. #Tr = 500 and #Cal = 250. The marginal test set includes 2000 observations. The mask-conditional test set includes 100 individuals for each missing data pattern size.

missing values. This can be explained because their length nearly does not vary with the size of the missing pattern, despite having the mask concatenated with the features;

- Both versions of CQR-MDA are marginally valid (Theorem 5.5) and mask(-size)-conditionally-valid (Theorem 5.4);
- **CQR-MDA-Exact** is exactly mask(-size)-conditionallyvalid (Theorem 5.4) and its length is close to the oracle ones. It has more variability for the patterns with few missing values as for these masks Cal<sup>(test)</sup> is smaller.

#### 7.2 Semi-synthetic experiments

We consider 6 benchmark real data sets for regression: meps\_19, meps\_20, meps\_21 (MEPS), bio, bike and concrete (Dua & Graff, 2017), where we introduce missing values in their quantitative features. Note that some patterns have a low (or null) frequency of appearance in the training sets of bio and concrete. The sample sizes for training, calibration, and testing, and simulation details are provided in Appendix F.3, along with results for smaller training and calibration sets.

Figure 4 depicts the results by combining *validity* and *efficiency* (length) for meps\_19, bio, concrete, and bike, where this graph follows the visualization used in Zaffran et al. (2022). The results for meps\_20 and meps\_21 are given in Appendix F.3, as they are similar to meps\_19. Each of the panels in Figure 4 summarizes the results for one data set, with the average coverage shown in the *x*-axis and the average length in the *y*-axis. A method is mask-conditionally-valid if all the markers of its color are at the right of the vertical dotted line (90%). We observe that:

• CQR is marginally valid (orange diamond ♦, Proposition 3.3). It is not mask-conditionally-valid as the lowest coverage (orange down triangle ▼) is far below 90% for the bio, concrete, and bike data sets;

• CQR-MDA-Exact is marginally valid (purple diamond

<sup>&</sup>lt;sup>5</sup>Note that in practice the relationship between the coverage and the number of missing values is not necessarily monotonic as a mask with only one missing value can lead to more uncertainty than a mask with many missing values, see Appendix C.



Figure 4: Validity and efficiency with missing values for 4 data sets (panels) with *d* features, including *l* quantitative ones in which missing values are introduced with probability 0.2. Colors represent the methods. Diamonds ( $\blacklozenge$ ) represent marginal coverage while the patterns giving the lowest and highest coverage are represented with triangles ( $\checkmark$  and  $\blacktriangle$ ). Vertical dotted lines represent the target coverage of 90%.

♦, Theorem 5.5). It is also exactly mask-conditionally-valid, as the lowest (purple down triangle ▼) and highest (purple up triangle ▲) coverages are about 90% (Theorem 5.4);
• CQR-MDA-Nested is marginally valid (blue diamond ♦, Theorem 5.5). It is also mask-conditionally-valid, as the lowest (blue down triangle ▼) and the highest (blue up triangle ▲) coverages are larger than 90% (Theorem 5.4).

#### 7.3 Predicting the level of platelets for trauma patients

We study the applicability and robustness of CPMDA on the critical care TraumaBase® data. We focus on predicting the level of platelets of severely injured patients upon arrival at the hospital. This level is directly related to the occurrence of hemorrhagic shock and is difficult to obtain in real-time: predicting it accurately could be crucial to anticipate the need for transfusion and blood resources. In addition, this prediction task appears to be challenging as Jiang et al. (2022) achieved an average relative prediction error ( $||\hat{y} - y||^2/||y||^2$ ) that is no lower than 0.23. This highlights the need for reliable uncertainty quantification.

After applying inclusion and exclusion criteria obtained by medical doctors and following the pipeline of Sportisse et al. (2020) described in Appendix F.4, we left with a subset of 28855 patients and 7 features. Missing values vary from 0% to 24% by features, with a total average of 7%.

**Results.** The results are summarized in Figure 5, where we use different markers to denote the different masks. To ensure a fair comparison between the conformal methods, we only keep the missing patterns for which there are more than 200 individuals; this excludes 7 patterns. Finally, since we found that the vanilla QR tends to be overly conservative, we refer to Appendix F.4 for its results. Figure 5 shows that all conformal approaches achieve marginal coverage higher than the desired 90% level (diamonds ♦). Furthermore, for each mask (each set of linked markers) **CQR**-**MDA** improves coverage compared to **CQR** by approaching

90%, and efficiency by reducing the average length. Noticeably, for the pattern corresponding to all features observed (squares ■), **CQR-MDA** has a coverage rate above 90% while **CQR** is below the target level. Therefore, we believe **CQR-MDA** should be recommended as it improves upon the vanilla impute-then-regress+CQR approach.



Figure 5: Average coverage and length on the TraumaBase® analysis. See the caption of Figure 4 for details. Other symbols than diamond correspond to computing the average per mask. Each individual's prediction is obtained by using 15390 observations for training, and 7694 for calibration.

#### Acknowledgements

We thank Baptiste Goujaud for interesting discussions. This work was supported by a public grant as part of the Investissement d'avenir project, reference ANR-11-LABX-0056-LMH, LabEx LMH. M. Zaffran has been awarded the 2022 Scholarship for Mathematics granted by the Séphora Berrebi Foundation which she gratefully thanks for its support. The work of A. Dieuleveut is partially supported by ANR-19-CHIA-0002-01/chaire SCAI and Hi! Paris. The work of J. Josse is partially supported by ANR-16-IDEX-0006. Y. Romano was supported by the ISRAEL SCIENCE FOUNDATION (grant No. 729/21). He also thanks the Career Advancement Fellowship, Technion, for providing additional research support.

#### References

- Angelopoulos, A. N. and Bates, S. A gentle introduction to conformal prediction and distribution-free uncertainty quantification, 2021. URL https://arxiv.org/ abs/2107.07511.
- Ayme, A., Boyer, C., Dieuleveut, A., and Scornet, E. Nearoptimal rate of consistency for linear models with missing values. In Chaudhuri, K., Jegelka, S., Song, L., Szepesvari, C., Niu, G., and Sabato, S. (eds.), Proceedings of the 39th International Conference on Machine Learning, volume 162 of Proceedings of Machine Learning Research, pp. 1211–1243. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/ v162/ayme22a.html.
- Barber, R. F., Candès, E. J., Ramdas, A., and Tibshirani, R. J. The limits of distribution-free conditional predictive inference. *Information and Inference: A Journal of the IMA*, 10(2):455–482, June 2021a. ISSN 2049-8764, 2049-8772. doi: 10.1093/ imaiai/iaaa017. URL https://academic.oup. com/imaiai/article/10/2/455/5896927.
- Barber, R. F., Candès, E. J., Ramdas, A., and Tibshirani, R. J. Predictive inference with the jackknife+. *The Annals of Statistics*, 49(1):486–507, 2021b.
- Barber, R. F., Candès, E. J., Ramdas, A., and Tibshirani, R. J. Conformal prediction beyond exchangeability, 2022. URL https://arxiv.org/abs/2202.13415.
- Dua, D. and Graff, C. UCI machine learning repository, 2017. URL http://archive.ics.uci.edu/ml.
- Eaton, M. L. *Multivariate statistics*. Probability & Mathematical Statistics S. John Wiley & Sons, Nashville, TN, September 1983.
- Jiang, W., Bogdan, M., Josse, J., Majewski, S., Miasojedow, B., Ročková, V., and TraumaBase® Group. Adaptive bayesian slope: Model selection with incomplete data. *Journal of Computational and Graphical Statistics*, 31(1): 113–137, 2022. doi: 10.1080/10618600.2021.1963263.
- Josse, J. and Reiter, J. P. Introduction to the Special Section on Missing Data. *Statistical Science*, 33(2):139 – 141, 2018. doi: 10.1214/18-STS332IN. URL https:// doi.org/10.1214/18-STS332IN.
- Josse, J., Prost, N., Scornet, E., and Varoquaux, G. On the consistency of supervised learning with missing values, 2019. URL https://arxiv.org/abs/1902. 06931.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization, 2014. URL https://arxiv.org/abs/ 1412.6980.
- Le Morvan, M., Josse, J., Moreau, T., Scornet, E., and Varoquaux, G. Neumiss networks: differentiable programming

for supervised learning with missing values. *Advances in Neural Information Processing Systems*, 33:5980–5990, 2020a.

- Le Morvan, M., Prost, N., Josse, J., Scornet, E., and Varoquaux, G. Linear predictor on linearly-generated data with missing values: non consistency and solutions. In Chiappa, S. and Calandra, R. (eds.), *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pp. 3165–3174. PMLR, 26– 28 Aug 2020b. URL https://proceedings.mlr. press/v108/morvan20a.html.
- Le Morvan, M., Josse, J., Scornet, E., and Varoquaux, G. What's a good imputation to predict with missing values? In Advances in Neural Information Processing Systems, volume 34, 2021. URL https://proceedings. neurips.cc/paper/2021/file/ 5fe8fdc79ce292c39c5f209d734b7206-Paper. pdf.
- Lei, J. and Wasserman, L. Distribution-free prediction bands for non-parametric regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):71–96, 2014. ISSN 1467-9868. doi: https://doi.org/10.1111/rssb.12021.
  URL https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/rssb.12021. \_eprint: https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/rssb.12021.
- Lei, J., G'Sell, M., Rinaldo, A., Tibshirani, R. J., and Wasserman, L. Distribution-Free Predictive Inference for Regression. *Journal of the American Statistical Association*, 113(523):1094–1111, 2018.
- Little, R. J. A. *Statistical analysis with missing data, third edition*. Wiley Series in Probability and Statistics. John Wiley & Sons, Nashville, TN, 3 edition, May 2019.
- Manokhin, V. Awesome conformal prediction, April 2022. URL https://doi.org/10.5281/zenodo. 6467205.
- Mayer, I., Sportisse, A., Josse, J., Tierney, N., and Vialaneix, N. R-miss-tastic: a unified platform for missing values methods and workflows, 2019. URL https://arxiv. org/abs/1908.04822.
- MEPS. Medical expenditure panel survey. https://meps.ahrq.gov/mepsweb/data\_ stats/data\_overview.jsp.
- Papadopoulos, H., Proedrou, K., Vovk, V., and Gammerman, A. Inductive Confidence Machines for Regression. In *Machine Learning: ECML 2002*, pp. 345–356. Springer, 2002.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

Romano, Y., Patterson, E., and Candès, E. Conformalized Quantile Regression. In Advances in Neural Information Processing Systems, volume 32, 2019. URL https: //papers.nips.cc/paper/2019/hash/ 5103c3584b063c431bd1268e9b5e76fb-Abstract. html.

- Romano, Y., Barber, R. F., Sabatti, C., and Candès, E. With Malice Toward None: Assessing Uncertainty via Equalized Coverage. *Harvard Data Science Review*, 2(2), April 2020. doi: 10.1162/99608f92.03f00592. Publisher: Pub-
- Rubin, D. B. Inference and missing data. *Biometrika*, 63(3): 581–592, 1976.

Pub.

- Sesia, M. and Candès, E. J. A comparison of some conformal quantile regression methods. *Stat*, 9(1): e261, 2020. doi: https://doi.org/10.1002/sta4.261. URL https://onlinelibrary.wiley.com/doi/ abs/10.1002/sta4.261. e261 sta4.261.
- Sesia, M. and Romano, Y. Conformal prediction using conditional histograms. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), Advances in Neural Information Processing Systems, volume 34, pp. 6304–6315. Curran Associates, Inc., 2021. URL https://proceedings. neurips.cc/paper/2021/file/ 31b3b31a1c2f8a370206f111127c0dbd-Paper. pdf.
- Sportisse, A., Boyer, C., Dieuleveut, A., and Josse, J. Debiasing averaged stochastic gradient descent to handle missing values. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 12957–12967. Curran Associates, Inc., 2020. URL https://proceedings. neurips.cc/paper/2020/file/ 972ededf6c4d7c1405ef53f27d961eda-Paper. pdf.
- Van Ness, M., Bosschieter, T. M., Halpin-Gregorio, R., and Udell, M. The missing indicator method: From low to high dimensions, 2022. URL https://arxiv.org/ abs/2211.09259.
- Vovk, V. Conditional Validity of Inductive Conformal Predictors. In Asian Conference on Machine Learning, pp. 475–490. PMLR, November 2012. URL http://proceedings.mlr.press/ v25/vovk12.html. ISSN: 1938-7228.
- Vovk, V., Gammerman, A., and Shafer, G. *Algorithmic Learning in a Random World*. Springer US, 2005.
- Zaffran, M., Féron, O., Goude, Y., Josse, J., and Dieuleveut,
  A. Adaptive conformal predictions for time series.
  In Chaudhuri, K., Jegelka, S., Song, L., Szepesvari,
  C., Niu, G., and Sabato, S. (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learn*-

ing Research, pp. 25834-25866. PMLR, 17-23 Jul 2022. URL https://proceedings.mlr.press/ v162/zaffran22a.html.

# Appendices

The appendices are organized as follows.

Appendix A describes CQR, used in the paper.

Appendix B provides an explicit description of impute-then-predict+conformalization (Appendix B.1), along with its proof of validity, that is the proofs for Section 3 (Appendix B.2).

Then, Appendix C contains the proofs for the Gaussian linear model oracle intervals presented in Section 4 (Appendix C.1), along with the discussion on how mean-based approaches fail (Appendix C.2).

Appendix D gives the general statement of CP-MDA-Exact (Appendix D.1), and the proofs of the validity theorems for CP-MDA-Exact (Appendix D.2), along with the theoretical study of CP-MDA-Nested (Appendix D.3).

Appendix E provides all the proofs about consistency and asymptotic conditional coverage presented in Section 6.

Finally, Appendix F contains all the details for the experimental study and additional results completing Section 7. More precisely, Appendix F.1 gives more details about the settings. Appendix F.2 contains results on synthetic data. Appendix F.3 describes the real data sets used for the semi-synthetic experiments, and presents the remaining results. Appendix F.4 presents the real medical data set (TraumaBase®), the pipeline and settings used and the results obtained by QR on this data set.

### A Illustration and details on CQR (Romano et al., 2019) procedure

Figure 6 provides a visualization and step by step description of CQR.

#### **B** Impute-then-predict+conformalization

#### **B.1** Description of the algorithm

Algorithm 3 SCP on impute-then-predict

**Input:** Imputation algorithm  $\mathcal{I}$ , predictive algorithm  $\mathcal{A}$ , conformity score function s, significance level  $\alpha$ , training set  $\{(X^{(1)}, M^{(1)}, Y^{(1)}), \dots, (X^{(n)}, M^{(n)}, Y^{(n)})\}.$ 

**Output:** Prediction interval  $\widehat{C}_{\alpha}(X, M)$ .

- 1: Randomly split  $\{1, \ldots, n\}$  into two disjoint sets Tr and Cal.
- 2: Fit the imputation function:  $\Phi(\cdot) \leftarrow \mathcal{I}\left(\left\{\left(X^{(k)}, M^{(k)}\right), k \in \mathrm{Tr}\right\}\right)$

3: Impute the data set: 
$$\left\{X_{imp}^{(k)}\right\}_{k=1}^{n} := \left\{\Phi\left(X^{(k)}, M^{(k)}\right)\right\}_{k=1}^{n}$$

4: Fit algorithm 
$$\mathcal{A}$$
:  $\hat{g}(\cdot) \leftarrow \mathcal{A}\left(\left\{\left(X_{imp}^{(k)}, Y^{(k)}\right), k \in \mathrm{Tr}\right\}\right)$ 

5: for  $k \in \text{Cal do}$ 

6: Set 
$$S^{(k)} = s\left(Y^{(k)}, \hat{g}\left(X_{imp}^{(k)}\right)\right)$$
, the conformity scores

- 7: end for
- 8: Set  $\mathcal{S}_{Cal} = \{S^{(k)}, k \in Cal\}$

9: Compute 
$$\widehat{Q}_{1-\alpha^{\text{SCP}}}(\mathcal{S}_{\text{Cal}})$$
, the  $1-\alpha^{\text{SCP}}$ -th empirical quantile of  $\mathcal{S}_{\text{Cal}}$ , with  $1-\alpha^{\text{SCP}} := (1-\alpha)(1+1/\#\text{Cal})$ .

10: Set  $\widehat{C}_{\alpha}(X, M) = \left\{ y \text{ such that } s(y, \widehat{g} \circ \Phi(X, M)) \le \widehat{Q}_{1-\alpha^{\text{SCP}}}(\mathcal{S}_{\text{Cal}}) \right\}.$ 

Similarly, Algorithm 1 can be written to include any underlying predictive algorithm (regression or classification) and any score function.



Create a proper training set, a calibration set, and keep your test set, by randomly splitting your data set.

On the proper training set:

► Learn  $\hat{q}_{\text{low}}$  and  $\hat{q}_{\text{upp}}$ 

On the calibration set:

- ▶ Predict with  $\hat{q}_{low}$  and  $\hat{q}_{upp}$
- ► Get the scores  $s^{(k)} = \max\left\{\hat{q}_{\text{low}}\left(x^{(k)}\right) - y^{(k)}, y^{(k)} - \hat{q}_{\text{upp}}\left(x^{(k)}\right)\right\}$
- ► Compute the (1 − α) × (1 + <sup>1</sup>/<sub>#Cal</sub>) empirical quantile of the s<sup>(k)</sup>, noted Q
  <sub>1−â</sub> (S)

On the test set:

- ▶ Predict with  $\hat{q}_{\text{low}}$  and  $\hat{q}_{\text{upp}}$
- ▶ Build  $\hat{C}_{\hat{\alpha}}(x)$ :  $[\hat{q}_{\text{low}}(x) \hat{Q}_{1-\hat{\alpha}}(S), \hat{q}_{\text{upp}}(x) + \hat{Q}_{1-\hat{\alpha}}(S)]$

Figure 6: Schematic illustration of Conformalized Quantile Regression (CQR) (Romano et al., 2019).

#### **B.2** Proof of exchangeability after imputation

In this subsection, we provide a more formal statement of Lemma 3.2 and Proposition 3.3 in respectively Lemma B.1 and Proposition B.2. To that end, we introduce a notion of symmetrical imputation *on a set*  $\mathcal{T}$ , for  $\mathcal{T} \subset [\![1, n + 1]\!]$ .

Assumption A5 (Symmetrical imputation on a set  $\mathcal{T}$ ). For a given set of points  $\{X^{(k)}, M^{(k)}, Y^{(k)}\}_{k \in \mathcal{T}}$  the imputation function  $\Phi$  is the output of an algorithm  $\mathcal{I}$  that treats the data points in  $\mathcal{T}$  symmetrically:  $\mathcal{I}(\{X^{(k)}, M^{(k)}, Y^{(k)}\}_{k \in \mathcal{T}}) \stackrel{(d)}{=} \mathcal{I}(\{X^{(\sigma(k))}, M^{(\sigma(k))}, Y^{(\sigma(k))}\}_{k \in \mathcal{T}} \text{ conditionally to } \{X^{(k)}, M^{(k)}, Y^{(k)}\}_{k \in \mathcal{T}} \text{ and for any permutation } \sigma \text{ on } [1, \#\mathcal{T}].$ 

**Lemma B.1** (Imputation preserves exchangeability). Let A1 hold. Then, for any missing mechanism, for any imputation function  $\Phi$  satisfying A5, the imputed random variables  $\left(\Phi\left(X^{(k)}, M^{(k)}\right), M^{(k)}, Y^{(k)}\right)_{k \in \mathcal{T}}$  are exchangeable.

**Proposition B.2** ((Exact) validity of impute-then-predict+conformalization). *If A1 is satisfied, then we have the following three results.* 

- 1. Full CP: if A5 is satisfied for  $\mathcal{T} = [1, n + 1]$  (i.e., the imputation algorithm treats all points symmetrically), then impute-then-predict+Full CP is marginally valid. If moreover the scores are almost surely distinct, it is exactly valid. OR
- 2. Jackknife+ if A5 is satisfied for  $\mathcal{T} = [\![1, n + 1]\!]$  (i.e., the imputation algorithm treats all points symmetrically), then impute-then-predict+Jackknife+ is marginally valid (of level  $1 2\alpha$ ).

OR

3. SCP with the split  $[1, n + 1] = \text{Tr} \bigcup \text{Cal} \bigcup \text{Test}$  and if A5 is satisfied for  $\mathcal{T} = \text{Cal} \bigcup \text{Test}$  (i.e., the imputation treats all points in Cal  $\bigcup$  Test symmetrically) then impute-then-predict+conformalization is marginally valid. If moreover the scores are almost surely distinct, it is exactly valid.

*Remark* B.3 (Imputation choices for SCP). In the latter case, for SCP, the coverage result can be derived conditionally on Tr, thus the coverage results holds for: (i) any deterministic imputation function (conditionally on Tr) (that is any arbitrary function of Tr), or (ii) any stochastic imputation function treating Cal and Test symmetrically (iii) any combination of both.

Proof of Lemma B.1.  $\Phi$  is the output of an imputing algorithm  $\mathcal{I}$  trained on  $\left\{ \left( X^{(k)}, M^{(k)}, Y^{(k)} \right)_{k \in \mathcal{I}} \right\}$ .

Assume  $(X^{(k)}, M^{(k)}, Y^{(k)})_{k \in \mathcal{T}}$  are exchangeable (A1).

Thus, if  $\mathcal{I}$  treats the data points in  $\mathcal{T}$  symmetrically,  $(\Phi(X^{(k)}, M^{(k)}), M^{(k)}, Y^{(k)})_{k \in \mathcal{T}}$  are exchangeable (see proof of Theorem 1b in (Barber et al., 2022) for example).

*Proof of Proposition B.2.* Proposition B.2 is a consequence of Lemma B.1 with different choices of  $\mathcal{T}$ , that enable to apply the following results:

- 1. Full CP: Vovk et al. (2005), also re-stated in Barber et al. (2022)
- 2. Jackknife+: Barber et al. (2021b)
- 3. SCP: Lei et al. (2018) or Papadopoulos et al. (2002) and Angelopoulos & Bates (2021) for a generic version with any score function (note that the coverage is proved conditionally on Tr).

#### C Gaussian linear model

C.1 Distribution of  $Y|(X_{obs(m)}, M)$  and oracle intervals

**Proposition C.1** (Distribution of  $Y|(X_{obs(M)}, M)$  (Le Morvan et al., 2020b)). Under Model 4.1, for any  $m \in \{0, 1\}^d$ :

$$Y|(X_{\text{obs}(\mathbf{m})}, M = m) \sim \mathcal{N}\left(\tilde{\mu}^m, \tilde{\Sigma}^m\right)$$

with:

$$\begin{split} \tilde{\mu}^m &= \beta^T_{\rm obs(m)} X_{\rm obs(m)} + \beta^T_{\rm mis(m)} \mu^m_{\rm mis|obs} \\ \mu^m_{\rm mis|obs} &= \mu^m_{\rm mis(m)} + \Sigma^m_{\rm mis(m),obs(m)} (\Sigma^m_{\rm obs(m),obs(m)})^{-1} (X_{\rm obs(m)} - \mu^m_{\rm obs(m)}), \\ \tilde{\Sigma}^m &= \beta^T_{\rm mis(m)} \Sigma^m_{\rm mis|obs} \beta_{\rm mis(m)} + \sigma^2_{\varepsilon} \\ \Sigma^m_{\rm mis|obs} &= \Sigma^m_{\rm mis(m),mis(m)} - \Sigma^m_{\rm mis(m),obs(m)} (\Sigma^m_{\rm obs(m),obs(m)})^{-1} \Sigma^m_{\rm obs(m),mis(m)} \end{split}$$

**Proposition C.2** (Oracle intervals). Under Model 4.1, for any  $m \in \{0, 1\}^d$ , for any  $\delta \in (0, 1)$ :

$$q_{\delta}^{Y|(X_{\text{obs}(\text{m})},M=m)} = \beta_{\text{obs}(\text{m})}^{T} X_{\text{obs}(\text{m})} + \beta_{\text{mis}(m)}^{T} \mu_{\text{mis}|\text{obs}}^{m} + q_{\delta}^{\mathcal{N}(0,1)} \sqrt{\beta_{\text{mis}(m)}^{T} \Sigma_{\text{mis}|\text{obs}}^{m} \beta_{\text{mis}(m)} + \sigma_{\varepsilon}^{2}},$$

and the oracle predictive interval length is given by:

$$\mathcal{L}^*_{\alpha}(m) = 2q_{1-\frac{\alpha}{2}}^{\mathcal{N}(0,1)} \sqrt{\beta_{\mathrm{mis}(m)}^T \Sigma_{\mathrm{mis}|\mathrm{obs}}^m \beta_{\mathrm{mis}(m)} + \sigma_{\varepsilon}^2}.$$
(5)

*Proof.* Using multivariate Gaussian conditioning (Eaton, 1983), for any subset of indices  $L \in [1, d]$ :

$$X_K|(X_L, M) \sim \mathcal{N}(\mu_{K|L}^M, \Sigma_{K|L}^M), \tag{6}$$

with  $K = \overline{L}$  (the complement indices) and:

$$\mu_{K|L}^{M} = \mu_{K}^{M} + \Sigma_{K,L}^{M} \Sigma_{L,L}^{M^{-1}} (X_{L} - \mu_{L}^{M}),$$
  
$$\Sigma_{K|L}^{M} = \Sigma_{K,K}^{M} - \Sigma_{K,L}^{M} \Sigma_{L,L}^{M^{-1}} \Sigma_{L,K}^{M}.$$

Given that  $Y = \beta^T X + \varepsilon$ , with  $\varepsilon \sim \mathcal{N}(0, \sigma_{\varepsilon}^2) \perp (X, M)$ , the following holds:

$$Y|(X_L, M) \stackrel{(d)}{=} (\beta^T X + \varepsilon)|(X_L, M) \stackrel{(d)}{=} \beta^T_L X_L + (\varepsilon + \beta^T_K X_K)|(X_L, M)$$

and by Equation (6),  $\beta_K^T X_K | (X_L, M) \sim \mathcal{N}(\beta_K^T \mu_{K|L}^M, \beta_K^T \Sigma_{K|L}^M \beta_K)$ , and  $(\varepsilon | (X_L, M)) \sim \mathcal{N}(0, \sigma_{\varepsilon}^2)$ , and  $(\beta_K^T X_K \perp \varepsilon) | (X_L, M)$ . Thus:

$$Y|(X_L, M) \sim \mathcal{N}(\beta_L^T X_L + \beta_K^T \mu_{K|L}^M, \beta_K^T \Sigma_{K|L}^M \beta_K + \sigma_{\varepsilon}^2)$$

Consequently, for any  $\delta \in (0, 1)$ :

$$q_{\delta}^{Y|(X_L,M)} = \beta_L^T X_L + \beta_K^T \mu_{K|L}^M + q_{\delta}^{\mathcal{N}(0,1)} \sqrt{\beta_K^T \Sigma_{K|L}^M \beta_K + \sigma_{\varepsilon}^2}.$$
(7)

For any pattern  $m \in \{0,1\}^d$ , applying Equation (7) with  $K = mis(m) = \overline{obs(m)}$ , L = obs(m), we have, for any  $\delta \in (0,1)$ :

$$q_{\delta}^{Y|(X_{\text{obs}(m)},M=m)} = \beta_{\text{obs}(m)}^{T} X_{\text{obs}(m)} + \beta_{\text{mis}(m)}^{T} \mu_{\text{mis}|\text{obs}}^{m} + q_{\delta}^{\mathcal{N}(0,1)} \sqrt{\beta_{\text{mis}(m)}^{T} \Sigma_{\text{mis}|\text{obs}}^{m} \beta_{\text{mis}(m)} + \sigma_{\varepsilon}^{2}}$$

and:

$$\mathcal{L}^*_{\alpha}(m) = 2 \times q_{1-\alpha/2}^{\mathcal{N}(0,1)} \times \sqrt{\beta_{\mathrm{mis}(m)}^T \Sigma_{\mathrm{mis}|\mathrm{obs}}^m \beta_{\mathrm{mis}(m)} + \sigma_{\varepsilon}^2}$$

with:

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

$$\Sigma_{\mathrm{mis}|\mathrm{obs}}^{m} = \Sigma_{\mathrm{mis}(m),\mathrm{mis}(m)}^{m} - \Sigma_{\mathrm{mis}(m),\mathrm{obs}(m)}^{m} (\Sigma_{\mathrm{obs}(m),\mathrm{obs}(m)}^{m})^{-1} \Sigma_{\mathrm{obs}(m),\mathrm{mis}(m)}^{m}.$$

#### C.2 Discussion on how mean-based approaches fail

Under Model 4.1, the Bayes predictor for a quadratic loss in presence of missing values –  $\mathbb{E}\left[Y \mid (X_{obs(M)}, M)\right]$  – is fully characterized (Le Morvan et al., 2020b,a; Ayme et al., 2022).

Figure 7 is obtained by generating the data according to Model 4.1 with d = 3,  $\beta = (1, 2, -1)^T$  and  $\sigma_{\varepsilon} = 1$ , with multivariate Gaussian X and MCAR mechanism  $(X \perp M)$  (which is a particular case of Model 4.1 with  $\mu^m \equiv \mu$  and  $\Sigma^m \equiv \Sigma$ ). The left panel represents the method *Oracle mean* + *SCP* where SCP is applied on the regressor being the Bayes predictor for the mean with absolute residuals as the score function. The first violin plot represents the marginal coverage whereas the other 7 represent conditional coverage with respect to the different possible patterns: conditional on observing all the variables, on observing all the variables except  $X_1$ , except  $X_2$  etc (see Section 7 for details on the simulation process).



Figure 7: Calibration set contains 500 points. Test size for each pattern is of 500 individuals and for marginal is of 2000. 200 repetitions allow to display violin plots, the horizontal black line representing the mean.

**SCP on a (oracle) mean regressor lacks of conditional coverage with respect to the mask.** Figure 7 (left) highlights that even with the best mean regressor (the Bayes predictor) and an homoskedastic noise, usual SCP intervals:

- over-cover when there are no missing values;
- cover less for a mask  $\breve{m}$  than for a mask  $\mathring{m}$  when  $\mathring{m} \subset \breve{m}$  (e.g.  $\mathring{m} = (1, 0, 0)$  only  $X_1$  is missing,  $\breve{m} = (1, 1, 0)$  that is  $X_1$  and  $X_2$  are missing);
- cover less when the most informative variable  $(X_2)$  is missing.

To tackle this issue, one could calibrate conditionally to the missing data patterns. This is in the same vein as calibrating conditionally to the categories of a categorical variable or to different groups (Romano et al., 2020). This strategy is not viable as there are  $2^d$  patterns: the number of subsets grows exponentially with the dimension, implying the creation of subsets with too little data to perform the calibration. As an alternative, one could consider to perform calibration conditionally to the pattern size (e.g. when d = 3, either 0 missing value, 1 or 2). This is possible as there are only d different pattern sizes.

**Calibrating by pattern size does not provide validity conditionally to the missing data patterns.** Figure 7 (right) shows the coverages of *Oracle mean* + *SCP per pattern size* where SCP is applied on the Bayes predictor for the mean and the calibration is protected by pattern size. The previous statements still hold with this strategy, even if the coverage disparities are smaller. Therefore, it is not enough to calibrate per pattern size.

#### **D** Finite sample algorithms

#### D.1 General statement of Algorithm 1

We provide in Algorithm 4 a general statement of CP-MDA-Exact handling any learning algorithm (both regression and classification) and any score function.

#### D.2 Mask-conditional valitidy of CP-MDA-Exact

Before proving the results, we introduce a slightly stronger notion of mask-conditional-validity, when the calibration set is itself of random cardinality.

**Definition D.1** (Mask-conditional-validity-random-calibration-size). A method is mask-conditionally-valid with a random calibration size #Cal if for any  $m \in \mathcal{M}$ , the lower bound is satisfied, and exactly mask-conditionally-valid if for any  $m \in \mathcal{M}$ ,  $1 \le c \le n$ , the upper bound is also satisfied:

$$1 - \alpha \leq \underset{\text{valid}}{\mathbb{P}} \left( Y^{(n+1)} \in \widehat{C}_{\alpha} \left( X^{(n+1)}, m \right) | M^{(n+1)} = m, \# \text{Cal} = c \right) \leq \underset{\text{exactly valid}}{\leq} 1 - \alpha + \frac{1}{c+1} + \frac$$

We start by proving Theorem D.2 that implies the result on CP-MDA-Exact in Theorem 5.4.

**Theorem D.2.** [Conditional validity of CP-MDA-Exact with calibration of random cardinality] Assume the missing mechanism is MCAR, and that Assumptions A1 to A3 hold. Then:

#### Algorithm 4 CP-MDA-Exact

**Input:** Imputation algorithm  $\mathcal{I}$ , predictive algorithm  $\mathcal{A}$ , conformity score function  $s_q$  parametrized by a model g, significance level  $\alpha$ , training set  $\{(X^{(k)}, M^{(k)}, Y^{(k)})\}_{k=1}^{n}$ , test point  $(X^{(\text{test})}, M^{(\text{test})})$ .

**Output:** Prediction interval  $\widehat{C}_{\alpha}(x^{(\text{test})}, m^{(\text{test})})$ .

- 1: Randomly split  $\{1, \ldots, n\}$  into two disjoint sets Tr and Cal.
- 2: Fit the imputation function:  $\Phi(\cdot) \leftarrow \mathcal{I}(\{(X^{(k)}, M^{(k)}), k \in \mathrm{Tr}\})$
- 3: Impute the training set:  $\left\{X_{imp}^{(k)}\right\}_{k\in\mathrm{Tr}} := \left\{\Phi\left(X^{(k)}, M^{(k)}\right)\right\}_{k\in\mathrm{Tr}}$ 4: Fit algorithm  $\mathcal{A}: \hat{g}(\cdot) \leftarrow \mathcal{A}\left(\left\{\left(X_{imp}^{(k)}, Y^{(k)}\right), k\in\mathrm{Tr}\right\}\right)$
- // Generate an augmented calibration set:
- 5:  $\operatorname{Cal}^{(\text{test})} = \{ k \in \operatorname{Cal} \text{ such that } M^{(k)} \subset M^{(\text{test})} \}$
- 6: for  $k \in \operatorname{Cal}^{(\text{test})}$  do
- 7:  $\widetilde{M}^{(k)} = M^{(\text{test})}$  Additional masking
- 8: end for

Augmented calibration set generated. //

9: Impute the calibration set:  $\left\{X_{imp}^{(k)}\right\}_{k \in Cal^{(test)}} := \left\{\Phi\left(X^{(k)}, \widetilde{M}^{(k)}\right)\right\}_{k \in Cal^{(test)}}$ 

- 10: for  $k \in \operatorname{Cal}^{(\text{test})} \operatorname{do}$
- Set  $S^{(k)} = s_{\hat{g}} \left( Y^{(k)}, X^{(k)}_{imp} \right)$ , the conformity scores 11:
- 12: end for
- 13: Set  $\mathcal{S}_{Cal} = \{S^{(k)}, k \in Cal^{(test)}\}$
- 14: Compute  $\widehat{Q}_{1-\tilde{\alpha}}(\mathcal{S}_{\text{Cal}})$ , the  $1-\tilde{\alpha}$ -th empirical quantile of  $\mathcal{S}_{\text{Cal}}$ , with  $1-\tilde{\alpha} := (1-\alpha)(1+1/\#\mathcal{S}_{\text{Cal}})$ .
- 15: Set  $\widehat{C}_{\alpha}\left(X^{(\text{test})}, M^{(\text{test})}\right) = \left\{y \text{ such that } s_{\widehat{g}}\left(y, \Phi\left(X^{(\text{test})}, M^{(\text{test})}\right)\right) \leq \widehat{Q}_{1-\widehat{\alpha}}\left(\mathcal{S}_{\text{Cal}}\right)\right\}.$

• if the scores  $S^{(k)}$  are almost surely distinct, CP-MDA-Exact is exactly mask-conditionally-valid with a random calibration size #Cal.

*Proof of Theorem D.2.* Let Tr and Cal be two disjoint sets on [1, n]. Let  $\hat{g}$  be some model. Given A1, the sequence  $\left\{ \left( X^{(k)}, M^{(k)}, Y^{(k)} \right)_{k \in \text{Cal}}, \left( X^{(\text{test})}, M^{(\text{test})}, Y^{(\text{test})} \right) \right\}$  is exchangeable. Therefore, the sequence  $\left\{ \left( X^{(k)}, Y^{(k)} \right)_{k \in \text{Cal}}, \left( X^{(\text{test})}, Y^{(\text{test})} \right) \right\}$  is also exchangeable.

Let m in  $\mathcal{M}$ . We define  $\operatorname{Cal}^{\mathrm{m}} = \{ k \in \operatorname{Cal} \text{ such that } M^{(k)} \subset m \}.$ 

Let  $c \in [1, \#Cal]$ .

As the  $M \perp X$  (missingness is MCAR) and  $(M \perp Y)|X$  (Assumption A3), then  $M \perp (X,Y)$ , and  $\# \operatorname{Cal}^m \perp$  $(X^{(k)}, Y^{(k)})_{k \in Cal}, (X^{(test)}, Y^{(test)})$ . It follows that the sequence  $\left\{ (X^{(k)}, Y^{(k)})_{k \in Cal^m}, (X^{(test)}, Y^{(test)}) \right\}$  is exchangeable. conditionally to  $\#Cal^m = c$ .

 $\begin{array}{lll} \text{Similarly,} & M^{(\text{test})} & \mathbb{I} & \left(X^{(k)}, Y^{(k)}\right)_{k \in \text{Cal}^m}, \left(X^{(\text{test})}, M^{(\text{test})}, Y^{(\text{test})}\right)_{k \in \text{Cal}^m}, \left(X^{(\text{test})}, M^{(\text{test})}, X^{(\text{test})}\right)_{k \in \text{Cal}^m}, \left(X^{(\text{test})}, M^{(\text{test})}, X^{(\text{test})}\right)_{k \in \text{Cal}^m}, \left(X^{(\text{test})}, X^{(\text{test})}\right)_{k \in \text{Cal}^m}$ 

Therefore, we can now invoke Proposition 3.3 in combination with Lemma 1 of Romano et al. (2020) to conclude the proof. But we can state a more rigorous version here, since in fact  $Cal^m$  is a random variable (as discussed in Definition D.1).

Since the algorithm  $\mathcal{I}$  treats the calibration and test data points symmetrically (A5 with  $\mathcal{T}$ Cal $\bigcup$ Test), A5 also holds for any  $\mathcal{T}'$  $\mathcal{T}$ . Therefore, by Lemma B.1 the sequence  $\subset$  $\left\{\left(\Phi(X^{(k)}, M^{(\text{test})}), M^{(\text{test})}, Y^{(k)}\right)_{k \in \text{Cal}^{\text{m}}}, \left(\Phi(X^{(\text{test})}, M^{(\text{test})}), M^{(\text{test})}, Y^{(\text{test})}\right)\right\} \text{ is exchangeable conditionally to a set of the se$  $#Cal^m = c \text{ and } M^{(test)} = m.$ 

The conclusion follows from usual arguments (Papadopoulos et al., 2002; Lei et al., 2018; Angelopoulos & Bates, 2021). Precisely,  $\left\{ \left( s_{\hat{g}}(Y^{(k)}, \Phi(X^{(k)}, M^{(\text{test})})) \right)_{k \in \text{Cal}^m}, s_{\hat{g}}(Y^{(\text{test})}, \Phi(X^{(\text{test})}, M^{(\text{test})})) \right\}$  is exchangeable conditionally to

<sup>•</sup> *CP-MDA-Exact is valid with a random calibration size* #Cal *conditionally to the missing patterns;* 

 $#Cal^m = c$  and  $M^{(test)} = m$ . Therefore,

$$\mathbb{P}\left(s_{\hat{g}}(Y^{(\text{test})}, \Phi(X^{(\text{test})}, M^{(\text{test})})) \leq \widehat{Q}_{1-\tilde{\alpha}}((s_{\hat{g}}(Y^{(k)}, \Phi(X^{(k)}, M^{(\text{test})})))_{k \in \text{Cal}^{m}}) \middle| M^{(\text{test})} = m, \#\text{Cal}^{m} = c\right) \geq 1-\alpha,$$

and if the  $\left(\left(s_{\hat{g}}(Y^{(k)}, \Phi(X^{(k)}, M^{(\text{test})}))\right)_{k \in \text{Cal}^{\text{m}}}, s_{\hat{g}}(Y^{(\text{test})}, \Phi(X^{(\text{test})}, M^{(\text{test})}))\right)$  are almost surely distinct (i.e. have a continuous distribution) then (Lei et al., 2018; Romano et al., 2019):

$$\mathbb{P}\left(s_{\hat{g}}(Y^{(\text{test})}, \Phi(X^{(\text{test})}, M^{(\text{test})})) \le \widehat{Q}_{1-\tilde{\alpha}}((s_{\hat{g}}(Y^{(k)}, \Phi(X^{(k)}, M^{(\text{test})})))_{k \in \text{Cal}^{m}}) \middle| M^{(\text{test})} = m, \#\text{Cal}^{m} = c\right) \le 1-\alpha + \frac{1}{c+1}$$

This proves the first two points (with respect to Definition D.1) of Theorem 5.4, by observing that  $\Big\{Y^{(\text{test})} \in \hat{C}_{\alpha}(X^{(\text{test})}, M^{(\text{test})})\Big\} = \Big\{s_{\hat{g}}(Y^{(\text{test})}, \Phi(X^{(\text{test})}, M^{(\text{test})})) \leq \hat{Q}_{1-\tilde{\alpha}}\left(\left(s_{\hat{g}}(Y^{(k)}, \Phi(X^{(k)}, M^{(\text{test})}))\right)_{k \in \text{Cal}^m}\right)\Big\}.$ 

Then, the proof of Theorem 5.5 (marginal validity of the CP-MDA-Exact) is direct by marginalizing the result of Theorem 5.4.

#### 

#### D.3 Validities of CP-MDA-Nested.

Next, we give more details on the results on CP-MDA-Nested.

#### D.3.1 MARGINAL VALIDITY OF CP-MDA-NESTED.

The proof of Theorem 5.3 (recalled below) is highly inspired from the Jackknife+ Barber et al. (2021b) proof.

**Theorem 5.3 (marginal validity of CP-MDA-Nested).** If A1 and A2 hold, then, for any missingness mechanism, CP-MDA-Nested outputs intervals  $\hat{C}_{\alpha}$  such that:  $\mathbb{P}(Y^{(n+1)} \in \hat{C}_{\alpha}(X^{(n+1)}, M^{(n+1)})) \ge 1 - 2\alpha$ .

Proof.

$$\begin{split} \left\{Y^{(n+1)} \notin \hat{C}_{\alpha}(X^{(n+1)}, M^{(n+1)})\right\} &= \left\{Y^{(n+1)} > \hat{Q}_{1-\hat{\alpha}}(Z_{1-\frac{\alpha}{2}}) \operatorname{or} Y^{(n+1)} < \hat{Q}_{\hat{\alpha}}(Z_{\frac{\alpha}{2}})\right\} \\ &= \left\{(1-\alpha)(n+1) \leq \sum_{k=1}^{n} \mathbbm{1}\left\{Y^{(n+1)} > Z_{1-\frac{\alpha}{2}}^{(k)}\right\}\right\} \\ &\quad \operatorname{cr} (1-\alpha)(n+1) \leq \sum_{k=1}^{n} \mathbbm{1}\left\{Y^{(n+1)} > Z_{1-\frac{\alpha}{2}}^{(k)} \operatorname{or} Y^{(n+1)} < Z_{\frac{\alpha}{2}}^{(k)}\right\}\right\} \\ &= \left\{\sum_{k=1}^{n} \mathbbm{1}\left\{Y^{(n+1)} > Z_{1-\frac{\alpha}{2}}^{(k)} \operatorname{or} Y^{(n+1)} < Z_{\frac{\alpha}{2}}^{(k)}\right\} \geq (1-\alpha)(n+1)\right\} \\ &\text{using CQR scores for simplicity} \rightarrow = \left\{\sum_{k=1}^{n} \mathbbm{1}\left\{Y^{(n+1)} > \hat{q}_{1-\frac{\alpha}{2}} \circ \Phi\left(X^{(n+1)}, \max(M^{(n+1)}, M^{(k)})\right) + S^{(k),n+1}\right\} \\ &\quad \operatorname{or} Y^{(n+1)} < \hat{q}_{\frac{\alpha}{2}} \circ \Phi\left(X^{(n+1)}, \max(M^{(n+1)}, M^{(k)})\right) - S^{(k),n+1}\right\} \\ &\geq (1-\alpha)(n+1)\right\} \\ &= \left\{\sum_{k=1}^{n} \mathbbm{1}\left\{Y^{(n+1)} - \hat{q}_{1-\frac{\alpha}{2}} \circ \Phi\left(X^{(n+1)}, \max(M^{(n+1)}, M^{(k)})\right) - S^{(k),n+1}\right\} \\ &\geq (1-\alpha)(n+1)\right\} \\ &= \left\{\sum_{k=1}^{n} \mathbbm{1}\left\{\max\left(Y^{(n+1)} - \hat{q}_{1-\frac{\alpha}{2}} \circ \Phi\left(X^{(n+1)}, \max(M^{(n+1)}, M^{(k)})\right) - Y^{(n+1)} > S^{(k),n+1}\right\} \\ &\geq (1-\alpha)(n+1)\right\} \\ &= \left\{\sum_{k=1}^{n} \mathbbm{1}\left\{\max\left(Y^{(n+1)} - \hat{q}_{1-\frac{\alpha}{2}} \circ \Phi\left(X^{(n+1)}, \max(M^{(n+1)}, M^{(k)})\right) - Y^{(n+1)}\right) > S^{(k),n+1}\right\} \\ &\geq (1-\alpha)(n+1)\right\} \\ &\geq (1-\alpha)(n+1)\right\} \\ &\geq (1-\alpha)(n+1)\right\} \end{split}$$

where we defined  $S^{(k),l} = \max\left(Y^{(k)} - \hat{q}_{1-\frac{\alpha}{2}} \circ \Phi\left(X^{(k)}, \max(M^{(k)}, M^{(l)})\right), \hat{q}_{\frac{\alpha}{2}} \circ \Phi\left(X^{(k)}, \max(M^{(k)}, M^{(l)})\right) - Y^{(k)}\right)$ or more generally  $S^{(k),l} = s\left((X^{(k)}, \max(M^{(k)}, M^{(l)})), Y^{(k)}\right)$ , that is the score of the point k when the mask of the point l is applied to it, on top of its own mask  $M^{(k)}$ .

Following Barber et al. (2021b), we now define a comparison matrix  $A \in \{0, 1\}^{(n+1)\times(n+1)}$ , such that for  $(k, l) \in [\![1, n+1]\!]^2$ :  $A_{k,l} = \mathbb{1} \{S^{(k),l} > S^{(l),k}\}$ . Hence, we now have (since by definition  $A_{n+1,n+1} = 0$ ):

$$\left\{Y^{(n+1)} \notin \widehat{C}_{\alpha}(X^{(n+1)}, M^{(n+1)})\right\} \subset \left\{\sum_{k=1}^{n+1} A_{n+1,k} \ge (1-\alpha)(n+1)\right\}.$$

Denote  $W(A) = \{l \in [\![1, n+1]\!] : \sum_{k=1}^{n+1} A_{l,k} \ge (1-\alpha)(n+1)\}$ . We can re-write:  $\left\{Y^{(n+1)} \notin \widehat{C}_{\alpha}(X^{(n+1)}, M^{(n+1)})\right\} \subset \{n+1 \in W(A)\}$ .

 $\text{Therefore } \mathbb{P}\left\{Y^{(n+1)} \notin \widehat{C}_{\alpha}(X^{(n+1)}, M^{(n+1)})\right\} \leq \mathbb{P}\left\{n+1 \in W(A)\right\}. \text{ Thus, we will now bound } \mathbb{P}\left\{n+1 \in W(A)\right\}.$ 

Remark that  $\#W(A) \le 2\alpha(n+1)$  deterministically: this is proven in Barber et al. (2021b) for any comparison matrix.

To conclude the proof, observe that the matrix A can be viewed as the output of a deterministic<sup>6</sup> function C of the exchangeable (by A1) sequence  $(X^{(k)}, M^{(k)}, Y^{(k)})_{k=1}^{n+1}$ :  $A = C\left((X^{(k)}, M^{(k)}, Y^{(k)})_{k=1}^{n+1}\right)$ .

Thus, for any permutation  $\sigma$  on  $[\![1, n+1]\!]$ ,  $C\left(\left(X^{(k)}, M^{(k)}, Y^{(k)}\right)_{k=1}^{n+1}\right) \stackrel{(d)}{=} C\left(\left(X^{(\sigma(k))}, M^{(\sigma(k))}, Y^{(\sigma(k))}\right)_{k=1}^{n+1}\right) := A^{\sigma}.$ 

It follows that for any  $k \in [\![1, n+1]\!]$ ,  $\mathbb{P}\{k \in W(A)\} = \mathbb{P}\{k \in W(A^{\sigma})\}$  for any permutation  $\sigma$ . Therefore  $\mathbb{P}\{k \in W(A)\}$  does not depend on k.

Finally:

$$\begin{split} \mathbb{P}(Y^{(n+1)} \notin \widehat{C}_{\alpha}(X^{(n+1)}, M^{(n+1)})) &\leq \mathbb{P}\{n+1 \in W(A)\} \\ &= \frac{1}{n+1} \sum_{k=1}^{n+1} \mathbb{P}\{k \in W(A)\} \\ &= \frac{1}{n+1} \mathbb{E}[\ \#W(A)] \\ &\leq \frac{1}{n+1} 2\alpha(n+1) \\ \mathbb{P}(Y^{(n+1)} \notin \widehat{C}_{\alpha}(X^{(n+1)}, M^{(n+1)})) &\leq 2\alpha. \end{split}$$

#### D.3.2 MASK-CONDITIONAL-VALIDITY OF CP-MDA-NESTED.

Let  $m \in \mathcal{M}$ .

We start by describing the links between CP-MDA-Nested and CP-MDA-Exact. CP-MDA-Exact can be re-written in the same way as CP-MDA-Nested, but keeping the subselection step of 1. 5.

Indeed, first mention that the output of Algorithm 1 can be written in the following ways:

• 
$$\widehat{C}_{\alpha}(X^{(\text{test})}, m^{(\text{test})}) = \left| \widehat{q}_{\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, m^{(\text{test})}) - \widehat{Q}_{1-\tilde{\alpha}}(S); \widehat{q}_{1-\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, m^{(\text{test})}) + \widehat{Q}_{1-\tilde{\alpha}}(S) \right|$$

• 
$$\widehat{C}_{\alpha}(X^{(\text{test})}, m^{(\text{test})}) = \left[ \widehat{Q}_{\tilde{\alpha}} \left( \hat{q}_{\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, m^{(\text{test})}) - S_{\text{Cal}^{(\text{test})}} \right); \widehat{Q}_{1-\tilde{\alpha}} \left( \hat{q}_{1-\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, m^{(\text{test})}) + S_{\text{Cal}^{(\text{test})}} \right) \right]$$
  
•  $\widehat{C}_{\alpha}(X^{(\text{test})}, m^{(\text{test})}) = \left[ \widehat{Q}_{\tilde{\alpha}} \left( Z_{\frac{\alpha}{2}}^{m^{(\text{test})}} \right); \widehat{Q}_{1-\tilde{\alpha}} \left( Z_{1-\frac{\alpha}{2}}^{m^{(\text{test})}} \right) \right].$ 

With  $Z_{\frac{\alpha}{2}}^m := \{z_{\frac{\alpha}{2}}^{(k)}, k \in \text{Cal and } \widetilde{M}^{(k)} = m\}$ , and similarly for the upper bag. Recall that we have:  $z_{\frac{\alpha}{2}}^{(k)} = \hat{q}_{\frac{\alpha}{2}} \circ \Phi(x^{(\text{test})}, \widetilde{m}^{(k)}) - s^{(k)}$ .

On the other hand, the output predictive interval of Algorithm 2 is then written as:

• 
$$\widehat{C}_{\alpha}\left(X^{(\text{test})}, m^{(\text{test})}\right) = [\widehat{Q}_{\tilde{\alpha}}\left(Z_{\frac{\alpha}{2}}\right); \widehat{Q}_{1-\tilde{\alpha}}\left(Z_{1-\frac{\alpha}{2}}\right)].$$

With these notations,  $Z_{\frac{\alpha}{2}}$  can be partitioned as

$$Z_{\frac{\alpha}{2}} = Z_{\frac{\alpha}{2}}^{m} \bigcup \left( \bigcup_{\widetilde{m}^{(k)} \supset m} Z_{\frac{\alpha}{2}}^{\widetilde{m}^{(k)}} \right).$$
(8)

With

$$\begin{split} & Z_{\frac{\alpha}{2}} = \{ Z_{\frac{\alpha}{2}}^{(k)}, k \in \text{Cal} \} \\ & Z_{\frac{\alpha}{2}}^{(k)} = \hat{q}_{\frac{\alpha}{2}} \circ \Phi\left( X^{(\text{test})}, \widetilde{M}^{(k)} \right) - S^{(k)} \\ & s^{(k)} = \max(\hat{q}_{\frac{\alpha}{2}}(x_{\text{imp}}^{(k)}) - y^{(k)}, y^{(k)} - \hat{q}_{1-\frac{\alpha}{2}}(x_{\text{imp}}^{(k)})) \end{split}$$

<sup>&</sup>lt;sup>6</sup>In fact, *C* is only required to be "independent of its input". If we denote by  $\xi$  the randomness of *C*, we can write that  $C(\cdot) = \mathcal{D}(\cdot; \xi)$  with  $\mathcal{D}$  a deterministic function and the requirement is that  $\xi$  and the argument of *C* are independent.

The result of Algorithm 1 implies that for any mask  $m \in \mathcal{M}$ , we have :

$$\mathbb{P}\left(Y^{(\text{test})} \in \widehat{C}_{\alpha}\left(X^{(\text{test})}, m\right) | M^{(\text{test})} = m\right) \ge 1 - \alpha,$$

i.e.

$$\mathbb{P}\left(Y^{(\text{test})} \notin \left[\hat{q}_{\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, m) - \hat{Q}_{1-\tilde{\alpha}}\left(S^{m}\right); \hat{q}_{1-\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, m) + \hat{Q}_{1-\tilde{\alpha}}\left(S^{m}\right)\right] | M^{(\text{test})} = m\right) \leq \alpha.$$
(9)

Where:  $Q_{1-\tilde{\alpha}}(S)$  is the  $(1-\alpha)(1+1/\#S)$ -quantile of S and  $S^m = \{s^{(k)} \text{ for } k \in \text{Cal and } \widetilde{M}^{(k)} = m\}$ . Equivalently:

$$\mathbb{P}\left(Y^{(\text{test})} \in \left[\widehat{Q}_{\tilde{\alpha}}\left(Z_{\frac{\alpha}{2}}^{m}\right); \widehat{Q}_{1-\tilde{\alpha}}\left(Z_{1-\frac{\alpha}{2}}^{m}\right)\right] | M^{(\text{test})} = m\right) \ge 1 - \alpha.$$
(10)

In the following Lemma, we show that for  $\tilde{m} \supset m$  the result extends under Assumption A4.

**Lemma D.3.** Assume Assumption A4. For any  $m \in M$ , for any  $\tilde{m} \supset m$ 

$$\mathbb{P}\left[\left(Y^{(\text{test})} \in \left[\widehat{Q}_{\tilde{\alpha}}\left(Z_{\frac{\alpha}{2}}^{\tilde{m}}\right); \widehat{Q}_{1-\tilde{\alpha}}\left(Z_{1-\frac{\alpha}{2}}^{\tilde{m}}\right)\right]\right) | M^{(\text{test})} = m\right] \ge 1 - \alpha.$$
(11)

This inequality shows the conservativeness of the quantiles of the bags resulting from larger missing patterns  $\tilde{m}$  than m when the construction of the output of Algorithm 2.

While inequality Equation (10) is "tight" in the sense that the probability is almost exactly  $1 - \alpha$  (item 2 of Theorem 5.4), the proof hereafter shows that Equation (11) can be pessimistic in terms of actual coverage, as one may have  $\mathbb{P}[(Y^{(\text{test})}\notin [\hat{Q}_{\tilde{\alpha}}(Z_{\frac{\tilde{m}}{2}}^{\tilde{m}}); \hat{Q}_{1-\tilde{\alpha}}(Z_{1-\frac{\alpha}{2}}^{\tilde{m}})])|M^{(\text{test})} = m] \ll \alpha.$ 

More precisely, we have the following inequality:

$$\mathbb{E}\left[\mathbb{P}\left(Y^{(\text{test})}\notin\left[\hat{q}_{\frac{\alpha}{2}}\circ\Phi(X^{(\text{test})},\tilde{m})-\hat{Q}_{1-\tilde{\alpha}}\left(S^{\tilde{m}}\right);\hat{q}_{1-\frac{\alpha}{2}}\circ\Phi(X^{(\text{test})},\tilde{m})+\hat{Q}_{1-\tilde{\alpha}}\left(S^{\tilde{m}}\right)\right]\middle|M^{(\text{test})}=m,X^{(\text{test})}_{\text{obs(m)}}\right)\middle|M^{(\text{test})}=m\right]\leq\alpha$$
(12)

The interpretation of that Lemma is that the intervals resulting from the prediction on  $x^{\text{test}}$ ,  $\tilde{m}$  (more data hidden) and corrected with the residuals of the calibration points  $(X^k, M^k = \tilde{m}, Y^k)$  have a *larger* probability of containing  $Y^{\text{test}}$ , conditionally to  $X_{\text{obs}(m)}$  than the interval built using prediction on  $x^{\text{test}}$ , m (more data available) and corrected with the residuals of the calibration points  $(X^k, M^k = m, Y^k)$  (more data available) and corrected with the residuals of the calibration points  $(X^k, M^k = m, Y^k)$  (more data available)

*Proof of Lemma D.3.* We start by invoking Equation (9) for  $\tilde{m}$ :

$$\mathbb{P}\left(Y^{(\text{test})} \notin \left[\hat{q}_{\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, \tilde{m}) - \hat{Q}_{1-\tilde{\alpha}}\left(S^{\tilde{m}}\right); \hat{q}_{1-\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, \tilde{m}) + \hat{Q}_{1-\tilde{\alpha}}\left(S^{\tilde{m}}\right)\right] | M^{(\text{test})} = \tilde{m}\right) \leq \alpha.$$
(13)

Consequently, by the tower property of conditional expectations:

$$\mathbb{E}\left[\mathbb{P}\left(Y^{(\text{test})} \notin \left[\hat{q}_{\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, \tilde{m}) - \hat{Q}_{1-\tilde{\alpha}}\left(S^{\tilde{m}}\right); \hat{q}_{1-\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, \tilde{m}) + \hat{Q}_{1-\tilde{\alpha}}\left(S^{\tilde{m}}\right)\right] \middle| M^{(\text{test})} = \tilde{m}, S^{(\tilde{m})}, X^{(\text{test})}_{\text{obs}(\tilde{m})}\right) \middle| M^{(\text{test})} = \tilde{m}\right] \leq \alpha$$

$$(14)$$

Observe that  $\hat{q}_{\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, \tilde{m}) - \hat{Q}_{1-\tilde{\alpha}}\left(S^{\tilde{m}}\right)$  is  $\{M^{(\text{test})} = \tilde{m}, S^{(\tilde{m})}, X^{(\text{test})}_{\text{obs}(\tilde{m})}\}$ -measurable.

Moreover, by Assumption A4, we have that for any  $\delta \in [0, 0.5]$ :

$$q_{1-\delta/2}^{Y|(X_{\text{obs}(m)}, M=m)} \le q_{1-\delta/2}^{Y|(X_{\text{obs}(\tilde{m})}, M=\tilde{m})}$$
(15)

$$q_{\delta/2}^{Y|(X_{\text{obs}(m)},M=m)} \ge q_{\delta/2}^{Y|(X_{\text{obs}(\tilde{m})},M=\tilde{m})}.$$
(16)

In other words the conditional distribution of Y given  $X_{obs(\tilde{m})}$  and  $M = \tilde{m}$  "stochastically dominates" the conditional distribution of Y given  $X_{obs(m)}$  and M = m.

We thus have, with  $F_Z$  denoting the cumulative distribution function of Z:  $F_{Y|(X_{obs(\tilde{m})}, M=\tilde{m})}$  the cumulative distribution function of  $Y|(X_{obs(\tilde{m})}, M=\tilde{m})$ :

$$\mathbb{P}\left(Y^{(\text{test})} \notin \left[\hat{q}_{\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, \tilde{m}) - \hat{Q}_{1-\tilde{\alpha}}\left(S^{\tilde{m}}\right); \hat{q}_{1-\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, \tilde{m}) + \hat{Q}_{1-\tilde{\alpha}}\left(S^{\tilde{m}}\right)\right] \left| M^{(\text{test})} = \tilde{m}, S^{(\tilde{m})}, X^{(\text{test})}_{\text{obs}(\tilde{m})} \right) \\ = 1 - \left[F_{Y|(X_{\text{obs}(\tilde{m})}, M=\tilde{m})}\left(\hat{q}_{1-\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, \tilde{m}) + \hat{Q}_{1-\tilde{\alpha}}(S^{\tilde{m}})\right) - F_{Y|(X_{\text{obs}(\tilde{m})}, M=\tilde{m})}\left(\hat{q}_{\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, \tilde{m}) - \hat{Q}_{1-\tilde{\alpha}}(S^{\tilde{m}})\right)\right] \\ \stackrel{(i)}{\geq} 1 - \left[F_{Y|(X_{\text{obs}(m)}, M=m)}\left(\hat{q}_{1-\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, \tilde{m}) + \hat{Q}_{1-\tilde{\alpha}}(S^{\tilde{m}})\right) - F_{Y|(X_{\text{obs}(m)}, M=m)}\left(\hat{q}_{\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, \tilde{m}) - \hat{Q}_{1-\tilde{\alpha}}(S^{\tilde{m}})\right)\right] \\ = \mathbb{P}\left(Y^{(\text{test})}\notin \left[\hat{q}_{\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, \tilde{m}) - \hat{Q}_{1-\tilde{\alpha}}\left(S^{\tilde{m}}\right); \hat{q}_{1-\frac{\alpha}{2}} \circ \Phi(X^{(\text{test})}, \tilde{m}) + \hat{Q}_{1-\tilde{\alpha}}\left(S^{\tilde{m}}\right)\right]\right| M^{(\text{test})} = m, S^{(\tilde{m})}, X^{(\text{test})}_{\text{obs}(m)}\right).$$
(17)

At (i) we use (16)  $F_{Y|(X_{obs(m)},M=m)}(\hat{q}_{\frac{\alpha}{2}} \circ \Phi(X^{(test)},\tilde{m}) - \hat{Q}_{1-\tilde{\alpha}}(S^{\tilde{m}})) \leq F_{Y|(X_{obs(\tilde{m})},M=\tilde{m})}(\hat{q}_{\frac{\alpha}{2}} \circ \Phi(X^{(test)},\tilde{m}) - \hat{Q}_{1-\tilde{\alpha}}(S^{\tilde{m}})),$ and (15):  $F_{Y|(X_{obs(m)},M=m)}(\hat{q}_{1-\frac{\alpha}{2}} \circ \Phi(X^{(test)},\tilde{m}) + \hat{Q}_{1-\tilde{\alpha}}(S^{\tilde{m}})) \geq F_{Y|(X_{obs(\tilde{m})},M=\tilde{m})}(\hat{q}_{1-\frac{\alpha}{2}} \circ \Phi(X^{(test)},\tilde{m}) + \hat{Q}_{1-\tilde{\alpha}}(S^{\tilde{m}}))$  by A4. Remark that here we assume that  $\left(\hat{q}_{1-\frac{\alpha}{2}} \circ \Phi(X^{(test)},\tilde{m}) + \hat{Q}_{1-\tilde{\alpha}}(S^{\tilde{m}})\right) \geq \operatorname{median}(Y^{(test)}|(X^{(test)}_{obs(\tilde{m})},M=\tilde{m}))$  and  $\left(\hat{q}_{\frac{\alpha}{2}} \circ \Phi(X^{(test)},\tilde{m}) - \hat{Q}_{1-\tilde{\alpha}}(S^{\tilde{m}})\right) \leq \operatorname{median}(Y^{test}|(X^{(test)}_{obs(\tilde{m})},M=\tilde{m})).$ 

We obtain Equation (12) in Lemma D.3 by plugging (17) in (14), then Equation (11) by the tower property.

**Theorem D.4.** Assume the missing mechanism is MCAR, and that Assumptions A1 to A3 hold. Additionally Assumption A4 is satisfied.

Consider the partition described in Equation (8), and consider CP-MDA-Nested running on a test point with missing pattern  $m^{(\text{test})}$ , with any of the following outputs, instead of l. 15  $\hat{C}_{\alpha} \left( x^{(\text{test})}, m^{(\text{test})} \right) = [\hat{Q}_{\tilde{\alpha}} \left( Z_{\frac{\alpha}{2}} \right); \hat{Q}_{1-\tilde{\alpha}} \left( Z_{1-\frac{\alpha}{2}} \right)]$ :

- $I. \ \widehat{C}_{\alpha}\left(x^{(\textit{test})}, m^{(\textit{test})}\right) = [\widehat{Q}_{\tilde{\alpha}}(Z_{\underline{\alpha}}^{\tilde{m}}); \widehat{Q}_{1-\tilde{\alpha}}(Z_{1-\underline{\alpha}}^{\tilde{m}})] \text{ where } \tilde{m} \supset m^{(\text{test})} \text{ is an arbitrary choice.}$
- 2.  $\widehat{C}_{\alpha}\left(x^{(test)}, m^{(test)}\right) = [\widehat{Q}_{\tilde{\alpha}}(Z_{\frac{\tilde{m}}{2}}^{\hat{m}}); \widehat{Q}_{1-\tilde{\alpha}}(Z_{1-\frac{\tilde{m}}{2}}^{\hat{m}})]$  where  $\hat{m}$  is a randomly selected pattern in  $\{\tilde{m}, \tilde{m} \supset m^{(test)}\}$ , possibly with varying probability depending on the cardinality of the sets  $Z_{\alpha/2}^{\tilde{m}}$ .

Then the resulting algorithm is mask-conditionally-valid.

<u>Proof of Theorem D.4</u>. The proof immediately follows from Equation (11), and gives the result without difficulty for any arbitrary pattern or random variable independent of all other randomness.

Extension to a choice that involves the cardinality of the sets  $Z_{\alpha/2}^{\tilde{m}}$ , leveraging the independence between these cardinals and the coverage properties (same as in the proof of Theorem D.2).

#### **E** Infinite data results

**Proposition 6.1** ( $\ell_{\beta}$ -consistency of an universal learner). Let  $\beta \in [0, 1]$ . If X admits a density on X, then, for almost all imputation function  $\Phi \in \mathcal{F}_{\infty}^{I}$ , the function  $g_{\ell_{\beta}, \Phi}^{*} \circ \Phi$  is Bayes optimal for the pinball risk of level  $\beta$ .

*Proof of Proposition 6.1.* The proof starts in the exact same way than Le Morvan et al. (2021), based on their Lemmas A.1 and A.2. For completeness, we copy here the statements of these lemmas without their proof and rewrite the two first parts of the main proof.

Let  $\Phi$  be an imputation function such that for each missing data pattern  $m, \phi^m \in \mathcal{C}^{\infty}(\mathbb{R}^{|obs(m)|}, \mathbb{R}^{|mis(m)|})$ .

**Lemma E.1** (Lemma A.1 in Le Morvan et al. (2021)). Let  $\phi^m \in \mathcal{C}^{\infty}(\mathbb{R}^{|\operatorname{obs}(m)|}, \mathbb{R}^{|\operatorname{mis}(m)|})$  be the imputation function for missing data pattern m, and let  $\mathcal{M}^m = \{x \in \mathbb{R}^d : x_{\operatorname{mis}(m)} = \phi^m(x_{\operatorname{obs}((m))})\}$ . For all m,  $\mathcal{M}^m$  is an  $|\operatorname{obs}((m))|$ -dimensional manifold.

In Lemma E.1,  $\mathcal{M}^m$  represents the manifold in which the data points are sent once imputed by  $\phi^m$ . Lemma E.1 states that this manifold is of dimension |obs(m)|.

**Lemma E.2** (Lemma A.2 in Le Morvan et al. (2021)). Let m and m' be two distinct missing data patterns with the same number of missing (resp. observed) values |mis| (resp |obs|). Let  $\phi^m \in C^{\infty}(\mathbb{R}^{|obs(m)|}, \mathbb{R}^{|mis(m)|})$  be the imputation function for missing data pattern m, and let  $\mathcal{M}^m = \{x \in \mathbb{R}^d : x_{mis(m)} = \phi^m(x_{obs(m)})\}$ . We define similarly  $\Phi^{(m')}$  and  $\mathcal{M}^{(m')}$ . For almost all imputation functions  $\phi^m$  and  $\Phi^{(m')}$ ,

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

Note that, as by Lemma E.1 dim  $(\mathcal{M}^m) = \dim \left(\mathcal{M}^{(m')}\right) = |\mathrm{obs}| = \mathrm{d} - |\mathrm{mis}|$ , Lemma E.2 states that dim  $\left(\mathcal{M}^m \cap \mathcal{M}^{(m')}\right) \leq \dim \left(\mathcal{M}^m\right) = \dim \left(\mathcal{M}^{(m')}\right)$ .

Now, to prove Proposition 6.1 the missing data patterns are ordered as in Le Morvan et al. (2021): the first one will be the one in which all the variables are missing, while the last one will be the one in which all the variables are observed. For two data patterns with the same number of missing variables, the ordering is picked at random. We denote by m(i) the *i*-th missing data pattern according to this ordering.

We are going to build a function  $g_{\Phi}$  which, composed with  $\Phi$ , will reach the  $\ell$ -Bayes risk.

For each missing data pattern, and starting by m(1) of all variables missing, we can define  $g_{\Phi}$  on the data points from the current missing data pattern. More precisely, for each i,  $g_{\Phi}$  is built for every imputed data point belonging to  $\mathcal{M}^{(m(i))}$  except for those already considered in previous steps (one imputed data point can belong to multiple manifolds):

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

That is,  $g_{\Phi} \circ \Phi(X, M)$  will equal  $\tilde{f}^*(X, M)$  except possibly if  $\Phi(X, M) = \Phi(\tilde{Y})$  for some  $\tilde{Y}$  that has more missing values than X, M. Therefore, for each missing data pattern  $m(i), g_{\Phi} \circ \Phi$  equals  $\tilde{f}^*$  except on  $\bigcup_{k < i} \mathcal{M}^{(m(k))}$ . The question that remains is: what is the dimension of  $\mathcal{M}^{(m(i))} \cap (\bigcup_{k < i} \mathcal{M}^{(m(k))})$ , these points for which there is no necessarily equality between  $g_{\Phi} \circ \Phi$  and  $\tilde{f}^*$ . First, note that  $\mathcal{M}^{(m(i))} \cap (\bigcup_{k < i} \mathcal{M}^{(m(k))}) = \bigcup_{k < i} (\mathcal{M}^{(m(i))} \cap \mathcal{M}^{(m(k))})$ . For each space in this reunion, there are two cases:

- either  $|\operatorname{obs}(\mathbf{m}(\mathbf{k}))| < |\operatorname{obs}(\mathbf{m}(\mathbf{i}))|$ : using Lemma E.1,  $\dim (\mathcal{M}^{(m(k))}) = |\operatorname{obs}(\mathbf{m}(\mathbf{k}))| < |\operatorname{obs}(\mathbf{m}(\mathbf{i}))| = \dim (\mathcal{M}^{(\mathbf{m}(\mathbf{i}))})$ . Thus,  $\mathcal{M}^{(m(i))} \cap \mathcal{M}^{(m(k))}$  is of measure zero in  $\mathcal{M}^{(m(i))}$ .
- either |obs(m(k))| = |obs(m(i))|: using Lemma E.2,  $\mathcal{M}^{(m(i))} \cap \mathcal{M}^{(m(k))}$  is of dimension 0 or smaller than  $\dim (\mathcal{M}^{(m(i))})$ , thus it is of measure zero in  $\mathcal{M}^{(m(i))}$ .

Therefore, the set of data points for which  $g_{\Phi} \circ \Phi$  does not equal the oracle is of measure 0 for each missing data pattern. Let  $\beta \in [0, 1]$ . We can now write down the  $\ell_{\beta}$ -risk of this built function:

$$\begin{split} \mathbb{E}\left[\ell_{\beta}\left(Y, g^{*} \circ \Phi(X, M)\right)\right] &= \mathbb{E}\left[\rho_{\beta}\left(Y - g^{*} \circ \Phi(X, M)\right)\right] \\ &= \mathbb{E}\left[\rho_{\beta}\left(Y - \tilde{f}^{*}(X, M) + \tilde{f}^{*}(X, M) - g^{*} \circ \Phi(X, M)\right)\right] \\ &(i) \leq \mathbb{E}\left[\rho_{\beta}\left(Y - \tilde{f}^{*}(X, M)\right)\right] + \mathbb{E}\left[\rho_{\beta}\left(\tilde{f}^{*}(X, M) - g^{*} \circ \Phi(X, M)\right)\right] \\ &\leq \mathcal{R}_{\ell_{\beta}}^{*} + \mathbb{E}\left[\rho_{\beta}\left(\tilde{f}^{*}(X, M) - g^{*} \circ \Phi(X, M)\right)\right], \end{split}$$

where (i) holds thanks to the shape of  $\rho_{\beta}$ . For any  $w \in \mathbb{R}$  and any  $\lambda \in \mathbb{R}_+$ :

$$\rho_{\beta} (\lambda w) = \beta \lambda |w| \mathbb{1}_{w \ge 0} + (1 - \beta) \lambda |w| \mathbb{1}_{w \le 0}$$
  
$$\rho_{\beta} (\lambda w) = \lambda \rho_{\beta} (w) .$$

Furthermore,  $\rho_{\beta}$  is convex, thus for any  $(u, v) \in \mathbb{R}^2$ :

$$\rho_{\beta}\left(\frac{1}{2}u+\frac{1}{2}v\right) \leq \frac{1}{2}\rho_{\beta}(u)+\frac{1}{2}\rho_{\beta}(v)$$
$$\frac{1}{2}\rho_{\beta}(u+v) \leq \frac{1}{2}\rho_{\beta}(u)+\frac{1}{2}\rho_{\beta}(v)$$
$$\rho_{\beta}(u+v) \leq \rho_{\beta}(u)+\rho_{\beta}(v).$$

As  $\tilde{f}^*$  and  $g^* \circ \Phi$  are equals almost everywhere on each missing subspace,  $\mathbb{E}\left[\rho_\beta\left(\tilde{f}^*(X,M) - g^* \circ \Phi(X,M)\right)\right] = 0$ . Indeed, decomposing by pattern one can write:

$$\mathbb{E}\left[\rho_{\beta}\left(\tilde{f}^{*}(X,M) - g^{*} \circ \Phi(X,M)\right)\right] = \sum_{M=m} \mathbb{P}(M=m)\mathbb{E}\left[\rho_{\beta}\left(\tilde{f}^{*}(X,M) - g^{*} \circ \Phi(X,M)\right) | M=m\right]$$

and thus by equality almost everywhere for each pattern every term in this sum is null.

Therefore one obtains:

$$\mathbb{E}\left[\ell_{\beta}\left(Y, g^{*} \circ \Phi(X, M)\right)\right] \leq \mathcal{R}_{\ell_{\beta}}^{*}$$

Thus:

$$\mathbb{E}\left[\ell_{\beta}\left(Y, g^{*} \circ \Phi(X, M)\right)\right] = \mathcal{R}_{\ell_{\beta}}^{*},$$

and  $g^* \circ \Phi$  is Bayes optimal. This implies that  $\mathcal{R}^*_{\ell_{\beta},\Phi} = \mathcal{R}^*_{\ell_{\beta}}$ . Thus, a universally consistent algorithm learning  $g_{\Phi}$  chained with  $\Phi$  will lead to a Bayes consistent function.

Proof of Corollary 6.2. Corollary 6.2 states that "For any missing mechanism, for almost all imputation function  $\Phi \in \mathcal{F}_{\infty}^{I}$ , if  $F_{Y|(X_{obs(M)},M)}$  is continuous, a universally consistent quantile regressor trained on the imputed data set yields asymptotic conditional coverage."

Let  $\beta \in [0, 1]$ .

Remark that Proposition 6.1 states that for any missing mechanism, for almost all imputation function  $\Phi \in \mathcal{F}_{\Delta}^{I}$ a universally consistent quantile regressor trained on the imputed data set achieves the Bayes risk asymptotically. We will thus show that any  $\ell_{\beta}$ -Bayes predictor  $f_{\beta}^{*}$  (any function achieving the  $\ell_{\beta}$ -Bayes-risk) is such that  $\mathbb{P}(Y \leq f_{\beta}^{*}(X, M)|X_{\text{obs}(M)}, M) = \beta$  if  $F_{Y|(X_{\text{obs}(M)}, M)}$  is continuous. Therefore, any two Bayes predictors  $f_{\alpha/2}^{*}$  and  $f_{1-\alpha/2}^{*}$  form an interval  $[f_{\alpha/2}^{*}(X, M); f_{1-\alpha/2}^{*}(X, M)]$  that achieves conditional coverage (conditionally to  $X_{\text{obs}(M)}$  and M).

Let  $f_{\beta}^*$  be a  $\ell_{\beta}$ -Bayes predictor. Then:

$$f_{\beta}^{*} \in \underset{f:\mathcal{X}\times\mathcal{M}\to\mathbb{R}}{\operatorname{arg\,min}} \mathbb{E}\left[\rho_{\beta}\left(Y-f\left(X,M\right)\right)\right]$$
$$=\mathbb{E}\left[\mathbb{E}\left[\rho_{\beta}\left(Y-f\left(X,M\right)\right)|X_{\operatorname{obs}(M)},M\right]\right].$$

Let  $(x,m) \in \mathcal{X} \times \mathcal{M}$ . Denote  $H_{x,m}(z) := \mathbb{E}\left[\rho_{\beta}\left(Y-z\right) | X_{\text{obs}(M)} = x_{\text{obs}(m)}, M = m\right]$ . As  $Y \neq z$  almost surely, we have:

$$\begin{aligned} H'_{x,m}(z) &= \mathbb{E} \left[ -\rho'_{\beta} \left( Y - z \right) | X_{\text{obs}(M)} = x_{\text{obs}(m)}, M = m \right] \\ &= \mathbb{E} \left[ -(-\beta \mathbb{1}_{Y-z \ge 0} + (1-\beta) \mathbb{1}_{Y-z \le 0}) | X_{\text{obs}(M)} = x_{\text{obs}(m)}, M = m \right] \\ &= \mathbb{E} \left[ \beta \mathbb{1}_{Y \ge z} - (1-\beta) \mathbb{1}_{Y \le z} | X_{\text{obs}(M)} = x_{\text{obs}(m)}, M = m \right] \\ &= \beta \mathbb{P} \left( Y \ge z | X_{\text{obs}(M)} = x_{\text{obs}(m)}, M = m \right) - (1-\beta) \mathbb{P} \left( Y \le z | X_{\text{obs}(M)} = x_{\text{obs}(m)}, M = m \right) \\ &= \beta \left( 1 - \mathbb{P} \left( Y \le z | X_{\text{obs}(M)} = x_{\text{obs}(m)}, M = m \right) \right) - (1-\beta) \mathbb{P} \left( Y \le z | X_{\text{obs}(M)} = x_{\text{obs}(m)}, M = m \right) \\ H'_{x,m}(z) &= \beta - \mathbb{P} \left( Y \le z | X_{\text{obs}(M)} = x_{\text{obs}(m)}, M = m \right). \end{aligned}$$

Therefore  $H'_{x,m}(z) \leq 0$  if and only if  $\beta \leq \mathbb{P}\left(Y \leq z | X_{\text{obs}(M)} = x_{\text{obs}(m)}, M = m\right)$ .

Thus, z minimizes  $H_{x,m}$  if and only if  $\beta = \mathbb{P}\left(Y \leq z | X_{\text{obs}(M)} = x_{\text{obs}(m)}, M = m\right)$ .

If  $F_{Y|(X_{obs(M)},M)}$  is continuous, there exists at least a solution, that might not be unique if it is not additionally strictly increasing. Therefore, if  $F_{Y|(X_{obs(M)},M)}$  is continuous, all the  $\ell_{\beta}$ -Bayes predictors can be written as  $f_{\beta}^{*}(x,m) = q_{x,m}$  with  $\mathbb{P}\left(Y \leq q_{x,m}|X_{obs(M)} = x_{obs(m)}, M = m\right) = \mathbb{P}\left(Y \leq f_{\beta}^{*}(x,m)|X_{obs(M)} = x_{obs(m)}, M = m\right) = \beta$ .

### F Experimental study

#### F.1 Settings detail

**Quantile Neural Network.** The architecture and optimization of the Quantile Neural Network used in the experiments is taken from Sesia & Romano (2021) (their code is freely available). This is the description provided in the original paper of the neural network: "The network is composed of three fully connected layers with a hidden dimension of 64, and ReLU activation functions. We use the pinball loss to estimate the conditional quantiles, with a dropout regularization of rate 0.1. The network is optimized using Adam Kingma & Ba (2014) with a learning rate equal to 0.0005. We tune the optimal number of epochs by cross validation, minimizing the loss function on the hold-out data points; the maximal number of epochs is set to 2000."

#### F.2 Gaussian linear results

Figure 8 is the analogous of Figure 1, but by evaluating the performances on the group leading to the highest coverage.



Figure 8: Coverage and interval's length for the easiest pattern. Model is NN. Calibration size fixed to 1000. The mask is concatenated in the features. Data is imputed using Iterative Ridge. 100 repetitions allow to display error bars, corresponding to standard error.

Hereafter, we present in Figure 9 the exact same figure than Figure 3 but with a panel (the first) for vanilla QR. The 3 other methods are displayed again to facilitate the comparison.

Finally, Figure 10 is the analogous of Figure 9, but for a training set containing 1000 observations and a calibration set containing 500 observations.

#### F.3 Semi-synthetic

In the smi-synthetic experiments, two settings are examined: one where the training size is small in comparison to the number of parameters of the Neural Network – "Medium" –, and one where the training size is even smaller so that some masks have a really low (or null) frequency of appearance in the training set – "Small". In both cases, the calibration size is approximately half the training size. Figure 4 presented the results for the "Medium" case.

Table 1: Semi-synthetic setting	s: training and calibration	n sizes for each of the 6	data sets depending on the setting
---------------------------------	-----------------------------	---------------------------	------------------------------------

		$ \begin{array}{c} \texttt{meps\_19} \\ d = 139, l = 5 \\ n = 15785 \end{array} $	meps_20 d = 139, l = 5 n = 17541	$ \begin{vmatrix} \text{meps}_{21} \\ d = 139, l = 5 \\ n = 15656 \end{vmatrix} $	bio d = 9, l = 9 n = 45730	bike d = 18, l = 4 n = 10886	$ \begin{array}{c} \text{concrete} \\ d = 8, l = 8 \\ n = 1030 \end{array} $
Small	Tr size	500	500	500	500	500	330
	Cal size	250	250	250	250	250	100
Medium	Tr size	1000	1000	1000	1000	1000	630
	Cal size	500	500	500	500	500	200



Figure 9: Average coverage (top) and length (bottom) as a function of the pattern size, i.e. the number of missing values (NA). First violin plot corresponds to marginal coverage. Stars correspond to the oracle length. Settings are: model is NN, train size is 500, calibration size is 250. The marginal test set includes 2000 observations. The conditional test set includes 100 individuals for each possible missing data pattern size. The mask is concatenated to the features. Data is imputed using Iterative Ridge. 100 repetitions are performed.



Figure 10: Model is NN. Train size is 1000. Calibration size fixed to 500. The marginal test set includes 2000 observations. The conditional test set includes 100 individuals for each possible missing data pattern size. The mask is concatenated in the features. Data is imputed using Iterative Ridge. 100 repetitions are performed.

Figure 11 respresents the results for both of these settings, using the same parameters than Figure 4. For the results on the two other meps data sets (meps\_20 and meps\_21) see Figure 12, which repeats the visualisation of meps\_19 to ease comparison.

#### F.4 Real data

**Data set description.** Sportisse et al. (2020) selected 7 variables to model the level of platelets, after discussion with medical doctors. Thus, we followed their pipeline. Here are the 7 variables used:

• Age: the age of the patient (no missing values);



Figure 11: Model is NN. The mask is concatenated in the features. Data is imputed using Iterative Ridge. 100 repetitions are performed, allowing to display the standard error as error bars. The vertical dotted lines represent the target coverage of 90%.



Figure 12: Model is NN. The mask is concatenated in the features. Data is imputed using Iterative Ridge. 100 repetitions are performed, allowing to display the standard error as error bars. The vertical dotted lines represent the target coverage of 90%.

- Lactate: the conjugate base of lactic acid, upon arrival at the hospital (17.66% missing values);
- Delta\_hemo: the difference between the hemoglobin upon arrival at hospital and the one in the ambulance (23.82% missing values);
- VE: binary variable indicating if a Volume Expander was applied in the ambulance. A volume expander is a type of intravenous therapy that has the function of providing volume for the circulatory system (2.46% missing values);
- RBC: a binary index which indicates whether the transfusion of Red Blood Cells Concentrates is performed (0.37% missing values);
- SI: the shock index. It indicates the level of occult shock based on heart rate (HR) and systolic blood pressure (SBP), that is  $SI = \frac{HR}{SBP}$ , upon arrival at hospital (2.09% missing values);
- HR: the heart rate measured upon arrival of hospital (1.62% missing values).

**Splitting strategy.** To study the coverage conditionally on the masks, we must handle the scarcity of some of them. For each individual in the data set, we make only one prediction, this way avoiding too many repetitions of the same test point when

computing the average. We split the data set into 5 folds, and predict on each fold by training the procedure on the 4 others, with 15390 observations for training, and 7694 for calibration.



Figure 13: Average coverage and length on the TraumaBase® data when predicting the platelets level. Colors correspond to the methods. Diamond ( $\blacklozenge$ ) corresponds to taking the average among all individuals. Other symbols correspond to computing the average among the individuals having a fixed mask. The vertical dotted line represents the target coverage of 90%. Model is NN. The mask is concatenated to the features. Imputation is Iterative Ridge. Each individual is predicted using 15390 observations for training, and 7694 for calibration.