

1 We thank the reviewers for their time and relevant comments.

2 **Why we use our architecture to approximate the self-masking Bayes predictor - to rev. #3 and #4:** Reviewer #4  
 3 mentioned a lack of clarity in the paragraph entitled "Differences between MNAR and M(C)AR predictors" (l.128-133).  
 4 This paragraph actually touches upon a point raised by Reviewer #3: why our method can be used for self-masking  
 5 missingness (also discussed l.190-192). These questions suggest that this paragraph should be more detailed. Hereafter,  
 6 we explain it in more details and give an answer to reviewer #3.

7 The expression of the M(C)AR Bayes predictors is given by (eq. 4 in the paper):

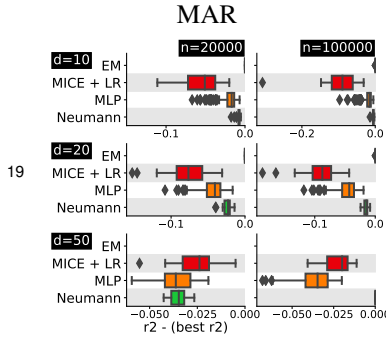
$$f^*(X_{obs}, M) = \beta_0^* + \langle \beta_{obs}^*, X_{obs} \rangle + \langle \beta_{mis}^*, \mu_{mis} + \Sigma_{mis,obs} (\Sigma_{obs})^{-1} (X_{obs} - \mu_{obs}) \rangle \quad (1)$$

8 The expression of the (MNAR) self-masking Bayes predictor is more complicated (eq. 5 in the paper). To study this  
 9 expression, we approximate  $D_{mis} \Sigma_{mis|obs}^{-1}$  by  $Id$ . Then, the self-masking Bayes predictor becomes:

$$f^*(X_{obs}, M) \approx \beta_0^* + \langle \beta_{obs}^*, X_{obs} \rangle + \langle \beta_{mis}^*, \frac{1}{2}(\tilde{\mu}_{mis} + \mu_{mis}) + \frac{1}{2} \Sigma_{mis,obs} (\Sigma_{obs})^{-1} (X_{obs} - \mu_{obs}) \rangle \quad (2)$$

10 Thus, under this approximation, the self-masking Bayes predictor can be modeled by our proposed architecture (just as  
 11 the M(C)AR Bayes predictor), the only difference being the targeted value for parameter  $\mu$  of the network (in blue  
 12 in the two models above) and a scaling factor of 1/2 for  $W_{mix}$  (in orange). A less coarse approximation also works:  
 13  $D_{mis} \Sigma_{mis|obs}^{-1} \approx \hat{D}_{mis}$  where  $\hat{D}$  is a diagonal matrix. In this case, the proposed architecture can perfectly model  
 14 the self-masking Bayes predictor: the parameter  $\mu$  of the network should target  $(Id + \hat{D})^{-1}(\tilde{\mu} + \hat{D}\mu)$  and  $W_{mix}$   
 15 should target  $(Id + \hat{D})^{-1} \hat{D} \Sigma$  instead of simply  $\Sigma$  in the M(C)AR case. Consequently, **our architecture can well**  
 16 **approximate the self-masking Bayes predictor by adjusting the values learned for the parameters  $\mu$  and  $W_{mix}$**   
 17 **if  $D_{mis} \Sigma_{mis|obs}^{-1}$  are close to diagonal matrices.** We will add this discussion to the Appendix.

18 **Experimental results under the MAR scenario - to rev.#3** The results are presented in the figure below:



The MAR data was generated as follows: first, a subset of variables with *no* missing values is randomly selected (10%). The remaining variables have missing values according to a logistic model with random weights, re-scaled so as to attain the desired proportion of missing values on those variables (50%). Note that we cannot compute the Bayes rate, so instead of showing (R2 - Bayes rate) we show (R2 - best R2).

As can be seen from the figure, the trends observed for MAR are the same as those for MCAR. We will add this figure to the appendix. We have not tested the scenario where MAR and self-masking missingness happen at the same time.

20 **Modeling non linear functions - to rev. #1 and #2** The reviewers pointed out that the proposed architecture is  
 21 limited to linear models. Indeed, our theoretical foundations derive from the study of linear models. However, as a  
 22 differentiable architecture, it can be readily included as a building block in more complex networks. For example, the  
 23 layer  $W_\beta$  can be replaced by a MLP.

24 **Robustness to non Gaussian data - to rev. #2 and #3** The crux of the method is to capture the covariance of the  
 25 data, which relates the different slopes of the models on incomplete data. This covariance will be relevant even on  
 26 non-Gaussian data, though we can so far only develop formal arguments under Gaussian assumptions.

27 **What is  $\nu$  in eq.(8) - to rev. #1** It is the smallest eigenvalue of the covariance matrix  $\Sigma$ . We apologize, the definition  
 28 was inadequately moved to the appendix; we will move it to the main text.

29 **What happens with few samples - to rev. #2** Reviewer #2 made a comment related to the number of samples. The  
 30 difficulty of a problem is linked to the ratio # dimensions/# samples, the higher the ratio, the more difficult the problem.  
 31 The highest ratio for which we presented experiments is for dimension 50 and 20000 samples. For such ratio and higher,  
 32 using a depth of 0 for our architecture is enough, more depth triggers overfitting. We showed theoretically in section  
 33 3.4 that the MLP can model a depth-0 Neumann network. Thus the MLP, whether in MCAR or selfmaking, is on par  
 34 or slightly better than the Neumann network for high ratios. As for MICE+LR, it has an advantage over MLP and  
 35 Neumann for high ratios in MCAR, but not in selfmasking, because Mice assumes the data to be M(C)AR.