NeuMiss networks: differentiable programming for supervised learning with missing values

Objectives

Supervised learning with missing values:

- both in the training set and at prediction time in the test set,
- under possibly non-ignorable missingness (Missing No At Random).

State-of-the-art and Challenges

- A rich literature on statistical inference but few works supervised learning with missing values.
- Most general methods for imputation are only valid if missingness is ignorable.
- Samples are represented by varying subsets of input variables \Rightarrow learn compensation mechanisms.
- The total number of possible missing data patterns is exponential in the dimension $(2^d) \Rightarrow$ keep the sample complexity polynomial.

Approach

- Derive the analytical expression of the optimal predictor under various missing data mechanisms.
- Propose a theoretically grounded neural network architecture (NeuMiss) designed to approximate thes optimal predictors.



Figure 1: Comparison of NeuMiss with wide and deep neural networks with ReLU activation, fed with the data imputed by 0 and the mask.

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

Notations and Assumptions

Rando	om variables	Ex. c
• $X \in \mathbb{R}^d$	^l : complete data	x = (1
• $\widetilde{X} \in \{\mathbb{R}$	$\mathbb{R} \cup \{\mathtt{NA}\}\}^d$: incomplete data	$\widetilde{x} = (1$
• $M \in \{0\}$	$0,1\}^d$: mask.	m = (0)
$\bullet obs(M)$: indices of the observed entries	$x_{obs(m)}$
Assum	ptions:	Optimal (E
Linear m	nodel: $Y = \beta_0^{\star} + \sum_{j=1}^d \beta_j^{\star} X_j + \epsilon$	$f^{\star} \in arg$
Gaussiar	n data: $X \sim \mathcal{N}\left(\mu, \Sigma\right)$	$f:(\mathbb{R}\cup\{$
Tł	ne optimal predictors u	nder variou
	mecha	nisms
Ignora	ble missing data mechanism	IS:
 Missin 	ng Completely At Random (MCA	AR): $P(M = i$
 Missin 	ng At Random (MAR): $P(M = $	m(X) = P(M
	M(C)AR Bay	ves predict
$f^{\star}(X_{obs}$	$(s, M) = \beta_0^{\star} + \langle \beta_{obs}^{\star}, X_{obs} \rangle + \langle \beta_{mi}^{\star} \rangle$	$\mu_{s}, \mu_{mis} + \Sigma_{mis}$
Non-ig	norable missing data mecha	ansim:
(Gaussian self-masking (MNAR) B
$f^{\star}(X$	$\begin{aligned} X_{obs}, M) &= \beta_0^{\star} + \langle \beta_{obs}^{\star}, X_{obs} \rangle + \langle \\ &\times (\tilde{\mu}_{mis} + D_{mis} \Sigma_{mis obs}^{-1} (\mu_{mis})) \end{aligned}$	$eta^{\star}_{mis}, (Id + D)$ $_{s} + \Sigma_{mis,obs} (\Sigma)$
Intuitie	on:	
	The optimal predictors	are linear
The slo ness of	pes of the obs. variables depend other correlated variables: $f^*($	I on M to con $X_{obs}, M) = eta_0$

References

[1] S van Buuren. Flexible Imputation of Missing Data. Boca Raton, FL: Chapman and Hall/CRC, 2018 [2] Marine Le Morvan et al. "Linear predictor on linearly-generated data with missing values: non consistency and solutions". In: vol. 108. AISTATS. 2020, pp. 3165–3174

Approximating the Bayes predictors

of realizations 1.1, 2.3, 3.1, 8, 5.271.1, NA, -3.1, 8, NA)0, 1, 0, 0, 1)= (1.1, 3.1, 8)

Bayes) predictor:

us missing data

$$m|X) = P(M = m)$$
$$I = m|X_{obs(m)})$$

tor

 $_{s,obs}(\Sigma_{obs})^{-1}(X_{obs}-\mu_{obs})\rangle$

Bayes predictor

 $\mathcal{D}_{mis} \Sigma_{mis|obs}^{-1}$ $\Sigma_{obs})^{-1} \left(X_{obs} - \mu_{obs} \right) \right) \rangle$

per pattern

npensate for the missing- $\beta_0(M) + \sum_{j \in \mathsf{obs}(\mathsf{M})} \beta_j(M) X_j$ **Main difficulty**: approx. of Σ_{obs}^{-1} , for any obs, i.e., any missing data pattern!



Figure 2:NeuMiss architecture for a depth of 4.

NeuMann iterations: approximate Σ_{obs}^{-1} by unrolling the order- ℓ truncation of a NeuMann series:

 $S_{obs(m)}^{(\ell)} = (Id - \Sigma_{obs(m)})S_{obs(m)}^{(\ell-1)} + Id.$

A new type of non-linearity: the multiplication entrywise by the mask.

Simulated data

- Gaussian covariates
- Response is a linear model
- 50% missing values



Figure 3:MCAR

- Robustness to the missing data mechanism.



Theoretically-grounded architecture: In M(C)AR and under a simplifying assumption in Gaussian selfmasking, NeuMiss with a depth $\ell + 1$ can exactly compute the order- ℓ approximation of the Bayes predictor.

Experimental results

Methods

- EM: Expectation Maximization.
- MICE [1] + LR: Conditional imputer followed by linear regression.
- MLP [2]: feedforward neural network

Figure 4: Gaussian self-masking

• Neumiss performances come close to the optimal performances.

• Suited for medium-sized datasets thanks to weights sharing across mdp.