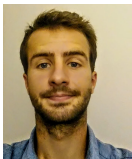


Linear predictor on linearly-generated data with missing values: non consistency and solutions

Marine Le Morvan

INRIA (Parietal), CNRS (IJCLab)



N. Prost



J. Josse

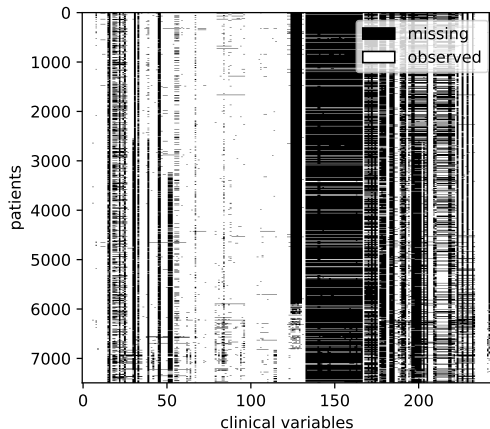


E. Scornet



G. Varoquaux

Missing values are ubiquitous in various fields



Traumabase clinical health records.

Most off-the-shelf supervised learning methods cannot be applied with missing values.

What to do:

- Complete-case analysis?
- Imputation prior to learning?
- Expectation Maximization?

We will study the case of **linear regression with missing values**, which has surprisingly received little attention up to now.

Content

- 1 Problem setting
- 2 The Bayes predictor
- 3 Linear approximation
- 4 Multilayer perceptron approximation
- 5 Empirical study

Outline

- 1 Problem setting
- 2 The Bayes predictor
- 3 Linear approximation
- 4 Multilayer perceptron approximation
- 5 Empirical study

Notation

- $\mathbf{x}_n \in \mathbb{R}^{n \times d}$: complete data (unavailable).
- $\mathbf{z}_n \in \{\mathbb{R} \times \text{na}\}^{n \times d}$: incomplete data (available).
- $\mathbf{m}_n \in \{0, 1\}^{n \times d}$: mask. 0s (1s) indicate the observed (missing) values.
- $\mathbf{y}_n \in \mathbb{R}^n$: the response vector.

$$\mathbf{z}_n = \begin{pmatrix} 9.1 & 8.5 \\ 2.1 & \text{na} \\ \text{na} & 9.6 \\ \text{na} & \text{na} \end{pmatrix}, \quad \mathbf{x}_n = \begin{pmatrix} 9.1 & 8.5 \\ 2.1 & 3.5 \\ 6.7 & 9.6 \\ 4.2 & 5.5 \end{pmatrix}, \quad \mathbf{m}_n = \begin{pmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{pmatrix}, \quad \mathbf{y}_n = \begin{pmatrix} 4.6 \\ 7.9 \\ 8.3 \\ 4.6 \end{pmatrix}$$

Each row of $\mathbf{x}_n, \mathbf{z}_n, \mathbf{m}_n, \mathbf{y}_n$ are realization of the generic random variable X, Z, M, Y .

The incomplete vector is related to X and M by:

$$Z = X \quad (1 - M) + \text{na} \quad M.$$

Problem setting

- **Working hypothesis:**

In this work, we assume that the response is linearly generated:

Assumption (Linear model)

$$Y = \beta_0 + X, \beta + \varepsilon, \quad X \in \mathbb{R}^d, \quad \varepsilon \sim N(0, \sigma^2).$$

- **Problem formulation:**

We wish to solve a least squares regression problem with missing values:

$$\min_{f: \{\mathbb{R} \times \text{na}\}^d \rightarrow \mathbb{R}} \mathbb{E} \left[(Y - f(Z))^2 \right],$$

Outline

- 1 Problem setting
- 2 The Bayes predictor**
- 3 Linear approximation
- 4 Multilayer perceptron approximation
- 5 Empirical study

Characterizing optimal regressors: the Bayes predictor

- A **Bayes predictor** f is the a minimizer of the loss (in our case least squares),

$$f^* = \underset{f: \{\mathbb{R} \times \mathbb{N}\}^d \rightarrow \mathbb{R}}{\operatorname{argmin}} \mathbb{E} \left[(Y - f(Z))^2 \right].$$

- For the least squares loss, we know it is the **conditional expectation of the response given the input**:
 - ✓ In the complete case: $f^* = \mathbb{E}[Y|X] = \beta_1 X + \beta_0$.
 - ✓ In the incomplete case: $f^* = \mathbb{E}[Y|Z] = \mathbb{E}[Y|M, X_{\text{obs}(M)}]$
- In the incomplete case, the Bayes predictor need not be linear.

Example

Let $Y = X_1 + X_2 + \varepsilon$, where $X_2 = \exp(X_1) + \varepsilon_1$. Now, assume that only X_1 is observed. Then the Bayes predictor is:

$$f(X_1) = X_1 + \exp(X_1).$$

The Bayes predictor for incomplete data

Assumption (Gaussian pattern mixture model)

$$X \mid (M = m) \sim \mathcal{N}(\mu^m, \Sigma^m).$$

Proposition (Expanded Bayes predictor)

Under our assumptions (linear model + Gaussian pattern mixture model), the Bayes predictor takes the form

$$f^*(Z) = W, \delta,$$

where the parameter $\delta \in \mathbb{R}^p$ is a function of β , $(\mu^m)_m$, $(\Sigma^m)_m$, and the random variable $W \in \mathbb{R}^p$ is the concatenation of $j = 1, \dots, 2^d$ blocks, each one being

$$(\mathbb{1}_{M=m_j}, X_{\text{obs}(m_j)} \mathbb{1}_{M=m_j}).$$

where W is an expansion of Z .

The Bayes predictor for incomplete data

Assumption (Gaussian pattern mixture model)

$$X / (M = m) \quad \mathcal{N}(\mu^m, \Sigma^m).$$

Proposition (Expanded Bayes predictor)

Under our assumptions (linear model + Gaussian pattern mixture model), the Bayes predictor takes the form

$$f^*(Z) = W, \delta ,$$

where (ex. $d=2$)

$$W = \left(\begin{array}{ccc|cc|cc|c} 1 & x_{1,1} & x_{1,2} & 0 & 0 & 0 & 0 & 0 \\ 1 & x_{2,1} & x_{2,2} & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 1 & x_{3,1} & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & x_{4,1} & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & 1 & x_{5,2} & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & x_{6,2} & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{array} \right)$$

Outline of the proof

Under the linear assumption we have:

$$\begin{aligned}f^*(Z) &= E[Y/Z] \\&= E[\beta_0 + \beta^T X \mid Z] \\&= E[\beta_0 + \beta^T X \mid M, X_{obs(M)}] \\&= \beta_0 + \beta_{obs(M)}^T X_{obs(M)} + \beta_{mis(M)}^T E[X_{mis(M)} \mid M, X_{obs(M)}]\end{aligned}$$

Moreover under the Gaussian per pattern assumption,

$$E[X_{mis(M)} \mid M, X_{obs(M)}] = \theta + \Gamma X_{obs(M)}$$

where θ and Γ depend on μ^M and Σ^M .

Thus,

$$f^*(Z) = \beta_0 + \beta_{mis(M)}^T \theta + (\beta_{obs(M)} + \Gamma)^T X_{obs(M)}$$

i.e., the Bayes predictor is **linear per pattern**.

The expanded linear model

$f(Z) = W, \delta$ where (example $d = 2$):

$$W = \left(\begin{array}{ccc|cc|cc|c} 1 & x_{1,1} & x_{1,2} & 0 & 0 & 0 & 0 & 0 \\ 1 & x_{2,1} & x_{2,2} & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 1 & x_{3,1} & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & x_{4,1} & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & 1 & x_{5,2} & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & x_{6,2} & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{array} \right)$$

Problem: the dimension of W is

$$p = \sum_{k=0}^d \binom{d}{k} \times (k+1) = 2^{d-1} \times (d+2).$$

Outline

- 1 Problem setting
- 2 The Bayes predictor
- 3 Linear approximation**
- 4 Multilayer perceptron approximation
- 5 Empirical study

The linear approximation model

The Bayes predictor can be expressed as a polynome of X and M , which can be truncated to a first order approximation.

Definition (Linear approximation)

We define the linear approximation of f^* as

$$f_{\text{approx}}^*(Z) = \beta_{0,0}^* + \sum_{j=1}^d \beta_{j,0}^* M_j + \sum_{j=1}^d \beta_j^* X_j (1 - M_j).$$

Estimation of the linear approximation model

- f_{approx}^* can be estimated by fitting a linear model on X imputed by 0 concatenated with the mask.
- This is equivalent to jointly fitting a linear model on X and optimizing an imputation constant for each variable.

$$\text{Given } \begin{pmatrix} X_1 & X_2 \\ 1.1 & 3.2 \\ \text{NA} & 0.1 \\ 4.6 & \text{NA} \\ 4.0 & 0.9 \\ \text{NA} & 2.2 \end{pmatrix}, \quad \begin{pmatrix} X_1 & X_2 \\ 1.1 & 3.2 \\ C_1 & 0.1 \\ 4.6 & C_2 \\ 4.0 & 0.9 \\ C_1 & 2.2 \end{pmatrix} \quad \begin{pmatrix} X_1 & M_1 & X_2 & M_2 \\ 1.1 & 0 & 3.2 & 0 \\ 0 & 1 & 0.1 & 0 \\ 4.6 & 0 & 0 & 1 \\ 4.0 & 0 & 0.9 & 0 \\ 0 & 1 & 2.2 & 0 \end{pmatrix}.$$

Indeed,

$$\beta_j \{X_j(1 - M_j) + c_j M_j\} = \beta_j X_j(1 - M_j) + \{\beta_j c_j\} M_j.$$

Finite sample bounds for linear predictors

The Bayes predictor and its linear approximation offer different bias-variance tradeoffs.

Assumption

- $Y = f_{\text{Bayes}}(Z) + \text{noise}(Z)$ where $\text{noise}(Z)$ is a centred noise conditional on Z and such that there exists $\sigma^2 > 0$ satisfying $\mathbb{V}[Y|Z] = \sigma^2$ almost surely,
- $f_{\text{Bayes}} \leq L$,
- $\text{Supp}(X) \subseteq [-1, 1]^d$.

This assumption is required for the next two results.

Finite sample bounds for linear predictors

Under these assumptions:

Theorem

- The risk of the OLS estimate clipped at L for the **expanded model** satisfies

$$\frac{2^d c_1}{n+1} R(T_L f_{\hat{\beta}_{\text{expanded}}}) - \sigma^2 \leq c \max\{\sigma^2, L^2\} \frac{2^{d-1}(d+2)(1+\log n)}{n}$$

- The risk of the OLS estimate clipped at L for the **linear approximation model** satisfies

$$R(T_L f_{\hat{\beta}_{\text{approx}}}) - \sigma^2 \leq c \max\{\sigma^2, L^2\} \frac{2d(1+\log n)}{n} + 64(d+1)^2 L^2$$

It follows that the risk of the expanded model is lower than that of the linear approximation model if:

$$n \geq \frac{2^d}{d}$$

Outline

- 1 Problem setting
- 2 The Bayes predictor
- 3 Linear approximation
- 4 Multilayer perceptron approximation**
- 5 Empirical study

Why a Multilayer perceptron?

A Multilayer Perceptron with:

- Rectified Linear Units activation functions for hidden units ($ReLU(x) = \max(0, x)$),
- Identity activation for the output unit,

produces a prediction function that is **piecewise affine**.

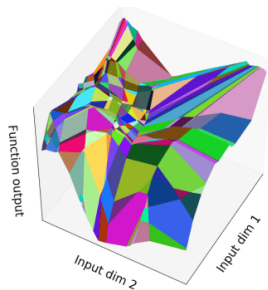
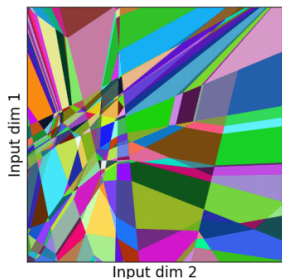


Figure from Hanin et al. 2019

Bayes consistency of the MLP

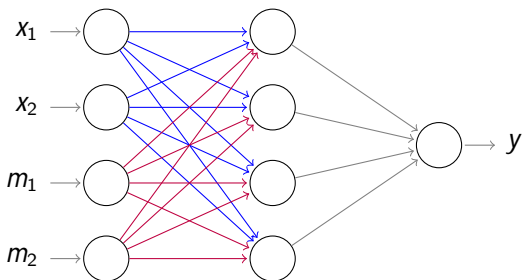
Theorem (MLP)

Assume that the Bayes predictor takes the form described earlier (expanded Bayes Predictor). A MLP:

- *with one hidden layer containing 2^d hidden units*
 - *ReLU activation functions*
 - *which is fed with the concatenated vector (X, M) where X is imputed by zero*
- is Bayes consistent.*

Proof: We show that there exists a configuration of the parameters of the MLP so that the resulting predictor is the Bayes predictor.

Proof 1/3 - Learned imputations



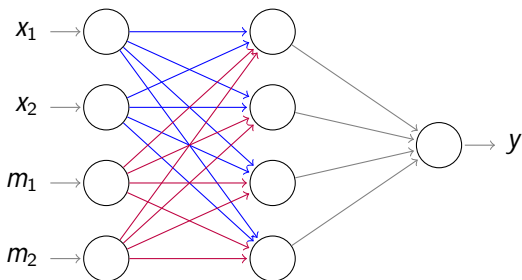
Parameters hidden layer:

$$W^{(1)} = [W^{(X)}, W^{(M)}] \quad \mathbb{R}^{4 \times 4}$$
$$b^{(1)} \quad \mathbb{R}^4$$

Parameters output layer:

$$W^{(2)} \quad \mathbb{R}^4$$
$$b^{(2)} \quad \mathbb{R}$$

Proof 1/3 - Learned imputations



Parameters hidden layer:

$$W^{(1)} = [W^{(X)}, W^{(M)}] \quad \mathbb{R}^{4 \times 4}$$
$$b^{(1)} \quad \mathbb{R}^4$$

Parameters output layer:

$$W^{(2)} \quad \mathbb{R}^4$$
$$b^{(2)} \quad \mathbb{R}$$

The activation of hidden unit k for input (x, m) is:

$$\begin{aligned} a_k &= W_{k,\cdot}^{(X)} x + W_{k,\cdot}^{(M)} m + b_k^{(1)} \\ &= W_{k,\cdot}^{(X)} x + W_{k,\cdot}^{(X)} G_{k,\cdot} m + b_k^{(1)} \\ &= W_{k,obs(m)}^{(X)} x_{obs(m)} + W_{k,mis(m)}^{(X)} G_{k,mis(m)} + b_k^{(1)} \end{aligned}$$

where G (reparametrization of $W^{(M)}$) can be seen as **learned imputations**.

Proof 2/3 - one-to-one mapping mdp/hidden unit

The proof shows that the parameters of the MLP can be chosen so that:

- 1 all points with missing data pattern m_k exclusively activate hidden unit k , and hidden unit k is exclusively activated by points with missing data pattern m_k .

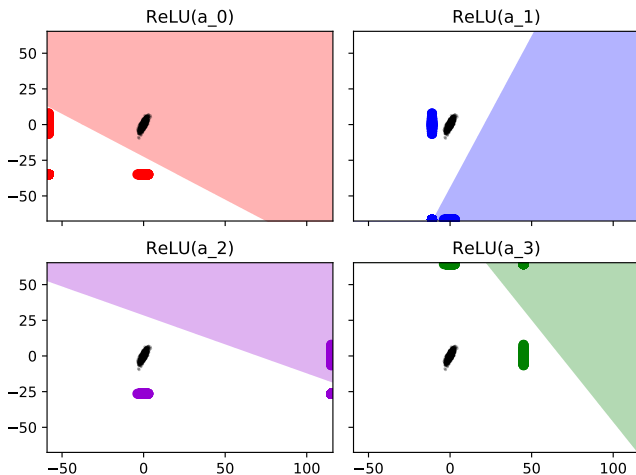
$$\begin{aligned}y(x, m_k) &= \sum_{h=1}^{2^d} W_h^{(2)} \text{ReLU}(a_h) + b^{(2)} \\&= \sum_{h=1}^{2^d} W_h^{(2)} \text{ReLU}(W_{h, \text{obs}(m_k)}^{(X)} x_{\text{obs}(m_k)} + W_{h, \text{mis}(m_k)}^{(X)} G_{h, \text{mis}(m_k)} + b_h^{(1)}) + b^{(2)} \\&= W_k^{(2)} \left(W_{k, \text{obs}(m_k)}^{(X)} x_{\text{obs}(m_k)} + W_{k, \text{mis}(m_k)}^{(X)} G_{k, \text{mis}(m_k)} + b_k^{(1)} \right) + b^{(2)}\end{aligned}$$

i.e, the MLP produces a predictor $y(x, m_k)$ that is linear per pattern.

- 2 The slopes and biases of $y(x, m_k)$ equal those of the Bayes predictor.

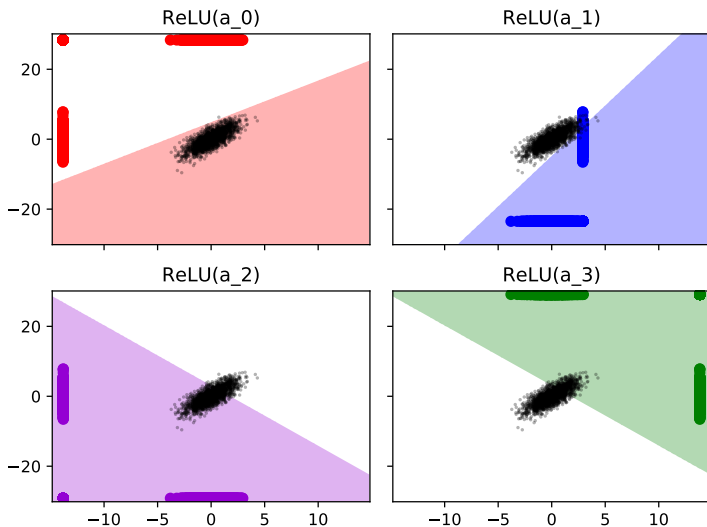
Proof 3/3 - visualisation of a bayes consistent MLP

We simulated data (X, M) in 2 dimensions, and based on our proof, built a MLP (with 4 hidden units) that is Bayes consistent.



$$y(x, m) = W_{1,\cdot}^{(2)} \text{ReLU}(a_0) + W_{1,\cdot}^{(2)} \text{ReLU}(a_1) + W_{2,\cdot}^{(2)} \text{ReLU}(a_2) + W_{3,\cdot}^{(2)} \text{ReLU}(a_3) + b^{(2)}$$

Example of an optimized MLP in two dimensions.



$$y(x, m) = W_{1,\cdot}^{(2)} \text{ReLU}(a_0) + W_{1,\cdot}^{(2)} \text{ReLU}(a_1) + W_{2,\cdot}^{(2)} \text{ReLU}(a_2) + W_{3,\cdot}^{(2)} \text{ReLU}(a_3) + b^{(2)}$$

Trading off estimation and approximation error

Number of parameters of:

- a MLP with one hidden layer and 2^d units:

$$(d + 1)2^{d+1} + 1$$

- the expanded linear model:

$$(d + 1)2^{d-1}$$

The MLP is slightly overparametrized, and the number of parameters is exponential in d .

However, contrary to the expanded linear model, **the MLP provides a natural way to reduce the model capacity** by reducing the number of hidden units.

Outline

- 1 Problem setting
- 2 The Bayes predictor
- 3 Linear approximation
- 4 Multilayer perceptron approximation
- 5 Empirical study**

Simulation models

The data (X, M) is generated according to 3 simulation models:

- **mixture 1:**

- ▶ $P(X) = \mathcal{N}(\mu, \Sigma)$
- ▶ $P(M) = \frac{1}{2^d}$
- ▶ Gaussian pattern mixture model with 1 component
- ▶ Corresponds to a Missing Completely At Random (MCAR) problem

- **mixture 3:**

- ▶ $P(X/M = m) = \mathcal{N}(\mu_m, \Sigma_m)$, with 3 distinct Gaussian components.
- ▶ $P(M) = \frac{1}{2^d}$
- ▶ Gaussian pattern mixture model (with 3 components)

- **selfmasking:**

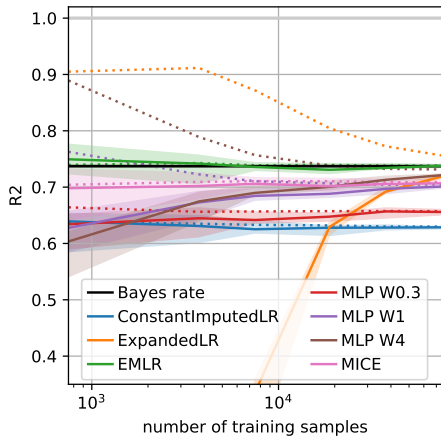
- ▶ $P(X) = \mathcal{N}(\mu, \Sigma)$
- ▶ $P(M = 1/X_j) = \text{Probit}(\lambda_j (X_j - \mu_0))$
- ▶ Not an instance of pattern mixture model! (Theory does not hold)
- ▶ Corresponds to a typical Missing Non At Random (MNAR) problem

Estimation Approaches

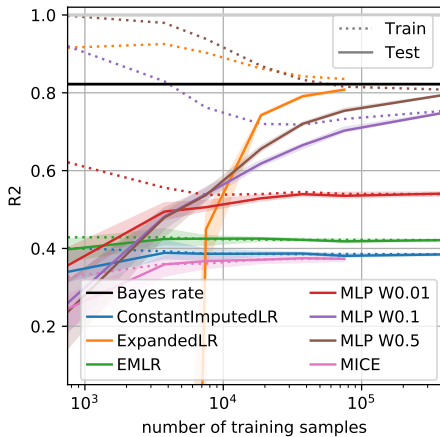
- **EMLR**: EM is used to fit a multivariate normal distribution for the $(p + 1)$ -dimensional random variable (X_1, \dots, X_p, Y) .
- **ConstantImputedLR**: Optimal imputation method.
- **MICE**: Conditional imputation with an iterative imputer (similar to the well known MICE) followed by linear regression.
- **ExpandedLR**: Expanded linear model.
- **MLP**: Multilayer perceptron with one hidden layer whose size is varied between and 1 and 2^d hidden units.

Learning curves: Gaussian mixtures

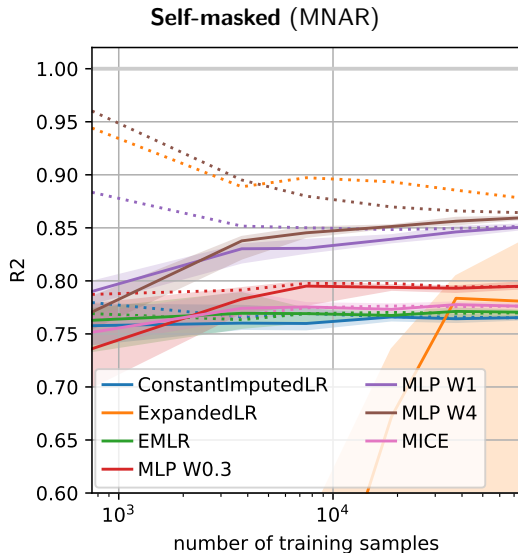
Mixture 1 (MCAR)



Mixture 3



Learning curves: self-masking



Conclusion

Conclusion:

- The Bayes-optimal predictor is no longer a linear function of the data.
- It is explicit under Gaussian assumptions, but high-dimensional.
- Possible approximations include constant imputation and MLP, which can be consistent.
- The MLP adapts naturally to the complexity of the data.
- Our risk-minimisation strategy is robust to the missing-value mechanism.