# Near-optimal rate of consistency for linear prediction with missing values

Alexis Ayme

Claire Boyer   Aymeric Dieuleveut  and  Erwan Scornet

# Background

o **Growing mass of data => NA (not attribut)/missing values**

o **Different sources:**

1. Bugs
2. Cost

| $1 | $10 | $100 | $0 |
|----|-----|------|----|

| | | | |
|--|--|--|--|
| | | NA | |
| | | | |
| | NA | NA | |

3. Multiplication of sources (i.e. merge)



4. Sensitive data

| Age | Job | Incomes |
|-----|-----|---------|

| | | |
|--|--|--|
| | | NA |
| | | |
| | | NA |
| NA | | NA |
| | | NA |

# Background

o **Growing mass of data => NA (not attribut)/missing values**

o **Different sources:**

1. Bugs
2. Cost

| $1 | $10 | $100 | $0 |
|----|-----|------|----|
|    |     | NA   |    |
|    |     |      |    |
|    | NA  | NA   |    |

3. Multiplication of sources (i.e. merge)



4. Sensitive data

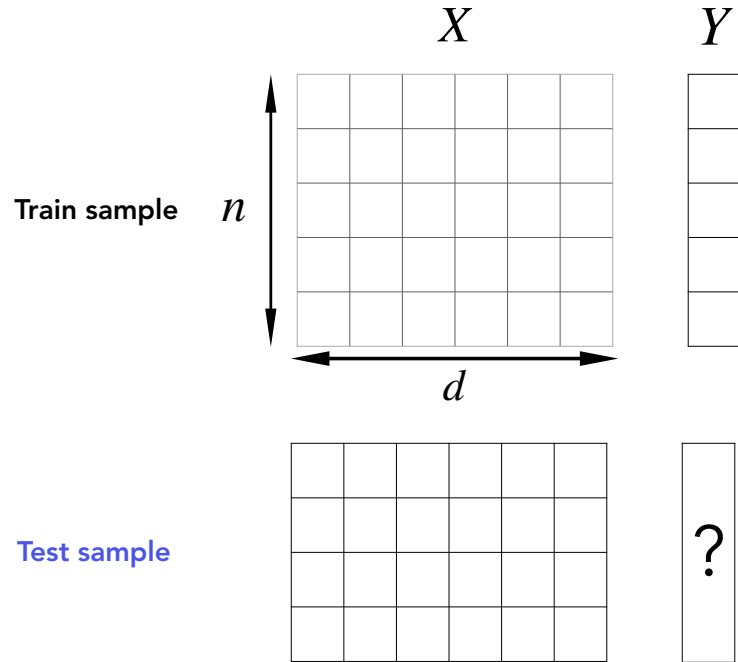| Age | Job | Income |
|-----|-----|--------|
|     |     | NA     |
|     |     |        |
|     |     | NA     |
| NA  |     | NA     |
|     |     | NA     |

o Any statistical analyses require **complete** data

o Strategy 1: **complete** the dataset before the ML process (e.g. by collaborative filtering)

o Strategy 2: **adapt** statistical analysis to handle missing values (e.g. EM algorithm to perform regression with NA)

# Background

○ **Growing mass of data => NA (not attribut)/missing values**

○ **Different sources:**

1. Bugs
2. Cost

| $1 | $10 | $100 | $0 |
|----|-----|------|----|
|    |     | NA   |    |
|    |     |      |    |
|    | NA  | NA   |    |

3. Multiplication of sources (i.e. merge)



4. Sensitive data

| Age | Job | **Income** |
|-----|-----|------------|
|     |     | NA         |
|     |     |            |
|     |     | NA         |
| NA  |     | NA         |
|     |     | NA         |

○ Any statistical analyses require **complete** data

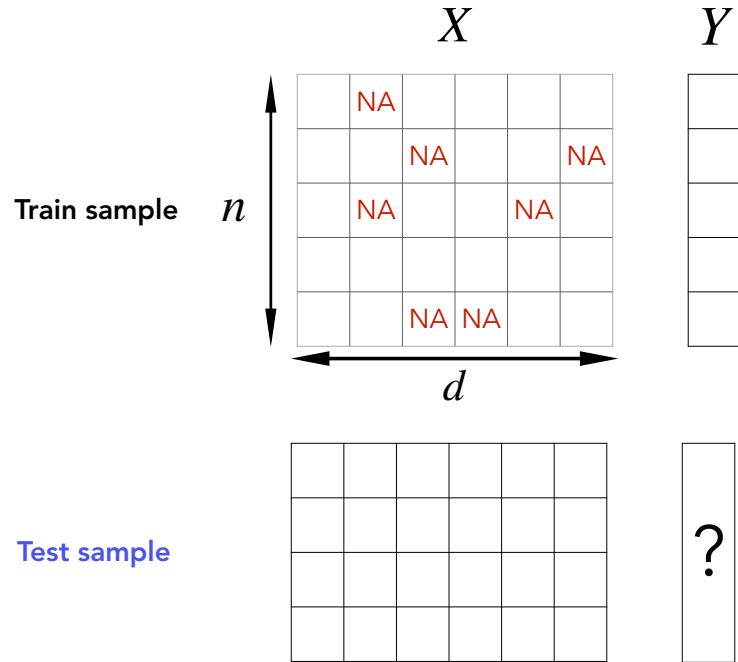  ○ Strategy 1: **complete** the dataset before the ML process (e.g. by collaborative filtering)

  ○ Strategy 2: **adapt** statistical analysis to handle missing values (e.g. EM algorithm to perform regression with NA)

> What about **supervised learning**?
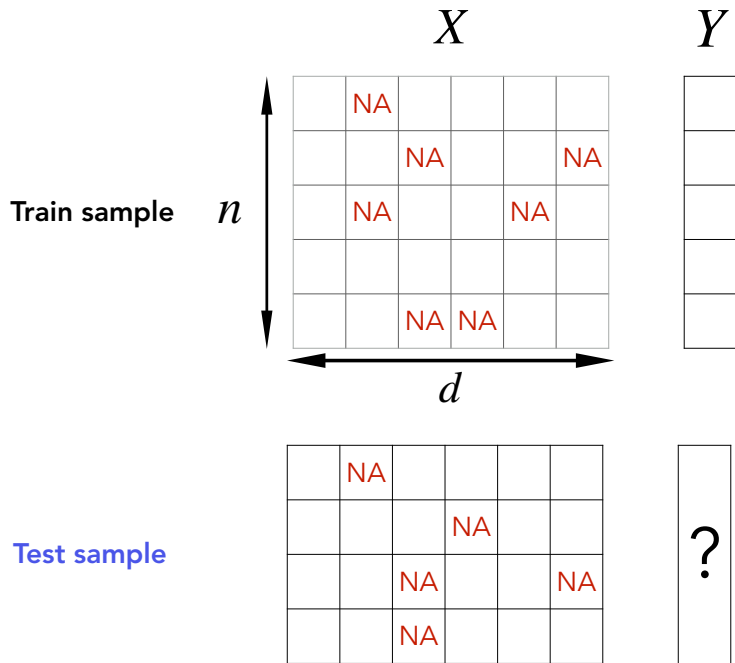> i.e. prediction with NAs

# Supervised learning with missing values (NA)

$X$  $Y$

**Train sample**  $n$

$d$

**Test sample**  ?

# Supervised learning with missing values (NA)

# Supervised learning with missing values (NA)
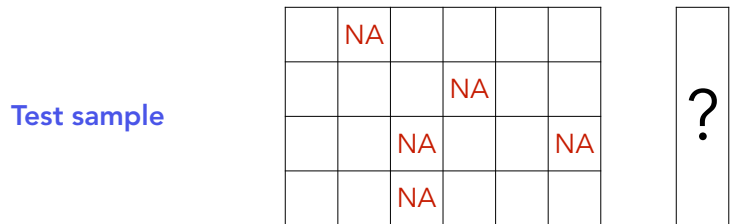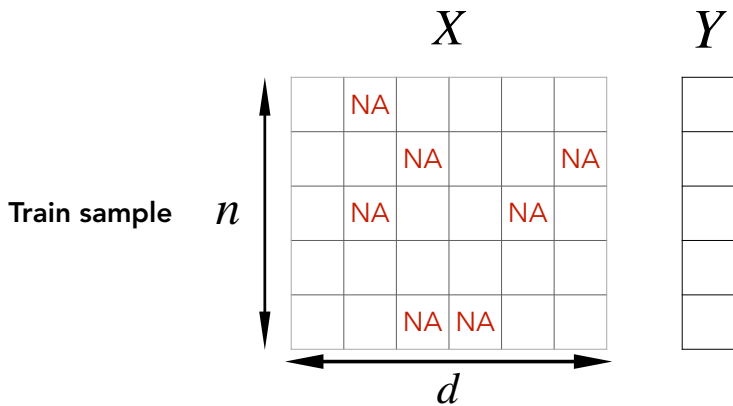


$X$

$Y$

Train sample  $n$

$d$

Test sample

?

$$Y = \beta^T X + \varepsilon$$

$$\|\hat{\beta} - \beta\|_2^2$$

$$\left(\hat{\beta}^T x - \beta^* x\right)^2$$

# Supervised learning with missing values (NA)



$X$      $Y$

Train sample   $n$

$d$

Test sample

○ Missing pattern: $M_i \in \{0,1\}^d$

$$X_i = \boxed{\text{NA}\ |\ 1\ |\ -5\ |\ \text{NA}\ |\ 0\ |\ 2}$$

$$M_i = (1,\ 0,\ 0,\ 1,\ 0,\ 0)$$

○ Input:     $Z = (X_{\text{obs}}, M)$

○ Output:   $Y \in \mathbb{R}$

?

# Supervised learning with missing values (NA)



$X$       $Y$

Train sample   $n$

$d$

Test sample

?

○ Missing pattern: $M_i \in \{0,1\}^d$

$$X_i = \boxed{\text{NA}}\,\boxed{1}\,\boxed{-5}\,\boxed{\text{NA}}\,\boxed{0}\,\boxed{2}$$

$$M_i = (1,\ 0,\ 0,\ 1,\ 0,\ 0)$$

○ Input:   $Z = (X_{\text{obs}}, M)$

○ Output:   $Y \in \mathbb{R}$

**Goal:** Predict on **test sample** minimizing

$$R(f) = \mathbb{E}_{Z,Y}\left[\left(Y - f(Z)\right)^2\right]$$

# Zoo of assumptions on NA

○ **First point of view :** $P(X, M) = P(M|X)P(X)$

1) Assumption on $P(X)$:
   Example: $X$ Gaussian Vector
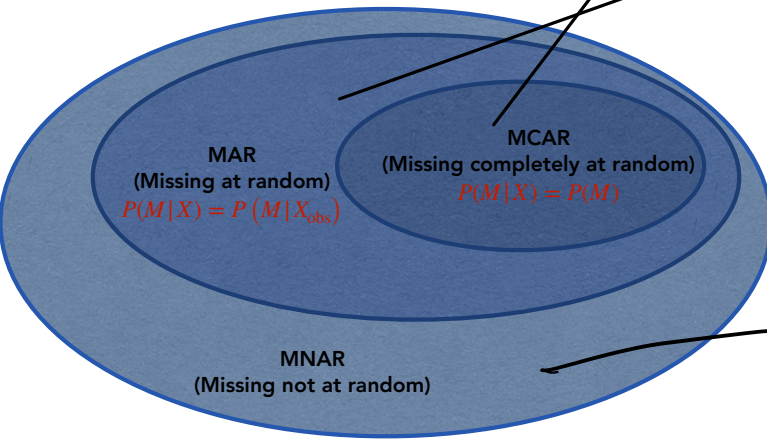
2) Assumption on $P(M|X)$:
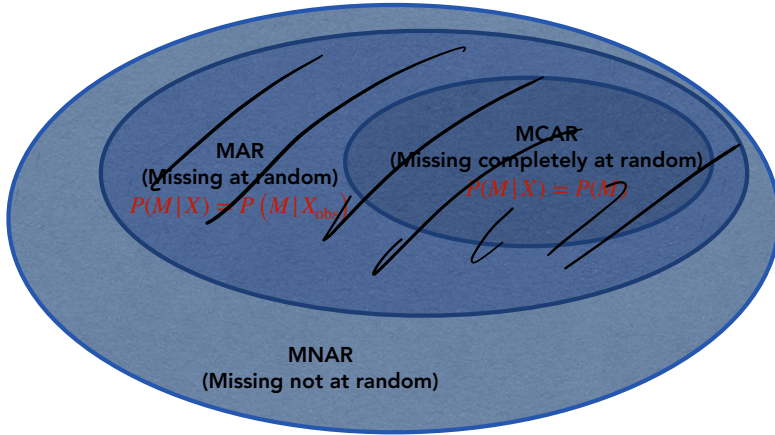
Bugs

Medical protocol

check up

Good        Bad

NA

High value censorship

**MAR**
**(Missing at random)**
$P(M|X) = P(M|X_{obs})$

**MCAR**
**(Missing completely at random)**
$P(M|X) = P(M)$

**MNAR**
**(Missing not at random)**

# Zoo of assumptions on NA

○ **First point of view :** $P(X, M) = P(M|X)P(X)$

1) Assumption on $P(X)$:
   Example: $X$ Gaussian Vector

2) Assumption on $P(M|X)$:



MAR
(Missing at random)
$P(M|X) = P(M|X_{obs})$

MCAR
(Missing completely at random)
$P(M|X) = P(M)$

MNAR
(Missing not at random)

○ **Second point of view :** $P(X, M) = P(M)P(X|M)$

1) Assumption on $P(M)$:
   Example: $P(M = m) = p_m$

2) Assumption on $P(X|M)$:

GPMM (Gaussian pattern mixture model):
$X|(M = m)$ Gaussian Vector

UJA          Frame

# Pattern-by-Pattern regression

o **Assumption**: linear model for complete inputs

$$y_i = \beta^\top X_i + \epsilon_i$$

/!\ With NA, the Bayes predictor does *not* necessarily remain linear

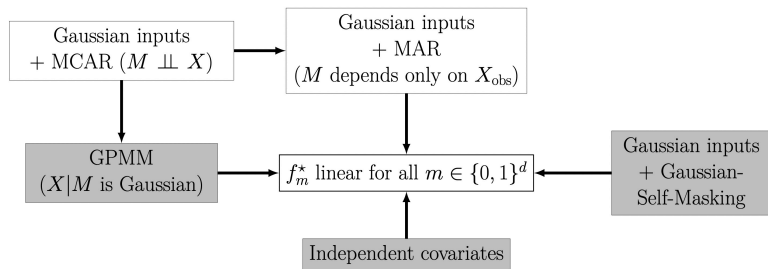o **Bayes predictor** (better prediction) decomposition

$$f^\star(Z) = \sum_{m \in \{0,1\}^d} f_m^\star(X_{obs(m)}) \mathbf{1}_{M=m}$$

Local **Bayes prediction** for the missing pattern ($M = m$)



Assumption to obtain linearity

**Proposition: (Le Morvan et al. 2020)**
Under **linear model** and several **missing data scenarios**
(including MNAR), $f_m^\star$ are **linear**

# Pattern-by-Pattern regression

o **Assumption**: linear model for complete inputs

$$y_i = \beta^\top X_i + \epsilon_i$$

/!\ With NA, the Bayes predictor does *not* necessarily remain linear

o **Bayes predictor** decomposition

$$f^\star(Z) = \sum_{m \in \{0,1\}^d} f_m^\star(X_{obs(m)}) \mathbf{1}_{M=m}$$

Local **Bayes prediction** for the missing pattern ($M = m$)

o **Pattern-by-pattern** predictor

$$\hat{f}(Z) = \sum_{m \in \{0,1\}^d} \hat{f}_m(X_{obs(m)}) \mathbf{1}_{M=m}$$

Local **Least-Square** regression on
$\{(X_{i,obs}, Y_i), M_i = m\}$

---

**Proposition: (Le Morvan et al. 2020)**
Under **linear model** and several **missing data scenarios**
(including MNAR), $f_m^\star$ are **linear**

# Pattern-by-Pattern regression

○ **Assumption**: linear model for complete inputs

$$y_i = \beta^\top X_i + \epsilon_i$$

/!\ With NA, the Bayes predictor does *not* necessarily remain linear

○ **Bayes predictor** decomposition

$$f^\star(Z) = \sum_{m \in \{0,1\}^d} f_m^\star(X_{obs(m)}) \mathbf{1}_{M=m}$$

Local **Bayes prediction** for the missing pattern ($M = m$)

---

**Proposition: (Le Morvan et al. 2020)**
Under **linear model** and several **missing data scenarios** (including MNAR), $f_m^\star$ are **linear**

---

○ **Pattern-by-pattern** predictor

$$\hat{f}(Z) = \sum_{m \in \{0,1\}^d} \hat{f}_m(X_{obs(m)}) \mathbf{1}_{M=m}$$

Local **Least-Square** regression on
$$\{(X_{i,obs}, Y_i), M_i = m\}$$

---

**Theorem 1:**
Under Lipschitz and sub-Gaussian assumptions,

$$\mathcal{E}(\hat{f}) := \mathbb{E}\left[\left(f^\star(Z) - \hat{f}(Z)\right)^2\right] \leq A \log(n) 2^d \frac{d}{n} + \text{Approx}$$

---

○ Optimal for equiprobable missing patterns $\left(p_m = \dfrac{1}{2^d}\right)$

○ Tight for the worst case of pattern-by-pattern predictors

○ Sub-optimal for other distributions?

# Thresholded Pattern-by-Pattern regression

○ Adaptivity to the missing pattern distribution to overcome the curse of dimensionality

○ Overfitting reduction

via **Thresholded** **P-by-P** predictor:

$$\hat{p}_m = \textbf{frequency of pattern } m$$

$$\hat{f}(Z) = \sum_{m \in \{0,1\}^d} \hat{f}_m(X_{obs(m)}) \mathbf{1}_{M=m} \mathbf{1}_{\hat{p}_m > \frac{d}{n}}$$

Local **Least-Square** regression on $\left\{ (X_{i,obs}, Y_i), M_i = m \right\}$

# Thresholded Pattern-by-Pattern regression

o Adaptivity to the missing pattern distribution to overcome the curse of dimensionality

o Overfitting reduction

via **Thresholded** **P-by-P** predictor:

$\hat{p}_m = $ **frequency** of pattern $m$

$$\hat{f}(Z) = \sum_{m \in \{0,1\}^d} \hat{f}_m(X_{obs(m)})\mathbf{1}_{M=m}\mathbf{1}_{\hat{p}_m > \frac{d}{n}}$$

Local **Least-Square** regression on $\left\{(X_{i,obs}, Y_i), M_i = m\right\}$

o Definition: missing pattern **complexity**

$$\mathfrak{C}_p\left(\frac{d}{n}\right) = \sum_{m \in \{0,1\}^d} p_m \wedge \frac{d}{n}$$

**Theorem 2: (Main result)**
Under Lipschitz and Sub-Gaussian assumptions

$$\mathscr{E}(\hat{f}) \leq A \log(n)\mathfrak{C}_p\left(\frac{d}{n}\right) + \text{Approx}$$

# Thresholded Pattern-by-Pattern regression

o Adaptivity to the missing pattern distribution to overcome the curse of dimensionality

o Overfitting reduction

via **Thresholded** **P-by-P** predictor:

$\hat{p}_m = $ **frequency** of pattern $m$

$$\hat{f}(Z) = \sum_{m \in \{0,1\}^d} \hat{f}_m(X_{obs(m)}) \mathbf{1}_{M=m} \mathbf{1}_{\hat{p}_m > \frac{d}{n}}$$

Local **Least-Square** regression on $\{(X_{i,obs}, Y_i), M_i = m\}$

o Definition:  missing pattern **complexity**

$$\mathfrak{C}_p\left(\frac{d}{n}\right) = \sum_{m \in \{0,1\}^d} p_m \wedge \frac{d}{n}$$

**Theorem 2: (Main result)**
Under Lipschitz and Sub-Gaussian assumptions

$$\mathscr{E}(\hat{f}) \le A \log(n) \mathfrak{C}_p\left(\frac{d}{n}\right) + \text{Approx}$$

**Examples:**

1. **Uniform** distribution:  $\mathfrak{C}_p\left(\frac{d}{n}\right) = 2^d \frac{d}{n}$

2. **Bernoulli** distribution: $M_j \sim \mathscr{B}(\epsilon)$ and $\epsilon \le \frac{d}{n}$

$$\mathfrak{C}_p\left(\frac{d}{n}\right) \le \frac{d^2}{n}$$

# The thresholded P-by-P predictor is near-optimal

o **Minimax risk**

**Worst case** on a class of problem $\mathcal{P}_p$

$$\mathcal{E}_{\text{mini}}\left(p\right) = \inf_{\tilde{f}} \sup_{\mathbb{P}\in\mathcal{P}_p} \mathbb{E}_{\mathbb{P}} \left[ \left( \tilde{f}(Z) - f^{\star}(Z) \right)^2 \right]$$

**Best algorithm**

where $\mathcal{P}_p$ represents a class of data distributions
o for which the missing pattern distribution is $p$
o under Lipschitz and Sub-Gaussian assumptions

# The thresholded P-by-P predictor is near-optimal

o **Minimax risk**

**Worst case** on a class of problem $\mathscr{P}_p$

$$\mathscr{E}_{\text{mini}}(p) = \inf_{\tilde{f}} \sup_{\mathbb{P} \in \mathscr{P}_p} \mathbb{E}_{\mathbb{P}}\left[\left(\tilde{f}(Z) - f^\star(Z)\right)^2\right]$$

**Best algorithm**

where $\mathscr{P}_p$ represents a class of data distributions
o for which the missing pattern distribution is $p$
o under Lipschitz and Sub-Gaussian assumptions

**Theorem 3:**

$$\sigma^2 \mathfrak{C}_p\left(\frac{1}{n}\right) \lesssim \mathscr{E}_{\text{mini}}(p) \leq \mathscr{E}(\hat{f}) \leq A \log(n) \mathfrak{C}_p\left(\frac{d}{n}\right)$$
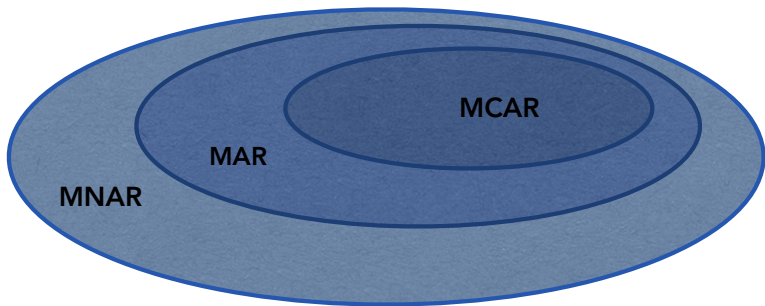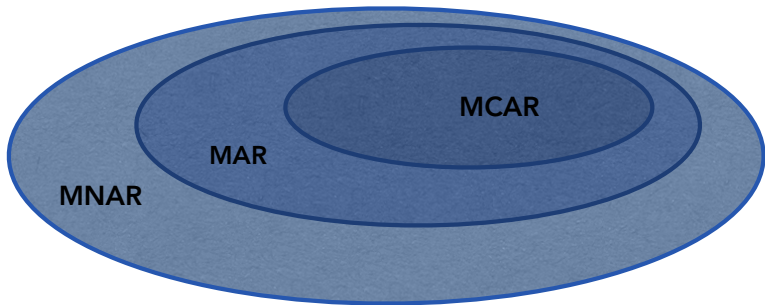
Theorem 2

o Lower bound still holds when $\mathscr{P}_p$ includes **MAR** missing values

**Examples**
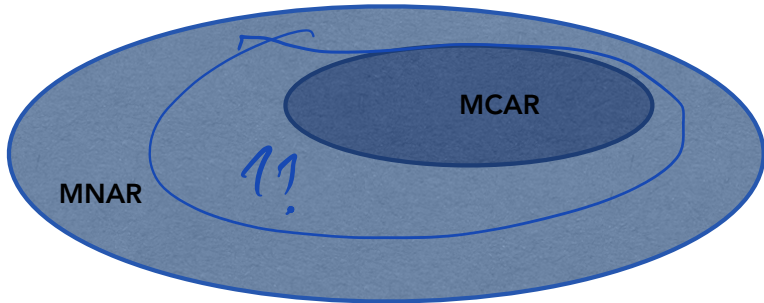
1. **Uniform** distribution: $\mathfrak{C}_p\left(\frac{1}{n}\right) = \frac{2^d}{n}$, $\mathfrak{C}_p\left(\frac{d}{n}\right) = 2^d \frac{d}{n}$

2. **Bernoulli** distribution: $\mathfrak{C}_p\left(\frac{1}{n}\right) = \frac{d}{n}$, $\mathfrak{C}_p\left(\frac{d}{n}\right) = \frac{d^2}{n}$

# The thresholded P-by-P predictor is near-optimal

o Inference POV



Theorem 3:
$$\sigma^2 \mathfrak{C}_p\left(\frac{1}{n}\right) \lesssim \underbrace{\mathscr{E}_{\mathrm{mini}}\left(p\right) \leq \mathscr{E}(\hat{f}) \leq A\log(n)\mathfrak{C}_p\left(\frac{d}{n}\right)}_{\text{Theorem 2}}$$

o Lower bound still holds when $\mathscr{P}_p$ includes **MAR** missing values

### Examples

1. **Uniform** distribution: $\mathfrak{C}_p\left(\frac{1}{n}\right) = \frac{2^d}{n}$, $\mathfrak{C}_p\left(\frac{d}{n}\right) = 2^d\frac{d}{n}$

2. **Bernoulli** distribution: $\mathfrak{C}_p\left(\frac{1}{n}\right) = \frac{d}{n}$, $\mathfrak{C}_p\left(\frac{d}{n}\right) = \frac{d^2}{n}$

# The thresholded P-by-P predictor is near-optimal

o Inference POV



Theorem 3:
$$\sigma^2 \mathfrak{C}_p\left(\frac{1}{n}\right) \lesssim \mathcal{E}_{\text{mini}}(p) \leq \mathcal{E}(\hat{f}) \leq A\log(n)\mathfrak{C}_p\left(\frac{d}{n}\right)$$
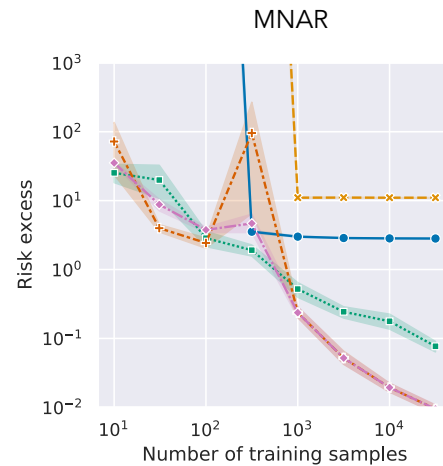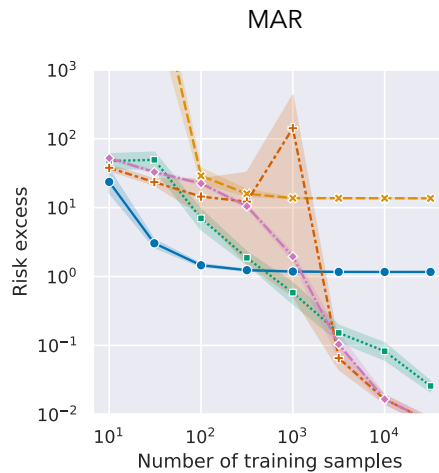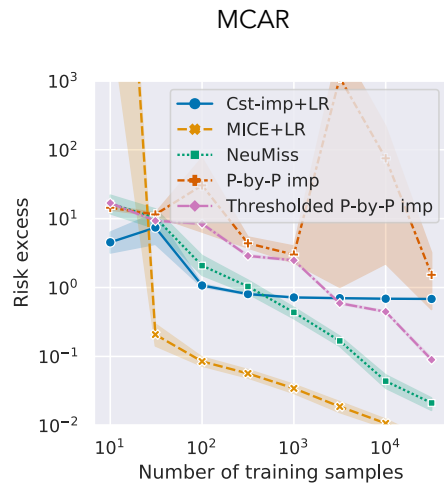$$\underbrace{\qquad\qquad\qquad\qquad}_{\text{Theorem 2}}$$

o Lower bound still holds when $\mathscr{P}_p$ includes **MAR** missing values

o Supervised learning POV



**Examples**

1. **Uniform** distribution: $\mathfrak{C}_p\left(\frac{1}{n}\right) = \frac{2^d}{n}$, $\mathfrak{C}_p\left(\frac{d}{n}\right) = 2^d\frac{d}{n}$

2. **Bernoulli** distribution: $\mathfrak{C}_p\left(\frac{1}{n}\right) = \frac{d}{n}$, $\mathfrak{C}_p\left(\frac{d}{n}\right) = \frac{d^2}{n}$
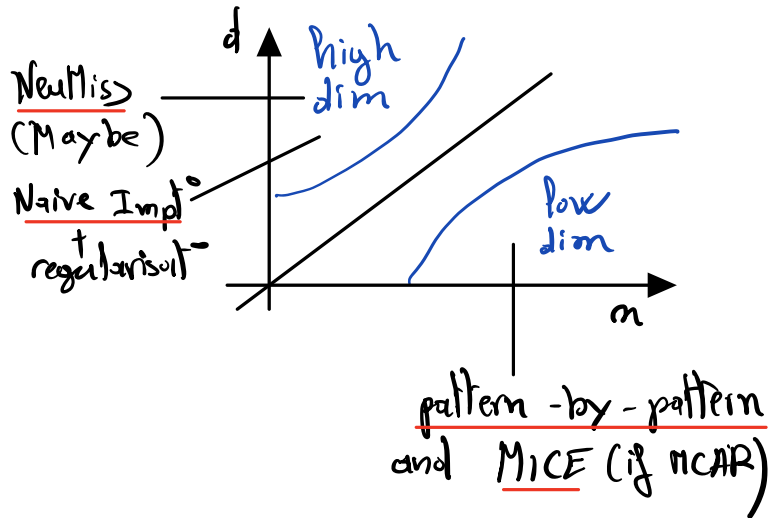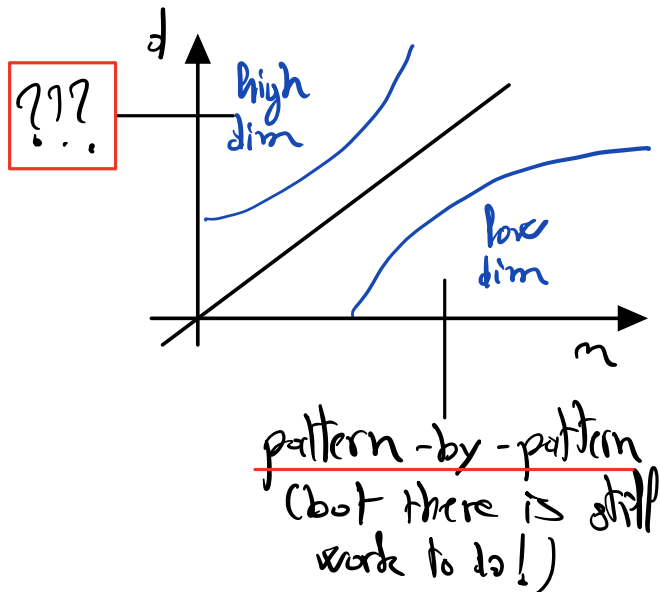
# Numerical experiments



**Excess risk** w.r.t. $n$ with $d = 8$

# What's work ?

## In practice



d

high dim

NewMiss (Maybe)

Naive Impl° + regularisat°

low dim

n

pattern-by-pattern and MICE (if MCAR)

## In Theory



???
...

d

high dim

low dim

n

pattern-by-pattern (but there is still work to do!)

# Conclusion

**Theoretical contributions**

o New thresholded predictor

o Adaptative upper bound

o Near optimal

**Numerical experiments**

o Thresholded P-by-P predictor:
  o reduced variance
  o consistent regardless of the missing scenario

$$\hat{f}(Z) = \sum_{m \in \{0,1\}^d} \hat{f}_m(X_{obs(m)}) \mathbf{1}_{M=m} \mathbf{1}_{\hat{p}_m > \frac{d}{n}}$$

$$\sigma^2 \mathfrak{C}_p\left(\frac{1}{n}\right) \lesssim \mathscr{E}_{\text{mini}}(\mathscr{P}) \leq A \log(n) \mathfrak{C}_p\left(\frac{d}{n}\right)$$



**Excess risk** w.r.t. $n$ with $d = 8$

# Near-optimal rate of consistency for linear prediction with missing values

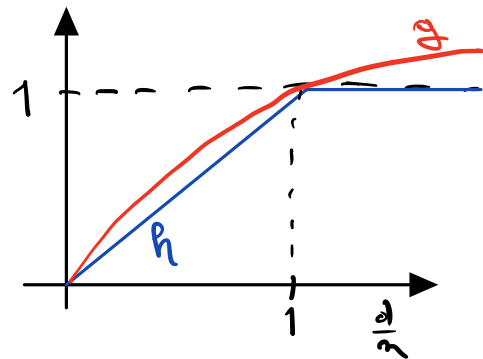**Corresponding Author:** alexis.ayme@sorbonne-universite.fr

$$\boxed{\begin{array}{l} \text{Why } c_p \\ \text{is adaptive?} \end{array}}$$

$$c_p\left(\frac{d}{n}\right) = \sum_{m \in \{0,1\}^d} p_m \wedge \frac{d}{n}$$

$$= \sum_m p_m \left(1 \wedge \frac{d}{n \, p_m}\right)$$

$$\underbrace{\phantom{\left(1 \wedge \frac{d}{n \, p_m}\right)}}_{h\left(\frac{d}{n \, p_m}\right)}$$

$$\leq \sum_m p_m \, g\left(\frac{d}{n \, p_m}\right)$$



**Two examples:**

1) $g = id \implies c_p\left(\frac{d}{n}\right) \leq \sum_m p_m \frac{d}{n \, p_m} = 2^d \frac{d}{n}$

2) $\begin{array}{l} g(x) = x^\alpha \\ \alpha \in ]0,1[ \end{array} \implies c_p\left(\frac{d}{n}\right) \leq \sum_m p_m^{1-\alpha} \left(\frac{d}{n}\right)^\alpha = \left(\frac{d}{n}\right)^\alpha e^{\alpha \, H_{1-\alpha}(p)}$

where $\boxed{H_{1-\alpha}(p) = \text{Renyi Entropy}}$