

# Linear prediction with **NA**, Imputation versus specific methods

Alexis Ayme

Under the supervision of:

Claire Boyer, Aymeric Dieuleveut and Erwan Scornet



# Background

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

○ Different sources:

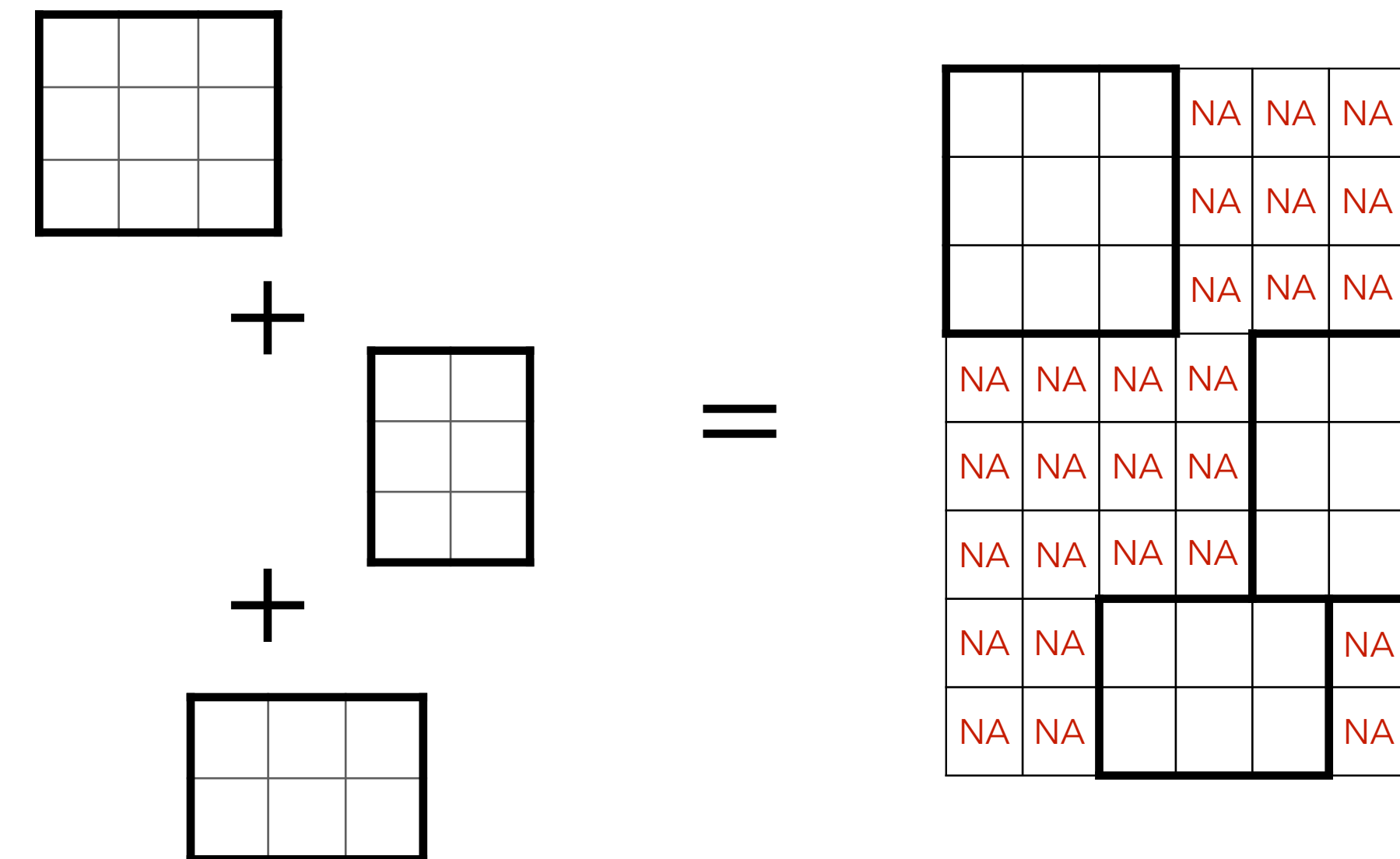
1. Bugs
2. Cost
3. Multiplication of sources (i.e. merging)
4. Sensitive data

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

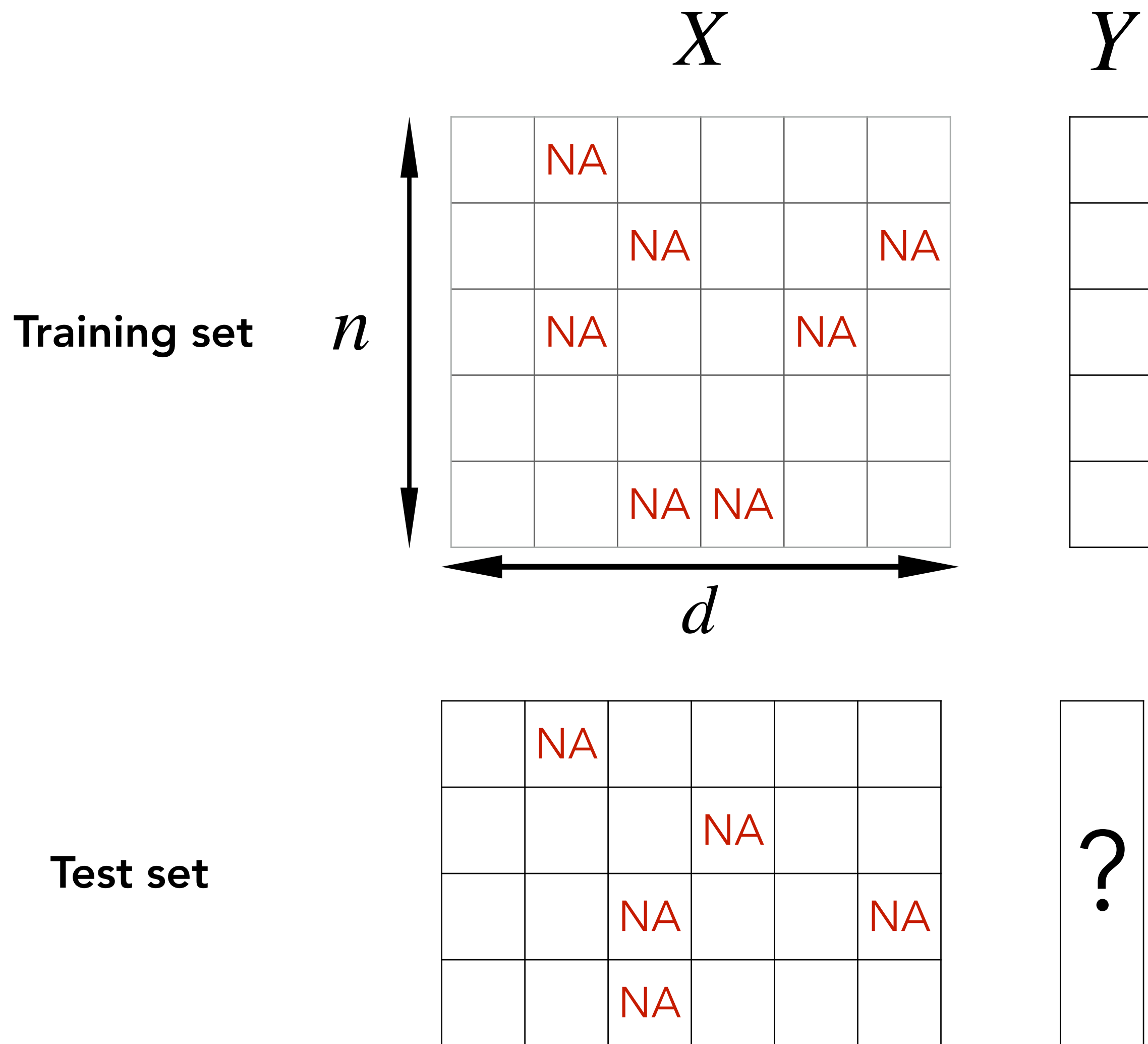
| \$1 | \$10 | \$100 | \$0 |
|-----|------|-------|-----|
|     |      | NA    |     |
|     |      |       |     |
|     | NA   | NA    |     |
|     |      | NA    |     |
|     | NA   | NA    |     |

○ Growing mass of data => **High-dimensional** dataset

1. Cost
2. Multiplication of sources (i.e. merging)
3. Genotype, text



# Introduction: Supervised learning with missing values (NA)

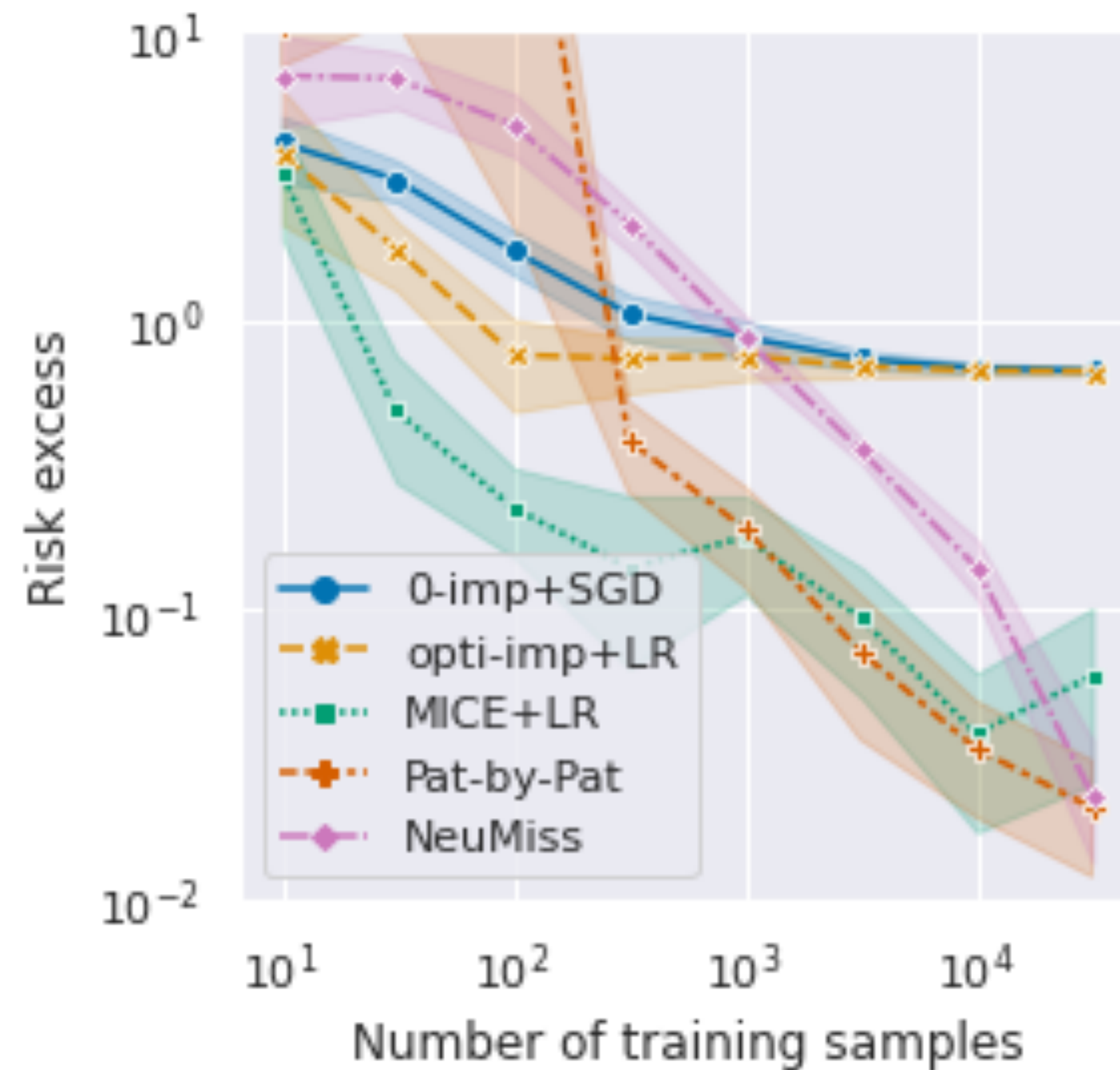


# Introduction: Supervised learning with missing values (NA)

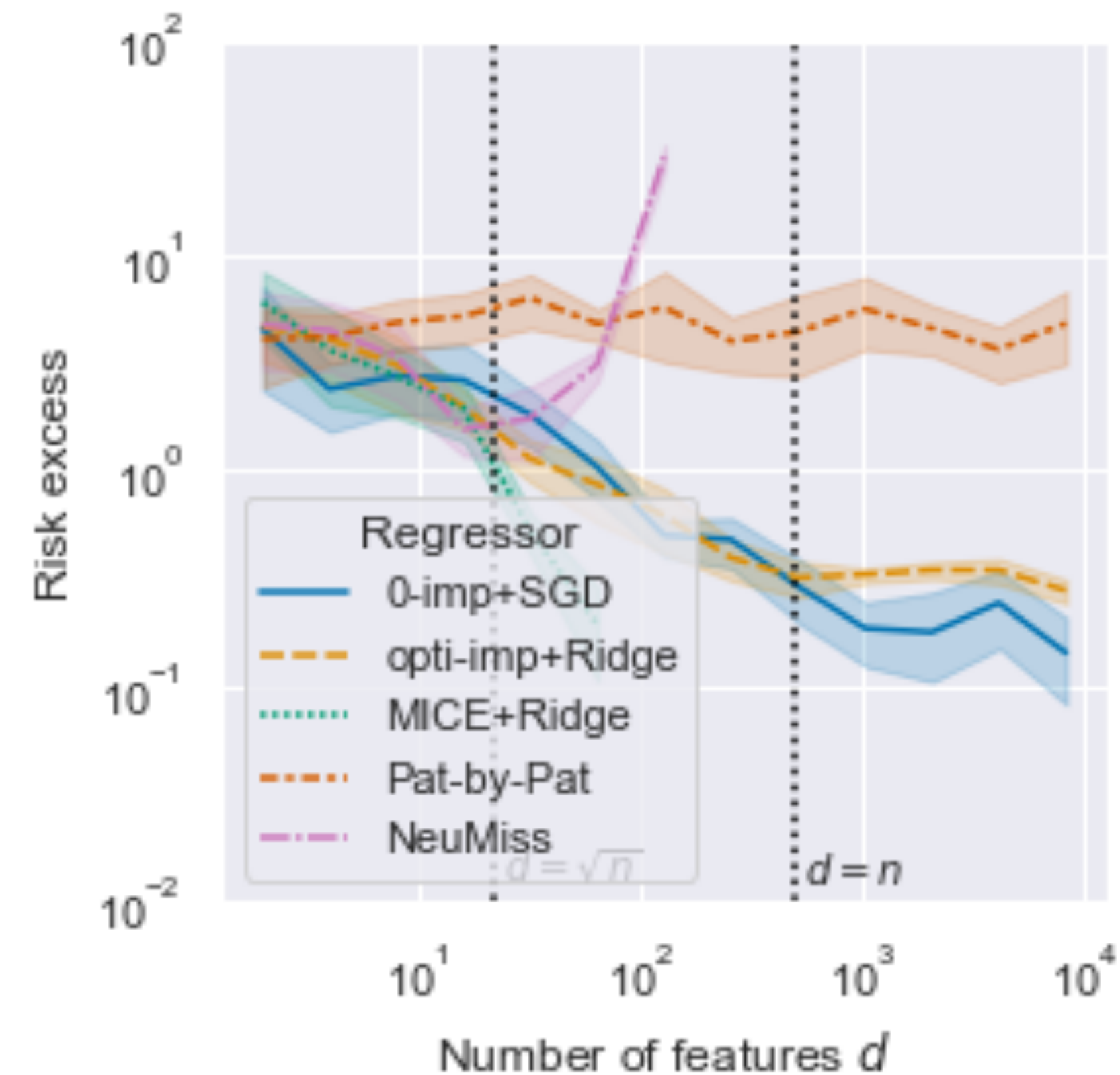
○ Handle missing values with:

1. Impute then regress procedure
2. Specific method

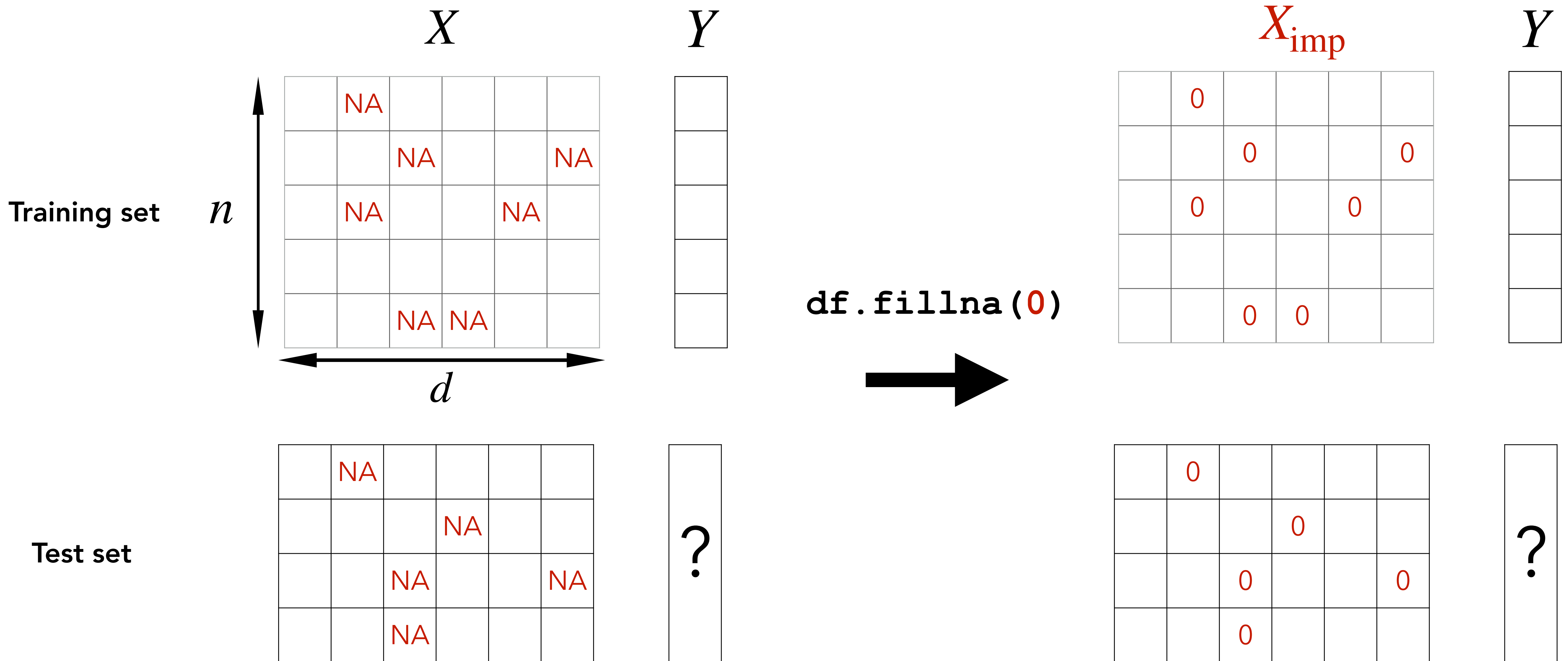
○ Low dimension  $n \rightarrow +\infty$



○ High dimension  $d \rightarrow +\infty$

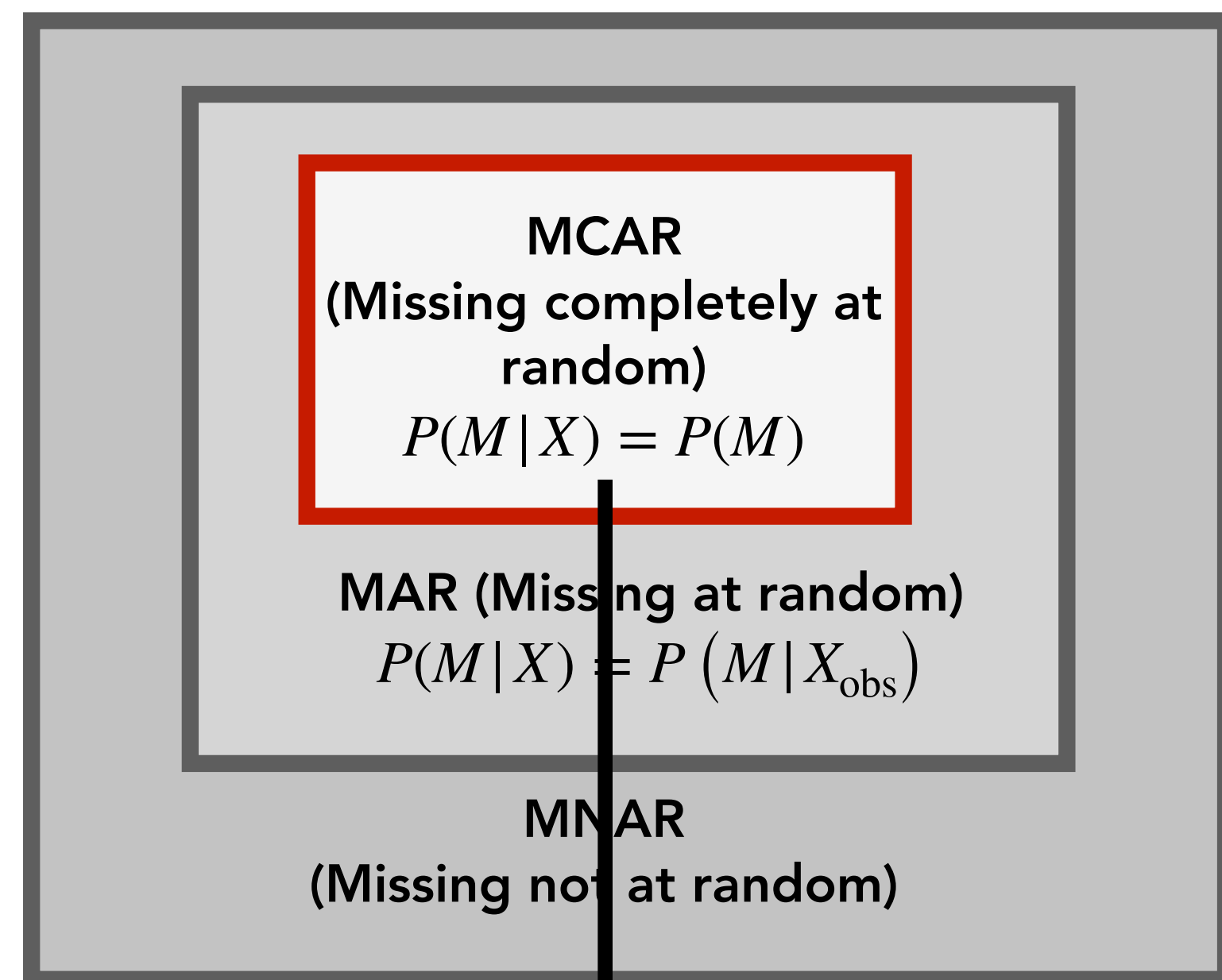


# Introduction: Naive imputation



# Introduction

- Missing values



**Bernoulli Model: Missing values**

$$\text{i.i.d } M_1, \dots, M_d \sim \mathcal{B}(1 - \rho)$$

- In this talk:

1) High-dimensional dataset

$$d \gg n$$

2) Predict on **imputed data**

$$R_{\text{imp}}(f) = \mathbb{E}_{X,Y} \left[ \left( Y - f(X_{\text{imp}}) \right)^2 \right]$$

3) Comparison with the complete case

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

# Definition

○ **Linear model** (well/miss-specified):  $\beta \in \mathbb{R}^d, \mathbb{E}[\epsilon X] = 0$

$$Y = \theta_{\star}^{\top} X + \epsilon$$

○ **Bayes risk:**

$$R^{\star} = \inf_{\theta} R(\theta)$$

$$R_{\text{imp}}^{\star} = \inf_{\theta} R_{\text{imp}}(\theta)$$

**Imputation bias:**

$$B_{\text{imp}} = R_{\text{imp}}^{\star} - R^{\star}$$

# Imputation by 0 = implicit ridge?

## ○ Ridge penalization

$$R_\lambda(\theta) = R(\theta) + \lambda \|\theta\|_2^2$$

**Theorem 2:** Under Bernoulli model and  $\sum_{j,j} = 1$  for all  $j \in [d]$ ,

$$R_{\text{imp}}(\theta) = R(\rho\theta) + \rho(1 - \rho)\|\theta\|_2^2$$

1. Optimal predictor has a small norm
2. **TAKE AT HOME: Imputation induce a Ridge penalization**

## ○ Ridge bias

$$B_{\text{ridge},\lambda} = \inf_{\theta} \{R(\theta) - R(\theta_\star) + \lambda \|\theta\|_2^2\}$$

**Theorem 2:** Under Bernoulli model and  $\sum_{j,j} = 1$  for all  $j \in [d]$ ,

$$B_{\text{imp}} = B_{\text{ridge},\lambda_{\text{imp}}}$$

where  $\lambda_{\text{imp}} = \frac{\rho}{1 - \rho}$

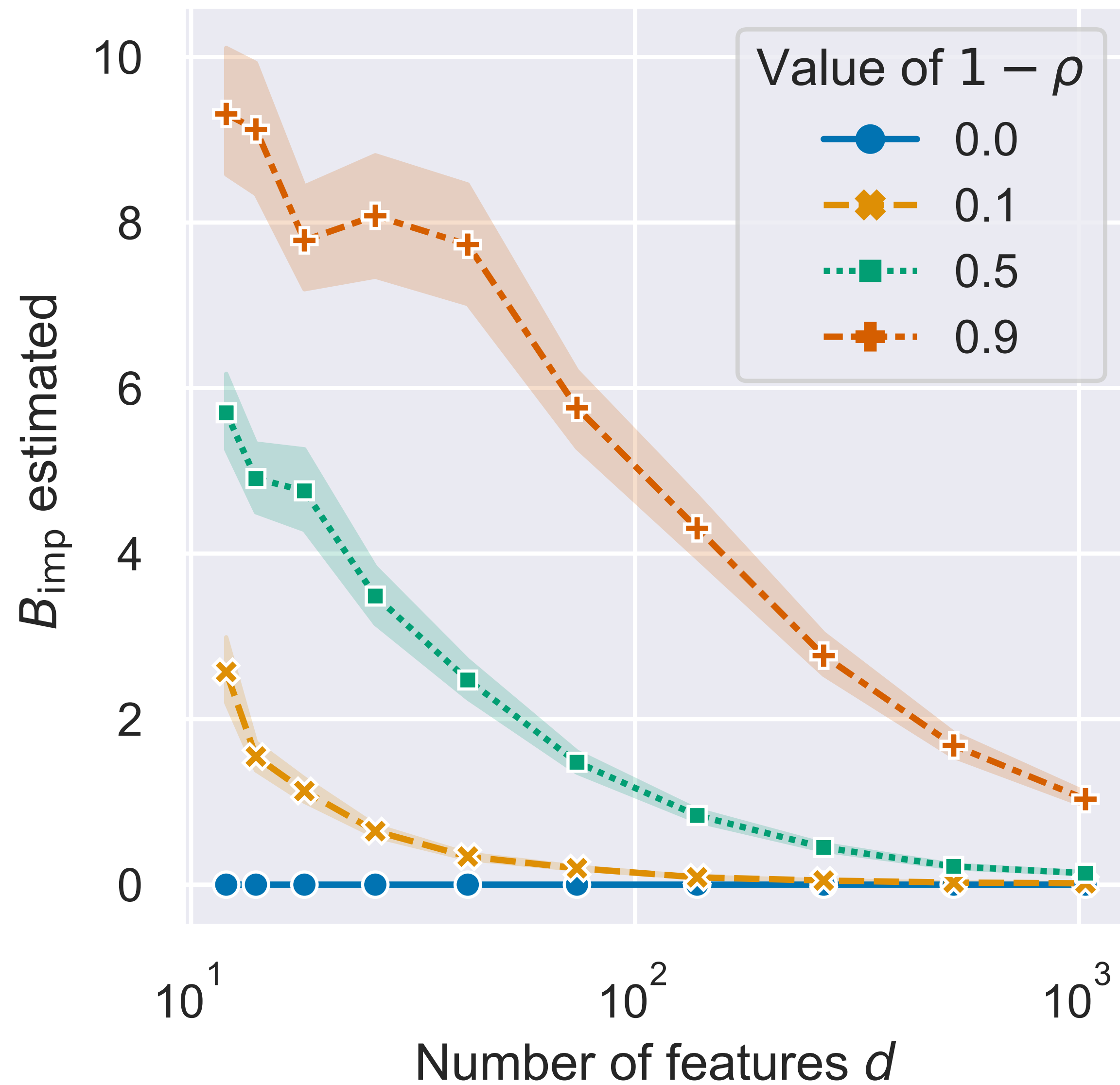
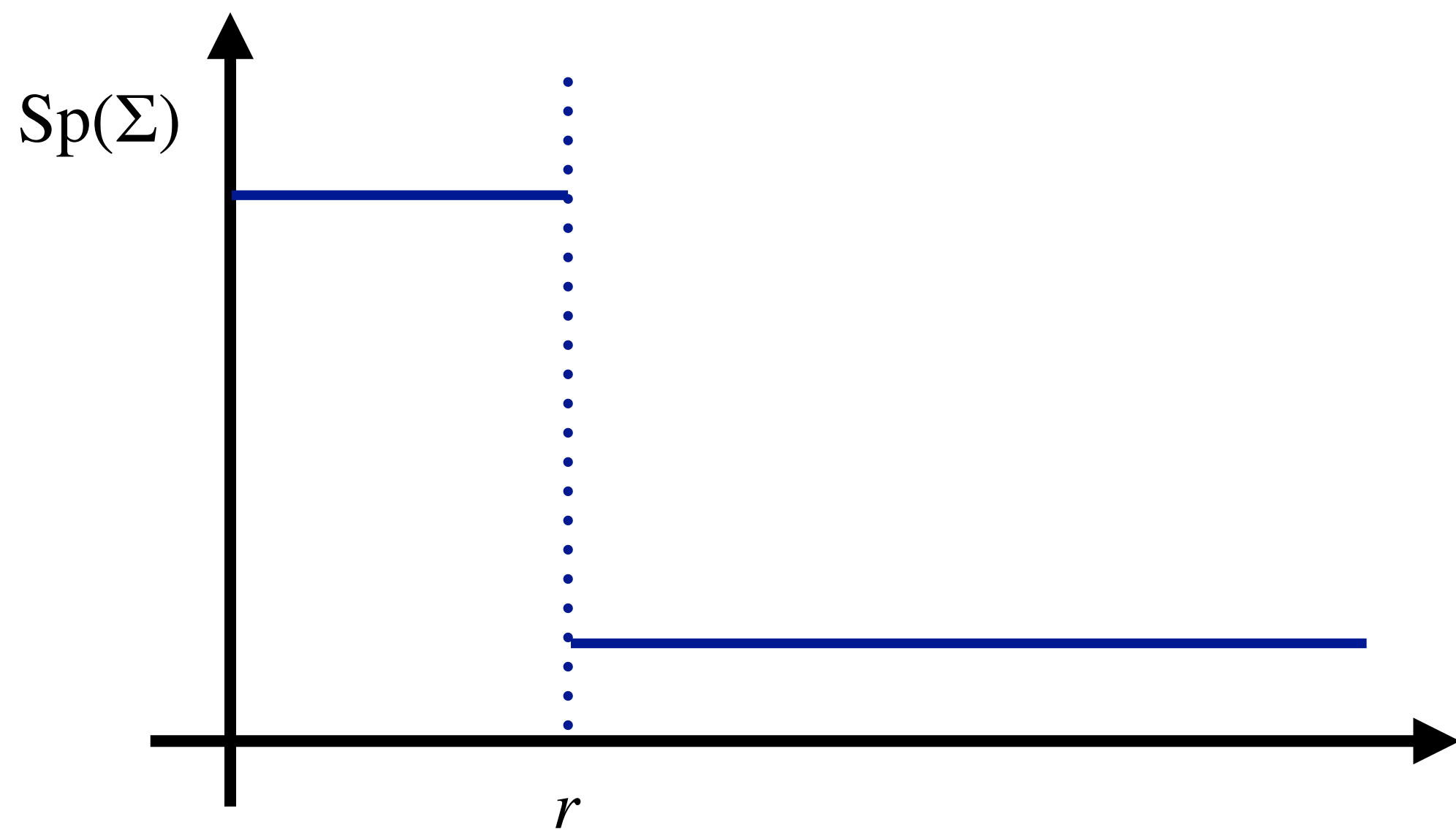
1. No strong assumptions on  $X$
2.  $\lambda_{\text{imp}}$  depends only on  $1 - \rho$  the proportion of missing values.
3. Available for all MCAR setting with another  $\lambda_{\text{imp}}$
4. Bias decreases with the dimension
5. **TAKE AT HOME: MCAR missing values seem to be at the same price of Ridge penalization**



# Illustration: on low rank data

○ Low rank data (or Spiked):  $\text{rank}(\Sigma) \approx r$

$$B_{\text{imp}} \lesssim \frac{r}{d} \mathbb{E}Y^2$$



# Learn imputed data with SGD

○ **SGD recursion:** with learning rate  $\gamma = \frac{1}{d\sqrt{n}}$

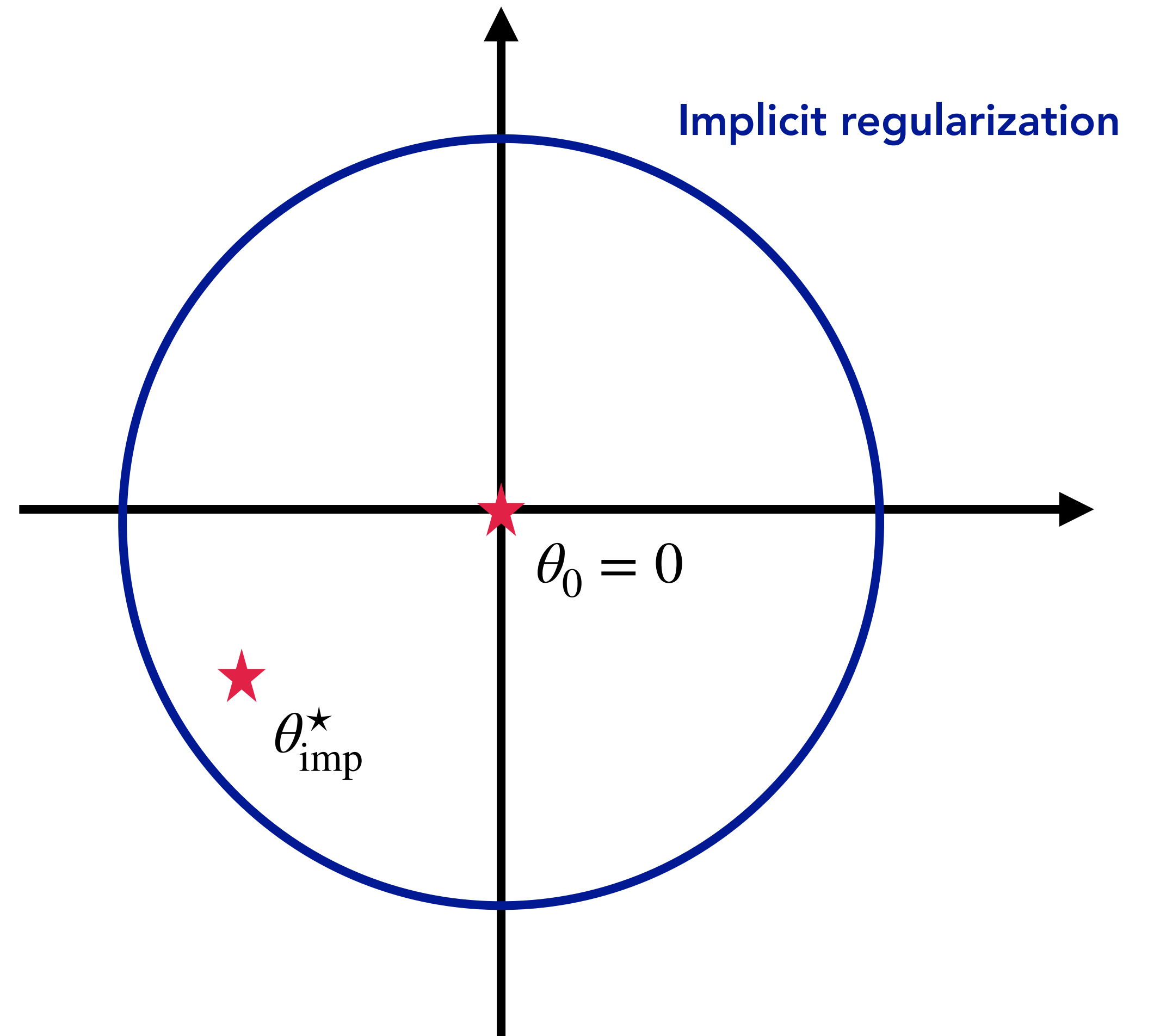
$$\begin{cases} \theta_0 = 0 \\ \theta_{\text{imp},t} = \left[ I - \gamma X_{\text{imp},t} X_{\text{imp},t}^\top \right] \theta_{\text{imp},t-1} + \gamma Y_t X_{\text{imp},t} \end{cases}$$

○ **Polyak Ruppert average:**

$$\bar{\theta}_{\text{imp},n} = \frac{1}{n+1} \sum_{t=1}^n \theta_{\text{imp},t}$$

**Theorem 2:** Under classical SGD assumptions,

$$\mathbb{E} \left[ R_{\text{imp}} \left( \bar{\theta}_{\text{imp},n} \right) \right] - R^* \leq B_{\text{imp}} + \frac{d}{\sqrt{n}} \|\theta_{\text{imp}}^*\|_2^2 + \frac{\sigma^2}{\sqrt{n}}$$



# Learn imputed data with SGD

## ○ SGD recursion:

$$\begin{cases} \theta_0 = 0 \\ \theta_{\text{imp},t} = \left[ I - \gamma X_{\text{imp},t} X_{\text{imp},t}^\top \right] \theta_{\text{imp},t-1} + \gamma Y_t X_{\text{imp},t} \end{cases}$$

## ○ Polyak Ruppert average:

$$\bar{\theta}_{\text{imp},n} = \frac{1}{n+1} \sum_{t=1}^n \theta_{\text{imp},t}$$

## ○ Illustration on low rank:

$$\mathbb{E} \left[ R_{\text{imp}} \left( \bar{\theta}_{\text{imp},n} \right) \right] - R^* \leq \left( \frac{1}{\rho\sqrt{n}} + \frac{1-\rho}{d} \right) \frac{r}{\rho} \mathbb{E} Y^2 + \frac{\sigma^2}{\sqrt{n}}$$

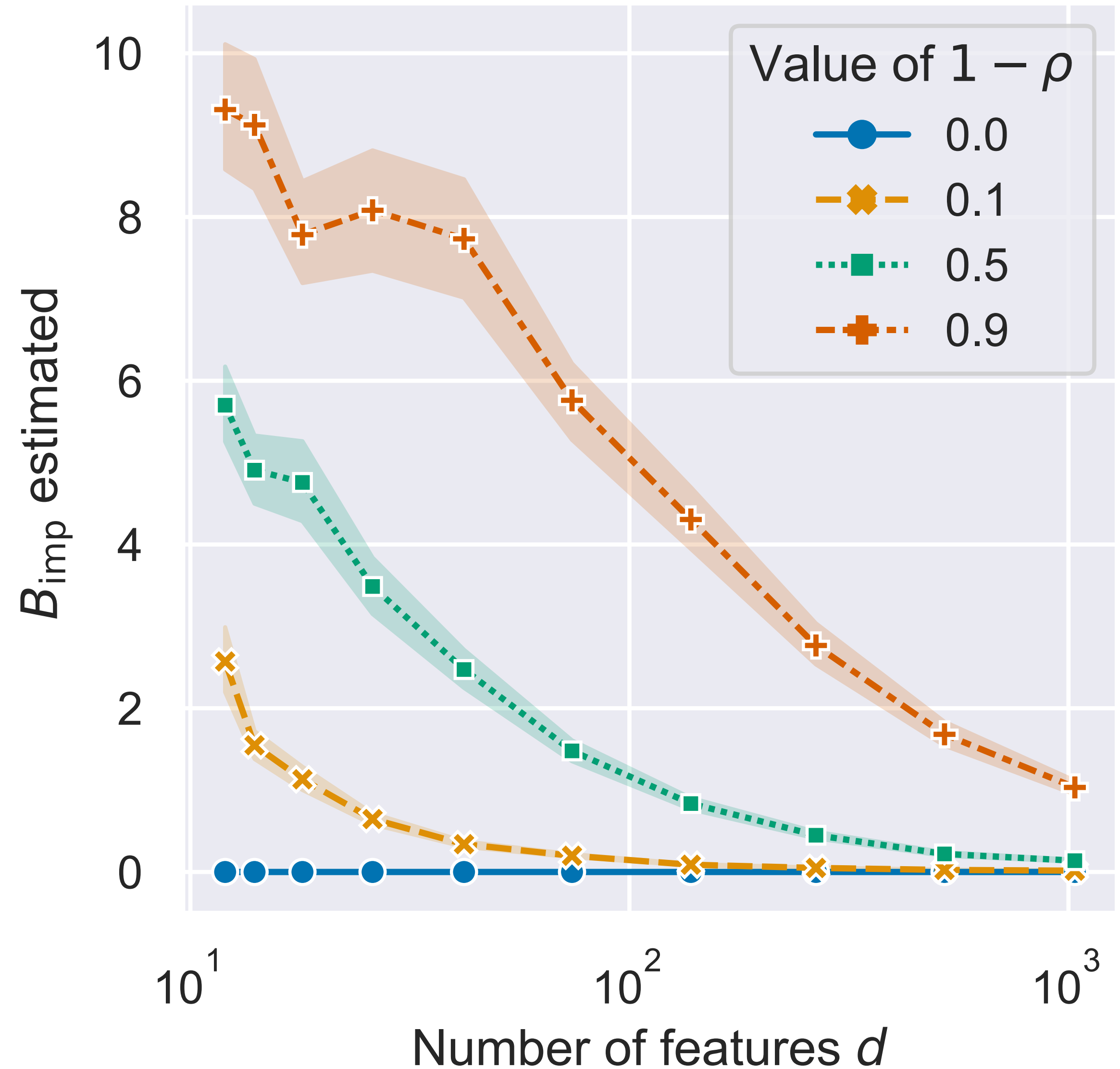
1. We leverage on implicit regularization.
2. Streaming online (one passe)
3. Trade-off between imputation bias and initial condition.
4. **TAKE AT HOME: Imputation bias vanishes for  $d \gg \sqrt{n}$**

**Theorem 2:** Under classical SGD assumptions,

$$\mathbb{E} \left[ R_{\text{imp}} \left( \bar{\theta}_{\text{imp},n} \right) \right] - R^* \leq B_{\text{imp}} + \frac{d}{\sqrt{n}} \|\theta_{\text{imp}}^2\|_2^2 + \frac{\sigma^2}{\sqrt{n}}$$

# Conclusion

1. Imputation (even very cheap) out-performs specific method to handle missing values in high-dimension.
2. Imputation by 0 induce a Ridge penalization.
3. Imputation bias vanishes with dimension as a consequence missing values are not an issue in high dimension.
4.  $d \gg \sqrt{n}$  regime leads to slow rates of consistency.



# Toy example

- Complete Model:

$$Y = X_1 .$$

$$X = (X_1, X_1, \dots, X_1)$$

$$M_1, \dots, M_d \sim \mathcal{B}(1/2).$$

$$R^* = 0$$

- With imputed missing values:

$$\theta_1 = (1, 0, \dots, 0)^\top$$

$$\theta_2 = 2(1/d, 1/d, \dots, 1/d)^\top$$

$$\theta_1^\top X_{\text{imp}} = X_1 M_1$$

$$\theta_2^\top X_{\text{imp}} = \frac{2X_1}{d} \sum_j M_j$$

$$R(\theta_1) = \frac{1}{2} \mathbb{E}[X_1^2]$$

$$R(\theta_2) = \frac{1}{d} \mathbb{E}[X_1^2]$$

$$B_{\text{imp}} = R^* - R_0^* \leq \frac{1}{d} \mathbb{E}[X_1^2]$$