

# Debiasing Stochastic Gradient Descent to handle missing values

Aude Sportisse<sup>1</sup> Claire Boyer<sup>1</sup> Aymeric Dieuleveut<sup>2</sup> Julie Josse<sup>3</sup>

<sup>1</sup>Sorbonne Université <sup>2</sup>Ecole Polytechnique <sup>3</sup>INRIA

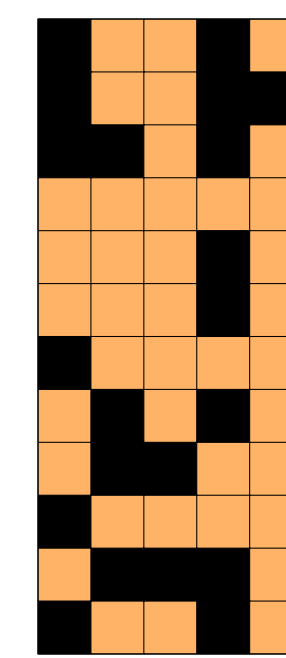


## Context

- **Large-scale** data analysis: large number of observations  $n$ , large dimension of the observations  $p$ .
- **Missing data** occurs more frequently. E.g. in clinical data: failure of the measuring device, no time to measure in emergency situations, aggregating datasets from multiple hospitals, etc.
- **Stochastic gradient descent** (SGD): key role in machine learning.

## Missing data

- Problem: Missing values in the covariates  $X_k$ .
- $D_k \in \{0, 1\}^d$ ,  $D_{kj} = \begin{cases} 0 & \text{if the var. } j \text{ of obs. } k \text{ missing} \\ 1 & \text{otherwise.} \end{cases}$
- **Heterogeneous MCAR** data:  $\neq$  missing proba. for each covariate.  $D = (\delta_{kj})$  with  $\delta_{kj} \sim \mathcal{B}(p_j)$ .
- Access to  $X_k^{\text{NA}} \in (\mathbb{R} \cup \{\text{NA}\})^d$  instead of  $X_k$ .



$$X_k^{\text{NA}} := X_k \odot D_k + \text{NA} \odot (\mathbf{1}_d - D_k),$$

## Setting: linear regression with missing covariates

$(X_k^{\text{NA}}, y_k) \in \mathbb{R}^d \times \mathbb{R}$  i.i.d. observations.

$$y_k = (X_k^{\text{NA}})^T \beta^* + \epsilon_k,$$

parametrized by  $\beta^* \in \mathbb{R}^d$ , with a noise term  $\epsilon_k \in \mathbb{R}$ .

## How to perform linear regression with missing covariates that handle large-scale or streaming data?

### Existing works in linear models

- Expectation Maximization algorithm [1].
- **X** parametric (Gaussian) assumption for the covariates.
- Naive imputation e.g. by the mean [2].
- **X** Bias in the estimate.
- **Imputing naively** by 0 and modifying SGD to **account for the imputation error** ([3]), also in [4], for homogeneous MCAR values.
- ⇒ **Our proposal**: debiased **averaged** SGD, better rate of convergence.

## Methodology

- Imputing the missing values by 0:  $\tilde{X}_k = X_k^{\text{NA}} \odot D_k = X_k \odot D_k$ .
- Using a **debiased gradient** for the **averaged** SGD:  

$$\tilde{g}_k(\beta_{k-1}) = P^{-1} \tilde{X}_k^T ( \tilde{X}_k^T P^{-1} \beta_{k-1} - y_k ) - (I - P) P^{-2} \text{diag}(\tilde{X}_k, \tilde{X}_k^T) \beta_{k-1},$$
 where  $P = \text{diag}((p_j)_{j \in \{1, \dots, d\}}) \in \mathbb{R}^{d \times d}$ .
- **Averaged iterates**:  $\tilde{\beta}_k = \frac{1}{k+1} \sum \beta_i$ .

## Theoretical results

- Goal: **establish a convergence rate**.
- Assumptions:  $(X_k, y_k) \in \mathbb{R}^d \times \mathbb{R}$  i.i.d.,  $\mathbb{E}[\|X_k\|^2]$  and  $\mathbb{E}[y_k^2]$  finite,  $H := \mathbb{E}_{(X_k, y_k)}[X_k X_k^T]$  invertible,  $\mathcal{F}_k = \sigma(X_1, y_1, D_1, \dots, X_k, y_k, D_k)$ .

### Lemma 2: structured noise induced by NA

- $\tilde{g}_k(\beta^*)$  is  $\mathcal{F}_k$ -measurable and  $\forall k \geq 0$ ,
- $\mathbb{E}[\tilde{g}_k(\beta^*) | \mathcal{F}_{k-1}] = 0$  a.s.,  $\mathbb{E}[\|\tilde{g}_k(\beta^*)\|^2 | \mathcal{F}_{k-1}]$  is a.s. finite.
- $\mathbb{E}[\tilde{g}_k(\beta^*) \tilde{g}_k(\beta^*)^T] \preceq C(\beta^*) = c(\beta^*) H$  a.s..

### Lemma 3: $(\tilde{g}_k(\beta^*))_{k \geq 0}$ a.s. co-coercive

- For any  $k \geq 0$ ,  $\tilde{g}_k$  is  $L_{k,D}$ -Lipschitz.
- There exists a primitive function  $\tilde{f}_k$  which is a.s. convex.

### Theorem 1: convergence rate, online streaming

Assume that for any  $k$ ,  $\|X_k\| \leq \gamma$  a.s. for some  $\gamma > 0$ . For **any constant step-size**  $\alpha \leq \frac{1}{2L}$  and for any  $k \geq 0$ , one has:

$$\mathbb{E}[R(\tilde{\beta}_k) - R(\beta^*)] \leq \frac{1}{2k} \cdot \left( \underbrace{\frac{\sqrt{c(\beta^*)d}}{1 - \sqrt{\alpha L}}}_{\text{variance term}} + \underbrace{\frac{\|\beta_0 - \beta^*\|}{\sqrt{\alpha}}}_{\text{bias term}} \right)^2,$$

- $L := \sup_{k,D}$  Lipschitz constants of  $\tilde{g}_k$ .
- $p_m = \min_{j=1, \dots, d} p_j$  minimal probability to be observed.
- $c(\beta^*) = \frac{\text{Var}(\epsilon_k)}{p_m^2} + \frac{\text{multiplicative noise (induced by naive imputation)}}{\left( \frac{(2 + 5p_m)(1 - p_m)}{p_m^3} \right) \gamma^2 \|\beta^*\|^2}$ .  
*increasing with the missing values rate*

- ✓ **Optimal rate for least-squares regression**:  $\mathcal{O}(k^{-1})$ .
- ✓ In the complete case ( $p_m = 1$ ): same bound as Bach and Moulines [5].

## Additional results

- **Finite-sample setting**:  $n$  is fixed.
- **True risk**: same convergence rate holds for **only one epoch**.  
**X** if we use the data more than once: bias in the gradient.
- **Empirical risk (open issue)**:  $\beta_n^* = \arg \min_{\beta \in \mathbb{R}^d} \{R_n(\beta) := \frac{1}{n} \sum f_i(\beta)\}$ .  
**X** data used several times or non-uniform sampling.
- **Using estimated missing probabilities**  $(\hat{p}_j)_j$  in our algorithm instead of  $(p_j)_j$  preserves the same order of convergence rate  $\mathcal{O}(k^{-1})$ .
- **Ridge Regression**:  $\beta \rightarrow R(\beta) + \lambda \|\beta\|^2$  is  $2\lambda$ -strongly convex: no change for the debiased gradient, convergence rate of  $\mathcal{O}((\lambda k)^{-1})$ .

## Consequences in practice

- **Before collecting data, fewer complete obs. is better than more incomplete ones**, e.g. variance bound for 200 incomplete obs. (50% NA) is twice as large as for 100 complete obs.
- **After collecting data with NA, obs. containing NA should not be removed**: the upper-bound is  $p^{d-3}$  smaller than the lower bound of any algorithm relying only on the complete observations.

## Experiments on synthetic data

- $X_i$  i.i.d.  $\mathcal{N}(0, \Sigma)$ , where  $\Sigma$  generated randomly with decreasing eigenvalues,  $y_i = X_i \beta + \epsilon_i$ , for  $\beta$  fixed and  $\epsilon_i \sim \mathcal{N}(0, 1)$ .
- $d = 10$ , 30% missing values,  $L$  oracle value.

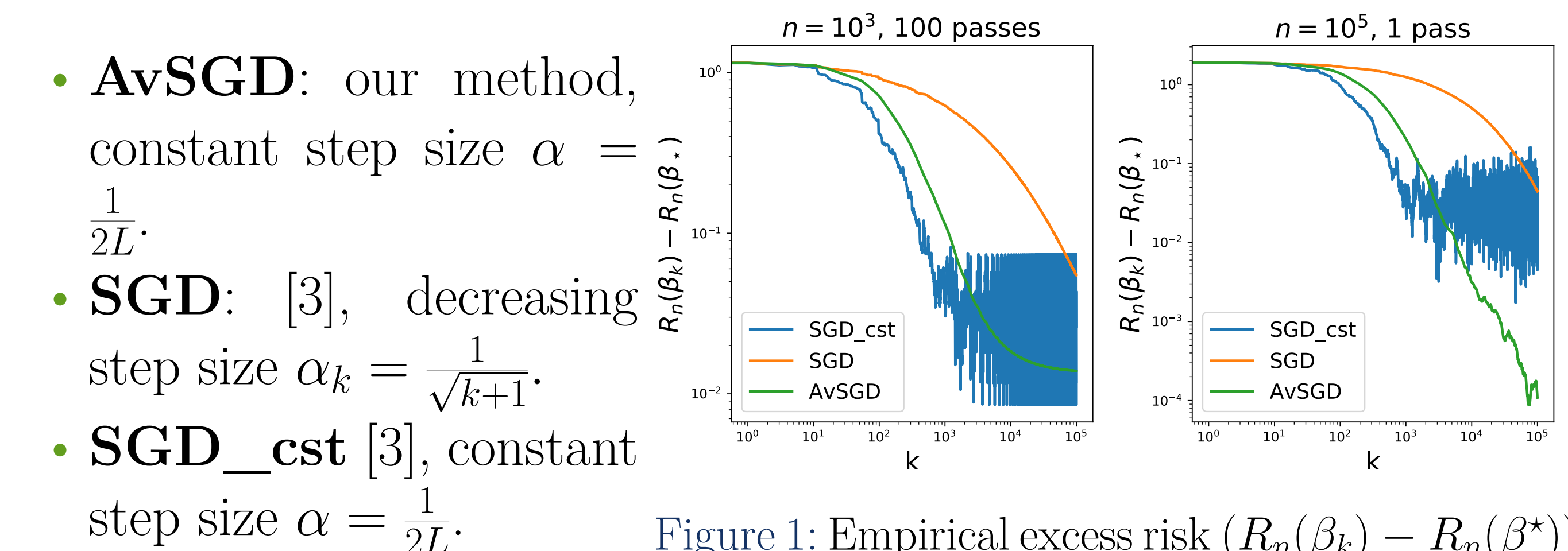


Figure 1: Empirical excess risk  $(R_n(\beta_k) - R_n(\beta^*))$ .

## Experiments on real data

- **Complete** dataset: 81 quant. features, 21263 superconductors.
- Introduction of 30% of heterogeneous MCAR values.
- Training/test split, with no NA in the test set.
- $\hat{y}_{n+1} = X_{n+1}^T \hat{\beta}$ , with  $\hat{\beta}$  computed on the training set, with **AvSGD** or with a two-steps procedure **Mean+AvSGD**.

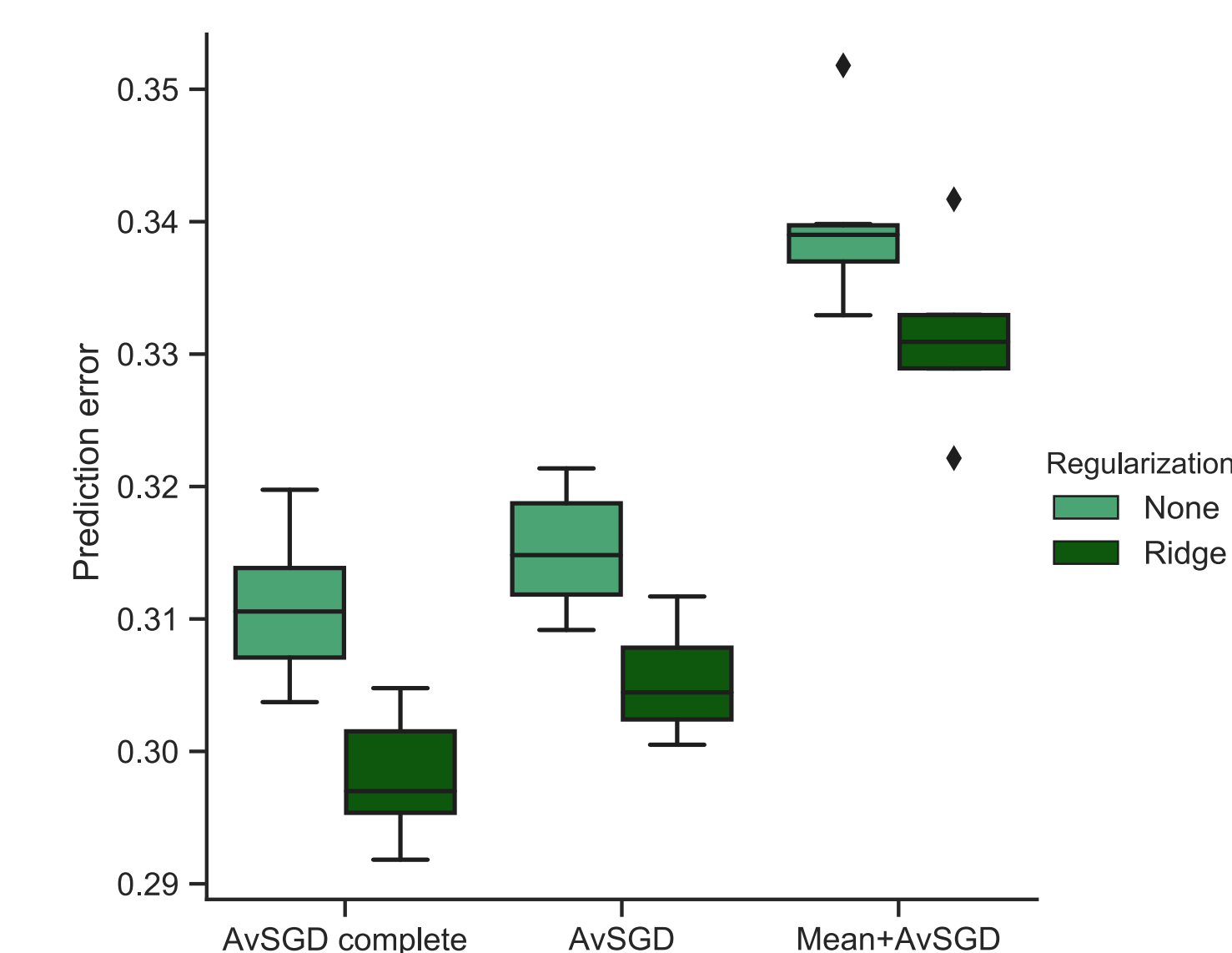


Figure 2: Prediction error  $\| \hat{y} - y \|^2 / \| y \|^2$ .

⇒ **Further research**: MAR values, GLM models.

## References

- [1] A. P. Dempster, N. M. Laird, and D. B. Rubin. Stochastic gradient descent for linear systems with missing data.
- [2] R. J. A. Little and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Statistical analysis with missing data*.
- [3] Anna Ma and Deanna Needell.
- [4] Po-Ling Loh and Martin J. Wainwright. High-dimensional regression with noisy and missing data: Provable guarantees with non-convexity.
- [5] F. Bach and E. Moulines. Non-strongly-convex smooth stochastic approximation with convergence rate  $\mathcal{O}(1/n)$ .