Breaking the Nonsmooth Barrier: A Scalable Parallel Method for Composite Optimization

Fabian Pedregosa[†], Rémi Leblond[†], Simon Lacoste–Julien^{*}









Summary

Optimization methods need to be adapted to the parallel setting to leverage modern computer architectures.

Highly efficient variants of stochastic gradient descent have been recently proposed, such as Hogwild [1], Kromagnon [2], ASAGA [3].

They assume that the objective function is smooth, so are inapplicable to problems such as Lasso, optimization with constraints, etc.

Contributions:

- . Sparse Proximal SAGA, a sparse variant of the linearly-convergent proximal SAGA algorithm.
- 2. ProxASAGA, the first parallel asynchronous variance-reduced method that supports composite objective functions.

Problem setting

Objective: develop parallel asynchronous method for problems of the

$$\underset{oldsymbol{x} \in \mathbb{R}^p}{\operatorname{minimize}} \frac{1}{n} \sum_{i=1}^n f_i(oldsymbol{x}) + h(oldsymbol{x}),$$

- f_i is differentiable with L-Lipschitz gradient.
- ullet is block-separable $(h(x) = \sum_B h_B([oldsymbol{x}]_B))$ and "simple" in the sense that we have access to $\mathbf{prox}_h \stackrel{\mathsf{def}}{=} \mathbf{arg} \, \mathbf{min}_{m{x}} \, h(m{x}) + \frac{1}{2} \|m{x} - m{z}\|^2$.
- includes Lasso, group Lasso or ERM with box constraints.

Variance-reduced stochastic gradient methods are natural candidates: state of the art performance and recent asynchronous variants. The **SAGA** algorithm [4] has an iteration of the form

choose random
$$i \in \{1, \dots, n\}$$

$$\boldsymbol{x}^+ = \mathbf{prox}_{\gamma g}(\boldsymbol{x} - \gamma(\nabla f_i(\boldsymbol{x}) - \boldsymbol{\alpha}_i + \overline{\boldsymbol{\alpha}})), \ \boldsymbol{\alpha}_i^+ = \nabla f_i(\boldsymbol{x}). \tag{SAGA}$$

Difficulty of a composite extension

- Existing parallel asynchronous variants of SGD crucially rely on sparse updates.
- Even in the presence of sparse gradients, the (SAGA) update is not sparse due to the presence of $\overline{\alpha}$ and prox.
- Existing convergence proofs crucially rely on the gradient smoothness

A new algorithm: Sparse Proximal SAGA

The algorithm relies on the definitions

- Extended support
- $lackbox{ } oldsymbol{D}_i$
- \varphi

$$oldsymbol{v}_i =
abla f_i(oldsymbol{x}) - oldsymbol{lpha}_i + oldsymbol{D}_i \overline{oldsymbol{lpha}} \,, \,\, oldsymbol{x}^+ = \mathbf{prox}_{\gamma arphi_i}(oldsymbol{x} - \gamma oldsymbol{v}_i)$$
 (SPS)

Linear convergence rate

Theorem. Let $\gamma = \frac{a}{5L}$ for any $a \leq 1$ and f be μ -strongly convex ($\mu > 1$ 0). Then Sparse Proximal SAGA converges geometrically in expectation with a rate factor of at least $\rho = \frac{1}{5} \min\{\frac{1}{n}, a_{\kappa}^{1}\}$. That is, for x_t obtained after t updates, we have the following bound:

$$\mathbb{E}\|m{x}_t - m{x}^*\|^2 \le (1 -
ho)^t C_0$$
, with $C_0 := \|m{x}_0 - m{x}^*\|^2 + \frac{1}{5L^2} \sum_{i=1}^n \|m{lpha}_i^0 -
abla f_i(m{x})\|^2$

Proximal Asynchronous SAGA (ProxASAGA)

[†] INRIA and École Normale Supérieure, Paris, France. * MILA and DIRO, Université de Montréal, Canada

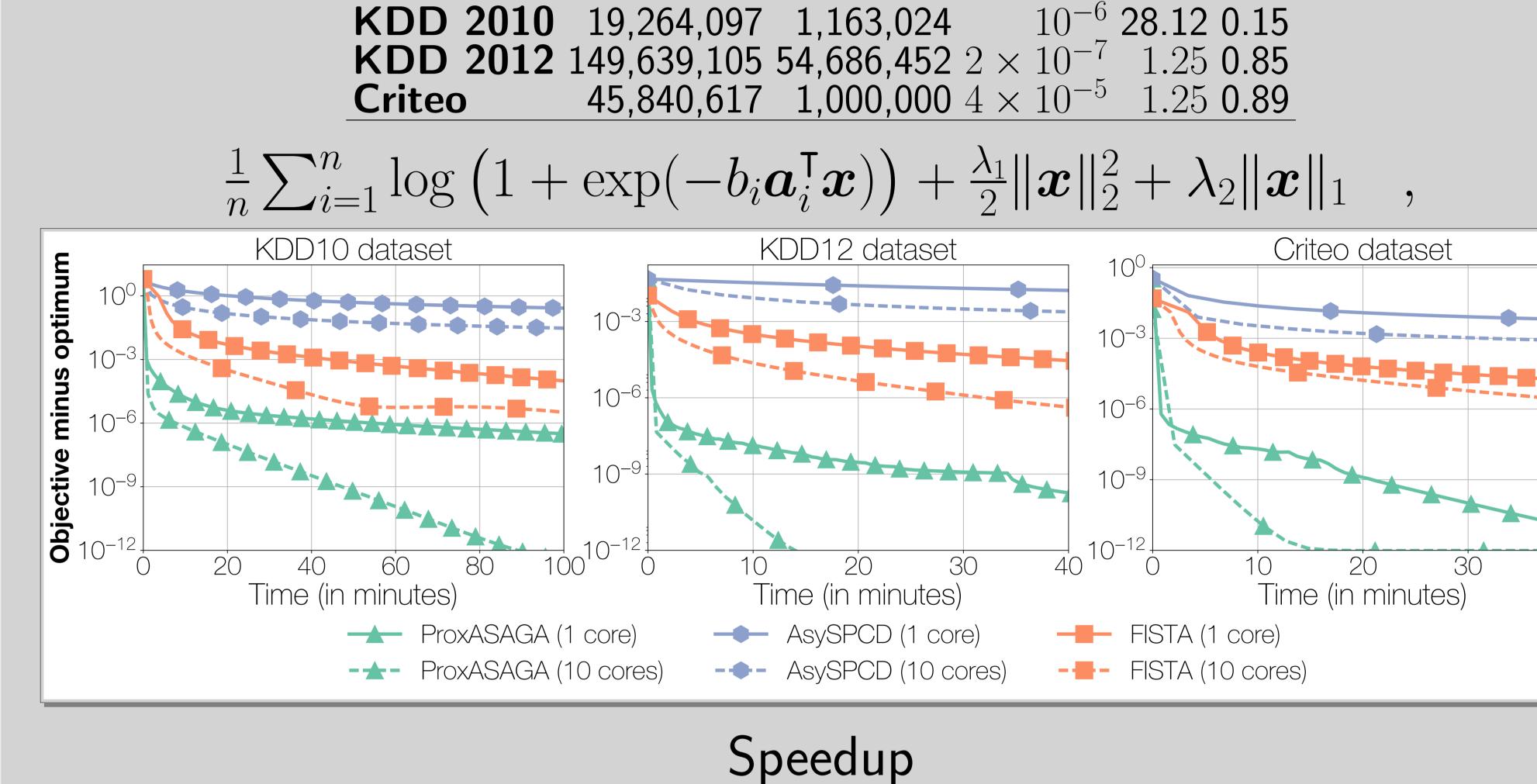
Algorithm 2 ProxASAGA (implemented) Algorithm 1 ProxASAGA (analyzed) Initialize shared variables $oldsymbol{x}$ and $(oldsymbol{lpha}_i)_{i=1}^n$ 1: Initialize shared variables \boldsymbol{x} , $(\boldsymbol{\alpha}_i)_{i=1}^n$, $\overline{\boldsymbol{\alpha}}$ loop 2: **loop** $\hat{m{x}}=\hat{m{x}}=\hat{m{x}}=\hat{m{x}}$ inconsistent read of $m{x}$: Sample i uniformly in $\{1, ..., n\}$ $\hat{m{lpha}} = ext{inconsistent read of } m{lpha}$ 4: $S_i := \text{support of } \nabla f_i$ Sample i uniformly in $\{1, ..., n\}$ 5: $T_i :=$ extended support of ∇f_i in $\mathcal B$ $S_i := \text{support of } \nabla f_i$ 6: $[\hat{m{x}}]_{T_i} = \text{inconsistent read of } {m{x}} \text{ on } T_i$ $T_i:= ext{extended support of } abla f_i:= ext{ in } \mathcal{B}$ $\hat{m{lpha}}_i = ext{inconsistent read of } m{lpha}_i$ 8: $[\overline{oldsymbol{lpha}}]_{T_i}=1/n\sum_{j=1}^n[\,\hat{oldsymbol{lpha}}_j\,]_{T_i}$ 8: $[\overline{m{lpha}}]_{T_i}=$ inconsistent read of $\overline{m{lpha}}$ on T_i $[\delta oldsymbol{lpha}]_{S_i} = [abla f_i(\hat{oldsymbol{x}})]_{S_i} - [\hat{oldsymbol{lpha}}_i]_{S_i}$ 9: $[\delta oldsymbol{lpha}]_{S_i} = [abla f_i(\hat{oldsymbol{x}})]_{S_i} - [\hat{oldsymbol{lpha}}_i]_{S_i}$ 10: $[\hat{oldsymbol{v}}]_{T_i} = [\delta oldsymbol{lpha}]_{T_i} + [oldsymbol{\Delta}_i \overline{oldsymbol{lpha}}]_{T_i}$ 10: $[\hat{oldsymbol{v}}]_{T_i} = [\delta oldsymbol{lpha}]_{T_i} + [oldsymbol{\Delta}_i \overline{oldsymbol{lpha}}]_{T_i}$ 11: $[\, \delta oldsymbol{x} \,]_{T_i} = [\mathbf{prox}_{\gamma arphi_i} (\hat{oldsymbol{x}} - \gamma \hat{oldsymbol{v}})]_{T_i} - [\hat{oldsymbol{x}}]_{T_i}$ 11: $[\delta oldsymbol{x}]_{T_i} = [\mathbf{prox}_{\gamma_{\mathcal{O}_i}}(\hat{oldsymbol{x}} - \gamma \hat{oldsymbol{v}})]_{T_i} - [\hat{oldsymbol{x}}]_{T_i}$ for B in T_i do 12: for B in T_i do for $b \in B$ do 13: for b in B do 14: $[\boldsymbol{x}]_b \leftarrow [\boldsymbol{x}]_b + [\delta \boldsymbol{x}]_b$ 14: $[\boldsymbol{x}]_b \leftarrow [\boldsymbol{x}]_b + [\delta \boldsymbol{x}]_b$ > atomic if $b \in S_i$ then 15: if $b \in S_i$ then 16: $[oldsymbol{lpha}_i]_b \leftarrow [abla f_i(\hat{oldsymbol{x}})]_b$ 16: $[\overline{\boldsymbol{\alpha}}]_b \leftarrow [\overline{\boldsymbol{\alpha}}]_b + 1/n[\delta \boldsymbol{\alpha}]_b$ > atomic end if 17: end if end for 18: end for 19: end for 19: end for 20: // (' \leftarrow ' denotes shared memory update.) (scalar update) ⊳ atomic 20: $oldsymbol{lpha}_i \leftarrow abla f_i(\hat{oldsymbol{x}})$ 1: end loop 21: end loop

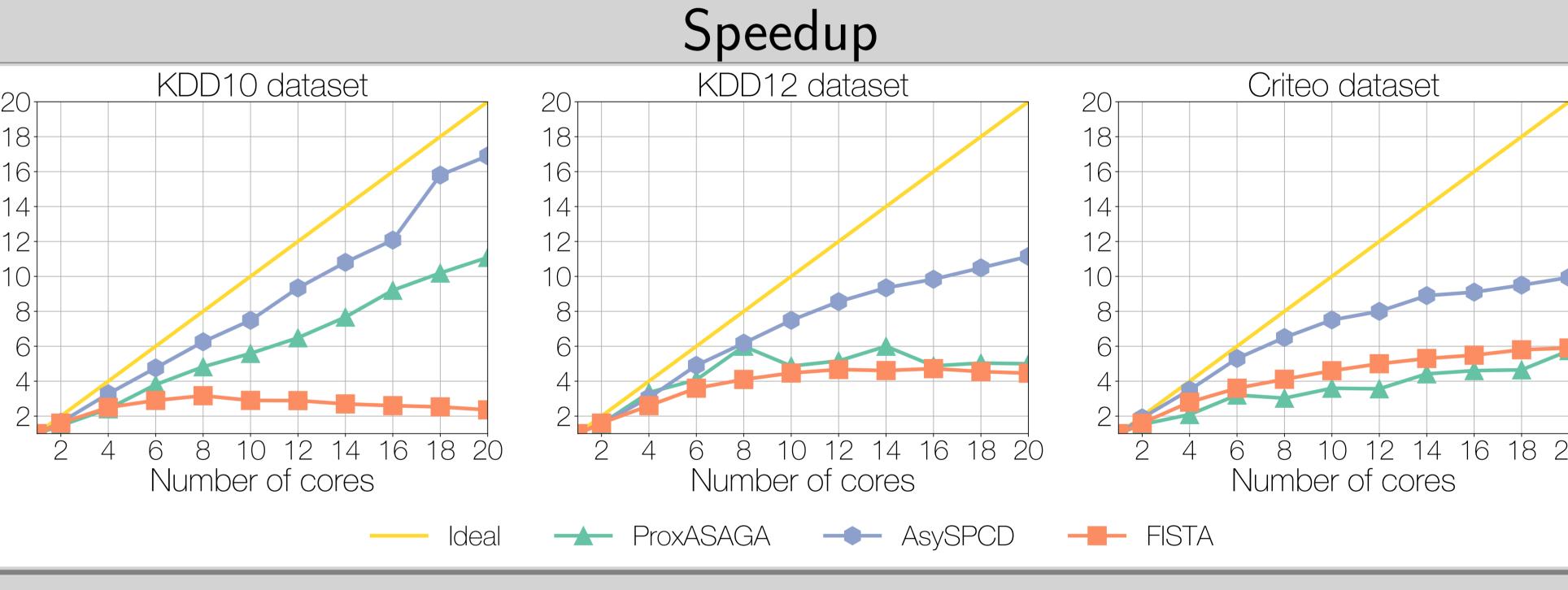
Experimental results

Comparison on 3 large-scale datasets

Dataset

Large-scale logistic regression





References

- Niu, F., Recht, B., Re, C. & Wright, S. Hogwild: A lock-free approach to parallelizing stochastic gradient descent.
- Mania, H. et al. Perturbed iterate analysis for asynchronous stochastic optimization. SIAM Journal on Optimization
- Leblond, R., Pedregosa, F. & Lacoste-Julien, S. ASAGA: asynchronous parallel SAGA. AISTATS (2017).
- Defazio, A., Bach, F. & Lacoste-Julien, S. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. in NIPS (2014).