

Doubly Sparse Asynchronous Learning for Stochastic Composite Optimization

Runxue Bao, Xidong Wu, Wenhan Xian, Heng Huang

Electrical and Computer Engineering, University of Pittsburgh, PA 15213, United States

IJCAI-ECAI 2022

Background

- ▶ We consider the composite optimization problem involving a data fitting function $\mathcal{F}(x) = \frac{1}{n} \sum_{i=1}^n f_i(a_i^\top x)$ plus a block-separable regularizer $\Omega(x) = \sum_{j=1}^q \Omega_j(x_{\mathcal{G}_j})$ as:

$$\min_{x \in \mathbb{R}^p} \mathcal{P}(x) := \mathcal{F}(x) + \lambda \Omega(x). \quad (1)$$

where $A = [a_1, \dots, a_n]^\top \in \mathbb{R}^{n \times p}$ is the design matrix, \mathcal{G} is the partition, λ is the regularization parameter and $x \in \mathbb{R}^p$ is the model coefficients.

- ▶ Most existing parallel learning methods focus on improving the algorithm efficiency in terms of sample complexity and thus suffer from high computation costs and memory burden in the high-dimensional setting.

Objective: To accelerate high-dimensional models by simultaneously enjoying the model sparsity and data sparsity.

Proposed Method

Algorithm 1 Sha-DSAL

- 1: **Input:** $x_{\mathcal{B}_0}^0 \in \mathbb{R}^p$, step size η , inner loops K .
- 2: **for** $s = 0$ **to** $S - 1$ **do**
- 3: All threads parallelly compute $\nabla \mathcal{F}(x_{\mathcal{B}_s}^0)$.
- 4: Compute dual variable y^s and update \mathcal{B}_{s+1} from \mathcal{B}_s by (2).
- 5: Update $A_{\mathcal{B}_{s+1}}, x_{\mathcal{B}_{s+1}}^0, \nabla \mathcal{F}(x_{\mathcal{B}_{s+1}}^0)$
- 6: For each thread, do:
- 7: **for** $t = 0$ **to** $K - 1$ **do**
- 8: Read $\hat{x}_{\mathcal{B}_{s+1}}^t$ from the shared memory.
- 9: Randomly sample i from $\{1, 2, \dots, n\}$.
- 10: Compute v_t^s by (3).
- 11: $\delta_t^s = \text{prox}_{\eta\lambda\phi_i}(\hat{x}_{\mathcal{B}_{s+1}}^t - \eta v_t^s) - \hat{x}_{\mathcal{B}_{s+1}}^t$.
- 12: $x_{\mathcal{B}_{s+1}}^{t+1} = x_{\mathcal{B}_{s+1}}^t + \delta_t^s$.
- 13: **end for**
- 14: $x_{\mathcal{B}_{s+1}} = x_{\mathcal{B}_{s+1}}^K, x_{\mathcal{B}_{s+1}}^0 = x_{\mathcal{B}_{s+1}}$.
- 15: **end for**

Proposed Method

Eliminating Step: We update \mathcal{B}_{s+1} for $\forall j \in \mathcal{B}_s$ as:

$$\Omega_j^D(A_j^\top y^s) + \Omega_j^D(A_j) \sqrt{2L(\mathcal{P}(x_{\mathcal{B}_s}^0) - D(y^s))} \geq n\lambda. \quad (2)$$

Variance-Reduced Sparse Gradient: Define Ψ_i as the set of blocks that intersect the nonzero coefficients of ∇f_i , let $n_{\mathcal{G}}$ be the number of occurrences that $\mathcal{G} \in \Psi_i$, if $n_{\mathcal{G}} > 0$, we define $d_{\mathcal{G}} = n/n_{\mathcal{G}}$. Thus, define diagonal matrix for block i as $[D_i]_{\mathcal{G},\mathcal{G}} = d_{\mathcal{G}} I_{|\mathcal{G}|}$, the gradient over \mathcal{B}_{s+1} can be computed as

$$v_t^s = \nabla f_i(a_{i,\mathcal{B}_{s+1}}^\top \hat{x}_{\mathcal{B}_{s+1}}^t) - \nabla f_i(a_{i,\mathcal{B}_{s+1}}^\top x_{\mathcal{B}_{s+1}}^0) + D_{i,\mathcal{B}_{s+1}} \nabla \mathcal{F}(x_{\mathcal{B}_{s+1}}^0). \quad (3)$$

Sparse Proximal Gradient Update: Define

$\phi_i(x) = \sum_{\mathcal{G} \in \Psi_i} d_{\mathcal{G}} \Omega_{\mathcal{G}}(x)$, the new proximal operator can be computed as

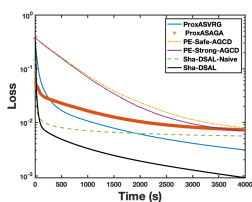
$$\text{prox}_{\eta\lambda\phi_i}(x') = \arg \min_x \frac{1}{2\eta} \|x - x'\|^2 + \lambda\phi_i(x). \quad (4)$$

Linear Convergence: Suppose $\tau \leq \frac{1}{10\sqrt{\Delta}}$, let step size $\eta = \min\{\frac{1}{24\kappa L}, \frac{\kappa}{2L}, \frac{\kappa}{10\tau L}\}$, inner loop size $K = \frac{4\log 3}{\eta\mu}$, we have

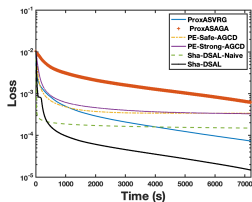
$$\mathbb{E} \|x_{\mathcal{B}_S} - x_{\mathcal{B}_S}^*\|^2 \leq (2/3)^S \|x_0 - x^*\|^2. \quad (5)$$

Elimination Ability: Equicorrelation set [4] is defined as $\mathcal{B}^* := \{j \in \{1, 2, \dots, q\} : \Omega_j^D(A_j^\top y^*) = n\lambda\}$. As DSAL converges, there exists an iteration number $S_0 \in \mathbb{N}$, s.t. $\forall s \geq S_0$, any variable block $j \notin \mathcal{B}^*$ is eliminated by DSAL almost surely.

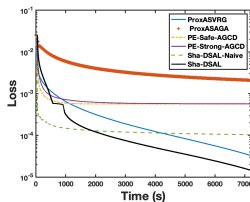
Convergence Results



(a) KDD2010



(b) Avazu-app

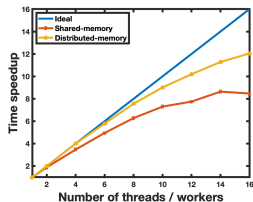


(c) Avazu-site

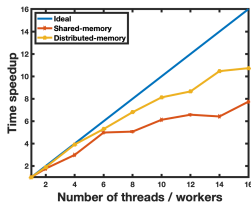
Figure: Convergence results on shared-memory architecture with 8 threads.

We compare six asynchronous methods: 1) PE-Strong-AGCD: parallel strong elimination in [1]; 2) PE-Safe-AGCD: parallel static safe elimination in [1]; 3) ProxASAGA [3]; 4) ProxASVRG [2]; 5) Sha-DSAL-Naive; 6) Our Sha-DSAL.

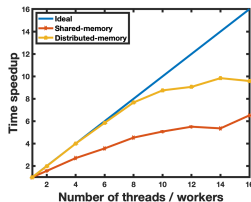
Linear Speedup Property



(a) KDD2010



(b) Avazu-app



(c) Avazu-site

Figure: Convergence results with different number of threads/workers.

Conclusion

- ▶ We propose a novel accelerated doubly sparse asynchronous learning method for stochastic composite optimization and apply it on shared-memory and distributed-memory architecture respectively.
- ▶ DSAL can simultaneously enjoy the model sparsity and data sparsity.
- ▶ We rigorously prove DSAL can achieve a linear convergence rate, reduce the per-iteration cost, and achieve a lower overall computational complexity under the strongly convex condition.
- ▶ We empirically show that DSAL can simultaneously achieve significant acceleration and linear speedup property.

Thank You!

-  Q. Li, S. Qiu, S. Ji, P. M. Thompson, J. Ye, and J. Wang.
Parallel lasso screening for big data optimization.
In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pages 1705–1714, 2016.
-  Q. Meng, W. Chen, J. Yu, T. Wang, Z.-M. Ma, and T.-Y. Liu.
Asynchronous stochastic proximal optimization algorithms with variance reduction.
In Proceedings of the AAAI Conference on Artificial Intelligence, volume 31, 2017.
-  F. Pedregosa, R. Leblond, and S. Lacoste-Julien.
Breaking the nonsmooth barrier: A scalable parallel method for composite optimization.
In NIPS, 2017.
-  R. J. Tibshirani et al.
The lasso problem and uniqueness.
Electronic Journal of statistics, 7:1456–1490, 2013.