

Adaptive Random Walk Gradient Descent for Decentralized Optimization

Tao Sun, Dongsheng Li, and Bao Wang

NUDT & U. of Utah

Table of contents

1. What is the problem (research area)
2. What we do
3. What about the new algorithm

What is the problem (research area)

The model

Given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with $\mathcal{V} = \{1, 2, \dots, n\}$ and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ being the set of agents and the set of edges that connect agents, respectively. We consider the **decentralized algorithm** for the minimization problem over graph \mathcal{G}

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}), \quad f_i(\mathbf{x}) = \mathbb{E}_{\xi \sim \mathcal{D}_i} F_i(\mathbf{x}; \xi), \quad (1)$$

where \mathcal{D}_i denotes the data distribution of the i -th client and $F_i(\mathbf{x}; \xi)$ is the loss function associated with the training data ξ .

The algorithm: random walk gradient descent

A token randomly walks over the graph \mathcal{G} to sample the data and updates the parameter.

Random walk gradient descent only involves one edge communication in each iteration, resulting in a minimum communication cost.

Another advantage is that it also applies to the directed graph setting.

What we do

The adaptive random walk

Motivated by the Adam-type algorithms, we propose the adaptive random walk algorithm as follows.

Algorithm 1 Adaptive Random Walk Gradient Descent

Require: parameters $\eta > 0, 0 \leq \theta < 1, \delta > 0$

Initialization: $\mathbf{g}^0 = \mathbf{0}, \mathbf{m}^0 = \mathbf{0}, \mathbf{v}^0 = \mathbf{0}$

for $k = 1, 2, \dots$

step 1: agent i_k calculates $\mathbf{g}^k = \nabla f_{i_k}(\mathbf{x}^k)$

step 2: $\mathbf{m}^k = \theta \mathbf{m}^{k-1} + (1 - \theta) \mathbf{g}^k$

step 3: $\mathbf{v}^k = \mathbf{v}^{k-1} + [\mathbf{g}^k]^2$

step 4: $\mathbf{z}^{k+1} = \mathbf{x}^k - \eta \mathbf{m}^k / (\mathbf{v}^k + \delta \mathbb{I})^{\frac{1}{2}}$

step 5: $\mathbf{x}^{k+1} = \arg \min_{\mathbf{x} \in \mathcal{K}} \|\mathbf{z}^{k+1} - \mathbf{x}\|_{(\mathbf{v}^k + \delta \mathbb{I})^{\frac{1}{2}}}^2$

step 6: uses random walk to choose a neighbor i_{k+1} and sends $(\mathbf{x}^k, \mathbf{m}^k, \mathbf{v}^k)$ via edge (i_k, i_{k+1}) to i_{k+1}

end for

The convergence

We investigate the adaptive random walk gradient descent and establish its theoretical performance bounds in both convex and nonconvex settings.

In the following, we present the convex one.

Theorem

Let Assumptions 1, 2, 3, 4, and condition

$$\mathbb{E}\|(\mathbf{v}^K + \delta\mathbb{I})^{\frac{1}{2}}\|_1 \leq CK^\alpha, \quad (2)$$

hold. Assume $\{\mathbf{x}^k\}_{k \geq 1}$ is generated by Algorithm 1. By setting $\eta = \min\left\{\frac{\ln(1/\sigma(\mathbf{P}))}{\ln(1/\epsilon)}, 1\right\}$, then

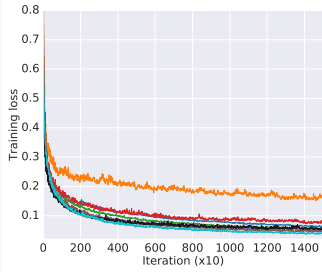
$$\mathbb{E}\left[f\left(\frac{\sum_{k=1}^K \mathbf{x}^k}{K}\right) - \min f\right] = \mathcal{O}(\epsilon), \quad (3)$$

with $K = \tilde{\mathcal{O}}\left(\max\left\{\frac{1}{\epsilon^{\frac{1}{1-\alpha}} [\ln(1/\sigma(\mathbf{P}))]^{\frac{1}{1-\alpha}}}, \frac{1}{\epsilon^{\frac{1}{1-\alpha}}}\right\}\right)$.

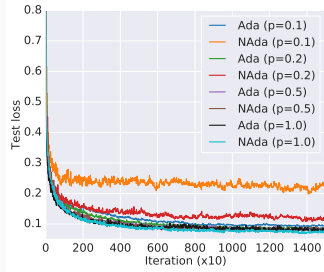
Due to the boundedness of stochastic gradients $\{\mathbf{g}^k\}_{k \geq 0}$, $\alpha = 1/2$ can hold in (2) without any extra assumption.

As the stochastic gradient decay fast (i.e., $\alpha < 1/2$), adaptive random walk gradient descent can be faster than the non-adaptive ones.

What about the new algorithm



Training



Test

Figure 1: Comparison of adaptive and non-adaptive random walk algorithms, with both sparse ($p < 1$) and non-sparse gradients ($p = 1$), for training an MLP model for MNIST classification. In this experiment, we have ten clients connected by a ring graph.

References i

1. Duchi, J. C., Agarwal, A., Johansson, M., and Jordan, M. I. Ergodic mirror descent. *SIAM Journal on Optimization*, 22(4):1549–1578, 2012.
2. Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.
3. Mao, X., Yuan, K., Hu, Y., Gu, Y., Sayed, A. H., and Yin, W. Walkman: A communication-efficient random-walk algorithm for decentralized optimization. *IEEE Transactions on Signal Processing*, 68:2513–2528, 2020.
4. Sun, T., Sun, Y., and Yin, W. On markov chain gradient descent. *Advances in neural information processing systems*, 2018.
5. Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *ICLR (Poster)*, 2015.