MARINA: Faster Non-Convex Distributed Learning with Compression

Eduard Gorbunov^{1,2,3,}, Konstantin Burlachenko³, Zhize Li³, Peter Richtárik³ ¹MIPT (Russia), ²Yandex (Russia), ³KAUST (Saudi Arabia)

1. The Problem



Problem: Non-convex distributed optimization / training, where *n* workers (devices/clients) jointly solve a problem by communicating with a central server

Assumptions: smoothness of local loss functions and lower-boundedness

$$\forall x, y \in \mathbb{R}^d \quad \|\nabla f_i(x) - \nabla f_i(y)\| \le L_i \|x - y\|$$
$$\exists f_* \in \mathbb{R} : \ \forall x \in \mathbb{R}^d \quad f(x) \ge f_*$$

Distributed methods often suffer from **communication bottleneck**

One can handle this issue via communication compression

2. Unbiased Compression

 $x \to \mathcal{Q}(x) \quad \mathbb{E}[\mathcal{Q}(x)] = x \quad \mathbb{E}\|\mathcal{Q}(x) - x\|^2 \le \omega \|x\|^2$

Example: RandK compression operator picks K components uniformly at random

$$\begin{array}{c} d=5\\ K=2 \end{array} \quad x = \begin{pmatrix} 1\\ -15\\ 0.2\\ -7\\ 10 \end{pmatrix} \quad \longrightarrow \quad \mathcal{Q}(x) = \frac{5}{2} \cdot \begin{pmatrix} 1\\ 0\\ 0\\ -7\\ 0 \end{pmatrix} \qquad \omega = \frac{d}{K} - 1 \end{array}$$

References: [1] Alistarh, Dan, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. "QSGD: Communication-efficient SGD via gradient quantization and encoding." In Advances in Neural Information Processing Systems, pp. 1709-1720, (2017). [2] Mishchenko, Konstantin, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. "Distributed learning with compressed gradient differences." arXiv:1901.09269 (2019). [3] Horváth, Samuel, Dmitry Kovalev, Konstantin Mishchenko, Sebastian Stich, and Peter Richtárik. "Stochastic distributed learning with gradient quantization and variance reduction." arXiv:1904.05115 (2019).





3. QGD and DIANA

learnable local shifts: $h_i^{k+1} = h_i^k + \alpha \mathcal{Q} \left(\nabla f_i(x^k) - h_i^k \right)$

of communication rounds to find \hat{x} such that $\mathbb{E}\left| \| \nabla f(\hat{x}) \|^2 \right| \leq \varepsilon^2$

QGD:

DIANA: $\mathcal{O}\left(\frac{1+(1+\omega)\sqrt{\omega/n}}{2}\right)$

 \geq

0.0 0.2

dependencies on numerical factors, smoothness constants and initial suboptimality are omitted

previous SOTA complexity

4. New Method: MARINA

$$g_i^k = \begin{cases} \nabla f_i(x^k) & \text{w.p. } p \\ g^{k-1} + \mathcal{Q}\left(\nabla f_i(x^k) - \nabla f_i(x^{k-1})\right) & \text{w.p. } 1 - p \end{cases}$$

of communication rounds to find an \mathcal{E} -stationary point



7. In the paper, we also have

of Epochs

(1) Mini-batch VR-MARINA, (2) VR-MARINA for expectation minimization, (3) MARINA with partial partcipation of clients, (4) analysis under Polyak-Lojasiewicz condition, (5) explicit dependencies on smoothness constants, non-uniform sampling, (6) experiments with neural networks

1e7

DIANA (K=10)

0.4 0.6 0.8 1.0 1.2

#bits/n

