

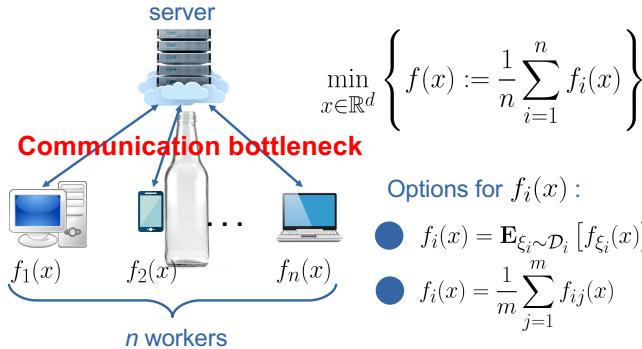
Linearly Converging Error Compensated SGD

Eduard Gorbunov^{1,2,3,4}, Dmitry Kovalev⁴, Dmitry Makarenko¹, Peter Richtárik⁴

¹MIPT (Russia), ²Yandex (Russia), ³Sirius (Russia), ⁴KAUST (Saudi Arabia)



1. The Problem



Problem: Distributed optimization / training, where n workers (devices / clients) jointly solve a problem by communicating with a central server.

Assumptions: Smoothness and convexity of local loss functions and quasi-strong convexity of their average:

$$\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\| \quad f_i(x) - f_i(y) \geq \langle \nabla f_i(y), x - y \rangle$$

$$f(x^*) \geq f(x) + \langle \nabla f(x), x^* - x \rangle + \frac{\mu}{2} \|x^* - x\|^2$$

The key bottleneck comes from the high cost of communication. However, one can handle this using compression of messages.

- ✓ There is good theory for methods with unbiased compressors: there exist linearly converging SGD-like methods [1, 2].
- ✗ However, Parallel-SGD with biased compressors diverges exponentially fast on some problems [3].
- ✓ One can fix this using error compensation (EC-SGD) [4]
- ✗ but EC-SGD fails to converge linearly even when workers compute full gradients.

2. Our Contributions

- ✓ The first linearly converging error compensated SGD method
- ✓ General theoretical framework covering error compensated methods, methods with delayed updates and non-distributed methods
- ✓ In one theorem, we recover the sharpest rates for all known special cases
- ✓ 16 new methods fitting our framework: 8 error-compensated methods, 7 methods with delayed updates and DIANA with bi-directional compression

3. Compression Operators

Unbiased compressors

$$x \rightarrow \mathcal{Q}(x) \quad \mathbb{E}[\mathcal{Q}(x)] = x$$

$$\mathbb{E}\|\mathcal{Q}(x) - x\|^2 \leq \omega\|x\|^2$$

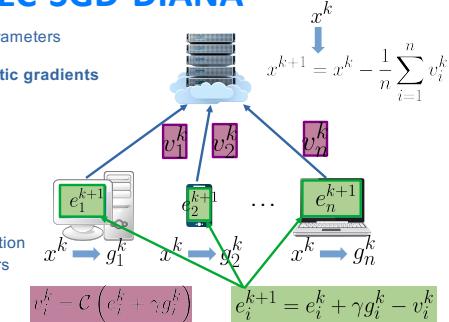
Biased compressors

$$x \rightarrow \mathcal{C}(x)$$

$$\mathbb{E}\|\mathcal{C}(x) - x\|^2 \leq (1 - \delta)\|x\|^2$$

4. EC-SGD-DIANA

- 1 Server broadcasts new parameters
- 2 Workers compute stochastic gradients in parallel
- 3 Compression
- 4 Devices send compressed vectors and update unsent information
- 5 Server gathers the information and updates the parameters
- 6 Repeat steps 1 – 5



The main innovation is reflected in the definition of g_i^k :

$$g_i^k = \hat{g}_i^k - h_i^k + h^k \quad \mathbb{E}[\hat{g}_i^k | x^k] = \nabla f_i(x^k)$$

"Learned shift vectors" are the key to get linear convergence:

$$h_i^{k+1} = h_i^k + \alpha \mathcal{Q}(\hat{g}_i^k - h_i^k) \quad h^k = \frac{1}{n} \sum_{i=1}^n h_i^k$$

5. General Framework

The assumption below covers a very broad class of methods:

$$g^k = \frac{1}{n} \sum_{i=1}^n g_i^k, \quad \mathbb{E}[g^k | x^k] = \nabla f(x^k) \quad \bar{g}_i^k = \mathbb{E}[g_i^k | x^k]$$

$$\frac{1}{n} \sum_{i=1}^n \|\bar{g}_i^k\|^2 \leq 2A(f(x^k) - f(x^*)) + B_1 \sigma_{1,k}^2 + B_2 \sigma_{2,k}^2 + D_1$$

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}[\|g_i^k - \bar{g}_i^k\|^2 | x^k] \leq 2\tilde{A}(f(x^k) - f(x^*)) + \tilde{B}_1 \sigma_{1,k}^2 + \tilde{B}_2 \sigma_{2,k}^2 + \tilde{D}_1$$

$$\mathbb{E}[|g^k|^2 | x^k] \leq 2A'(f(x^k) - f(x^*)) + B'_1 \sigma_{1,k}^2 + B'_2 \sigma_{2,k}^2 + D'_1$$

$$\mathbb{E}[\sigma_{1,k+1}^2 | \sigma_{1,k}^2, \sigma_{2,k}^2] \leq (1 - \rho_1) \sigma_{1,k}^2 + 2C_1(f(x^k) - f(x^*)) + G\rho_1 \sigma_{2,k}^2 + D_2$$

$$\mathbb{E}[\sigma_{2,k+1}^2 | \sigma_{2,k}^2] \leq (1 - \rho_2) \sigma_{2,k}^2 + 2C_2(f(x^k) - f(x^*))$$

- Reflects smoothness properties of the problem and noises introduced by compressions and stochastic gradients
- Describes the process of variance reduction of the variance coming from compressions
- Describes the process of variance reduction of the variance coming from stochastic gradients

Table 1: Complexity of Error-Compensated SGD methods established in this paper. Symbols: ε = error tolerance; δ = contraction factor of compressor \mathcal{C} ; ω = variance parameter of compressor \mathcal{Q} ; $\kappa = \mathcal{L}/\mu$; \mathcal{L} = expected smoothness constant; σ_*^2 = variance of the stochastic gradients in the solution; $\zeta_*^2 = \text{average of } \|\nabla f_i(x^*)\|^2$; σ^2 = average of the uniform bounds for the variances of stochastic gradients of workers. EC-GDstar, EC-LSVRGstar and

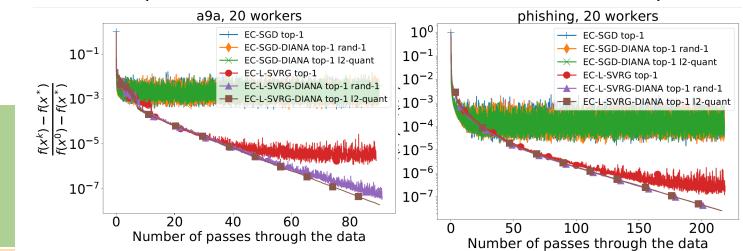
Probl.	Method	Citation	Rate (constants ignored)
Σ	EC-SGDsr	new	$\tilde{\mathcal{O}}\left(\frac{\mathcal{L}}{\mu} + \frac{L + \sqrt{\delta L\mathcal{L}}}{\delta\mu} + \frac{\sigma_*^2}{n\mu\varepsilon} + \frac{\sqrt{L(\sigma_*^2 + \zeta_*^2/\delta)}}{\mu\sqrt{\delta\varepsilon}}\right)$
\mathbf{E}	EC-SGD	[4]	$\tilde{\mathcal{O}}\left(\frac{\kappa}{\delta} + \frac{\sigma_*^2}{\delta\varepsilon} + \frac{\sqrt{L(\sigma_*^2 + \zeta_*^2/\delta)}}{\delta\mu\sqrt{\varepsilon}}\right)$
Σ	EC-GDstar	new	$\mathcal{O}\left(\frac{\kappa}{\delta} \log \frac{1}{\varepsilon}\right)$
\mathbf{E}	EC-SGD-DIANA	new	Opt. I: $\tilde{\mathcal{O}}\left(\omega + \frac{\kappa}{\delta} + \frac{\sigma_*^2}{n\mu\varepsilon} + \frac{\sqrt{L(\sigma_*^2 + \zeta_*^2/\delta)}}{\delta\mu\sqrt{\varepsilon}}\right)$ Opt. II: $\tilde{\mathcal{O}}\left(\frac{1+\omega}{\delta} + \frac{\kappa}{\delta} + \frac{\sigma_*^2}{n\mu\varepsilon} + \frac{\sqrt{L(\sigma_*^2 + \zeta_*^2/\delta)}}{\mu\sqrt{\delta\varepsilon}}\right)$
Σ	EC-SGDsr-DIANA	new	Opt. I: $\tilde{\mathcal{O}}\left(\omega + \frac{\kappa}{\mu} + \frac{\sqrt{L\mathcal{L}}}{\delta\mu} + \frac{\sigma_*^2}{n\mu\varepsilon} + \frac{\sqrt{L(\sigma_*^2 + \zeta_*^2/\delta)}}{\delta\mu\sqrt{\varepsilon}}\right)$ Opt. II: $\tilde{\mathcal{O}}\left(\frac{1+\omega}{\delta} + \frac{\kappa}{\mu} + \frac{\sqrt{L\mathcal{L}}}{\delta\mu} + \frac{\sigma_*^2}{n\mu\varepsilon} + \frac{\sqrt{L(\sigma_*^2 + \zeta_*^2/\delta)}}{\mu\sqrt{\delta\varepsilon}}\right)$
\mathbf{E}	EC-GD-DIANA [†]	new	$\mathcal{O}\left((\omega + \frac{\kappa}{\delta}) \log \frac{1}{\varepsilon}\right)$
Σ	EC-LSVRG	new	$\tilde{\mathcal{O}}\left(m + \frac{\kappa}{\delta} + \frac{\sqrt{L\mathcal{L}^2}}{\delta\mu\varepsilon}\right)$
Σ	EC-LSVRGstar	new	$\tilde{\mathcal{O}}\left((m + \frac{\kappa}{\delta}) \log \frac{1}{\varepsilon}\right)$
Σ	EC-LSVRG-DIANA	new	$\mathcal{O}\left((m + \omega + \frac{\kappa}{\delta}) \log \frac{1}{\varepsilon}\right)$

See more examples of methods (with delayed updates and without error feedback) fitting our framework together with the rates in weakly convex case in our paper.

6. Numerical Experiments

We conducted several numerical experiments on logistic regression problem with L_2 -regularization:

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) = \frac{1}{N} \sum_{i=1}^N \log(1 + \exp(-y_i \cdot (Ax)_i)) + \frac{\mu}{2} \|x\|^2 \right\}$$



References

- [1] Mishchenko, Konstantin, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. "Distributed learning with compressed gradient differences." arXiv preprint arXiv:1901.09269 (2019).
- [2] Horváth, Samuel, Dmitry Kovalev, Konstantin Mishchenko, Sebastian Stich, and Peter Richtárik. "Stochastic distributed learning with gradient quantization and variance reduction." arXiv preprint arXiv:1904.05115 (2019).
- [3] Beznosikov, Aleksandr, Samuel Horváth, Peter Richtárik, and Mher Safaryan. "On biased compression for distributed learning." arXiv preprint arXiv:2002.12410 (2020).
- [4] Stich, Sebastian U., and Sai Praneeth Karimireddy. "The error-feedback framework: Better rates for SGD with delayed gradients and compressed communication." arXiv preprint arXiv:1909.05350 (2019).