

Task

- We are interested in minimizing a loss function $f(x)$, $x \in \mathbb{R}^d$ defined as

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x).$$

The minimizer of this function is written as $x^* = \arg \min_x f(x)$.

Stochastic Gradient Descent (SGD)

- Gradient descent used to tackle this problem uses the following update rule:

$$x^{t+1} = x^t - \eta^t f'(x^t), \quad (1)$$

where η^t is a chosen step-size for iteration t .

- SGD picks a random i and uses a stochastic update:

$$x^{t+1} = x^t - \eta^t f'_i(x^t). \quad (2)$$

Convergence (in expectation) is guaranteed if $\mathbf{E}_i[f'_i(x^t)] = f'(x^t)$ but SGD suffers from a **high variance**.

Variance reduction

- SAGA [1], SVRG [2], S2GD [3],... belong to a family of generalized SGD algorithms that exhibit lower variance and use the following update:

$$x^{t+1} = x^t - \eta v^t, \quad \text{where } \mathbf{E}[v^t] = f'(x^t) \quad (3)$$

- Convergence analysis considers

$$\begin{aligned} \mathbf{E}\|x^{t+1} - x^*\|^2 &= \mathbf{E}\|x^t - \eta v^t - x^*\|^2 \\ &= \|x^t - x^*\|^2 - 2\eta \langle x^t - x^*, f'(x^t) \rangle + \eta^2 \underbrace{\mathbf{E}\|v^t\|^2}_{\text{Variance}} \end{aligned} \quad (4)$$

The (negative) middle term is what guarantees progress for any gradient descent procedure on a (strongly) convex objective.

Bound for SAGA-style updates

- Introduce a **correction term** for f'_i denoted $g_i(\phi^t)$, so that $g(\phi^t) := \frac{1}{n} \sum_i g_i(\phi^t)$ denotes its expectation. We write

$$v^t := f'_i(x^t) - g_i(\phi^t) + g(\phi^t). \quad (5)$$

- We use the shorthand notation $\delta h(x) := h(x) - h(x^*)$. The variance term can be decomposed and bounded as follows:

$$\begin{aligned} \mathbf{E}\|v^t\|^2 &= \mathbf{E}\|v^t\|^2 + \mathbf{E}\|v^t - \mathbf{E}[v^t]\|^2 = \mathbf{E}\|v^t - f'(x^t)\|^2 + \|f'(x^t)\|^2 \\ &\leq (1 + \beta) \mathbf{E}\|\delta f'_i(x^t)\|^2 + (1 + \beta^{-1}) \mathbf{E}\|\delta g_i(\phi^t)\|^2 \\ &\quad - \beta \|f'(x^t)\|^2 - (1 + \beta^{-1}) \|g(\phi^t)\|^2, \end{aligned} \quad (6)$$

where $\beta > 0$.

- The variance term vanishes as we convergence to the optimum.**

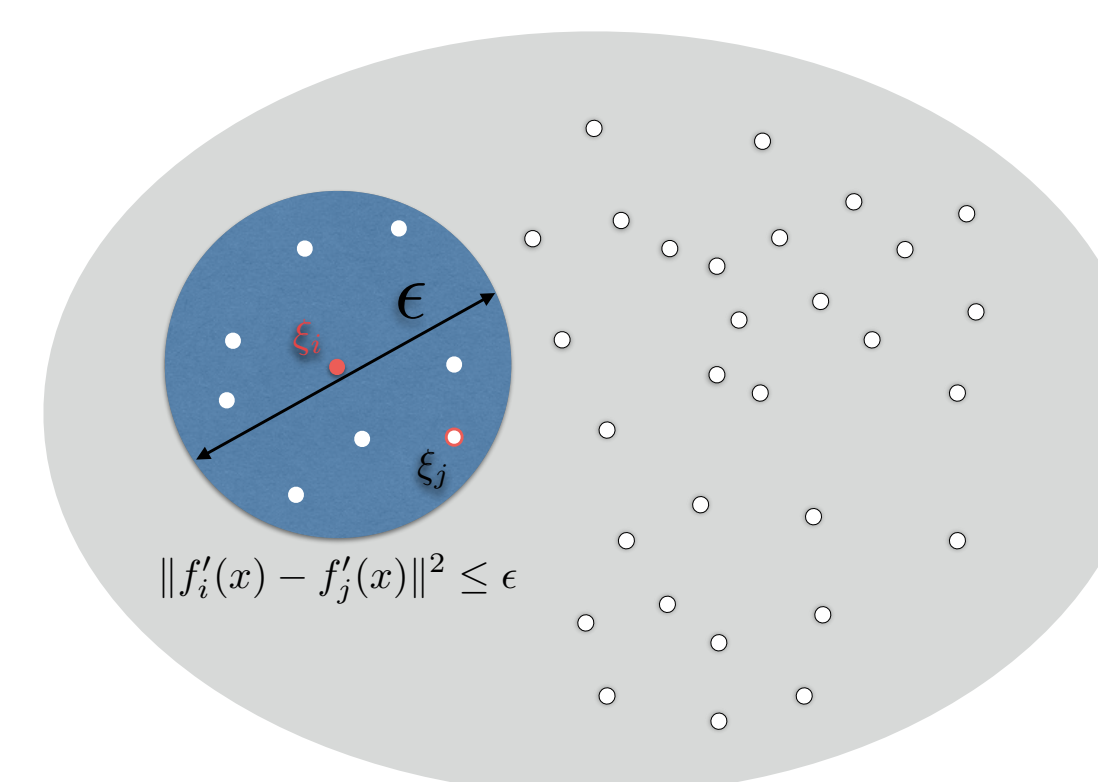
Generalizing SAGA updates

- Propose to generalize the correction to a weighted (convex) sum of the following type (exploiting some clustering structure in how to chose τ):

$$g_i(\phi^t) = \sum_j \tau_{ij} f'_j(\phi_j^t), \quad \text{with } \sum_j \tau_{ij} = 1 (\forall i), \quad \tau_{ij} \geq 0 (\forall i, j), \quad (7)$$

where $\tau_{ij} = \delta_{ij}$ is the special case of the original SAGA algorithm

- In our algorithm, we use $\tau_{ij} \in \{0, 1\}$, i.e. **selection of one "neighbor" j** whose ℓ_2 -distance from i is less than ϵ



Algorithm

- INPUTS :**
- \mathcal{D} : Training set of n examples.
- η : Step size
- ϵ : Neighborhood size
- OUTPUT :** x^T
- Cluster datapoints in \mathcal{D} whose distance is less than ϵ
- for** $t = 1 \dots T$ **do**
- Randomly pick $i \in 1 \dots n$
- $x^t = x^{t-1} - \eta(\nabla f_i(x^{t-1}) - (g_i(\phi^t) - \mathbf{E}[g_i(\phi^t)]))$
- end for**

Convergence properties

- Pivot the analysis around $f'_j(x^*)$ instead of $f'_i(x^*)$.
- The variance converges to $2(1 + \beta)(\mathbf{E}\|f'_i(x^*) - f'_j(x^*)\|^2)$ as $x_t \rightarrow x^*$
- Proof convergence sketch:

– Define the following Lyapunov function:

$$T^t := \|x^t - x^*\|^2 + \alpha [\delta f(\phi^t) - \langle f'_i(x^*), \phi^t - x^* \rangle] \quad (8)$$

– Show that

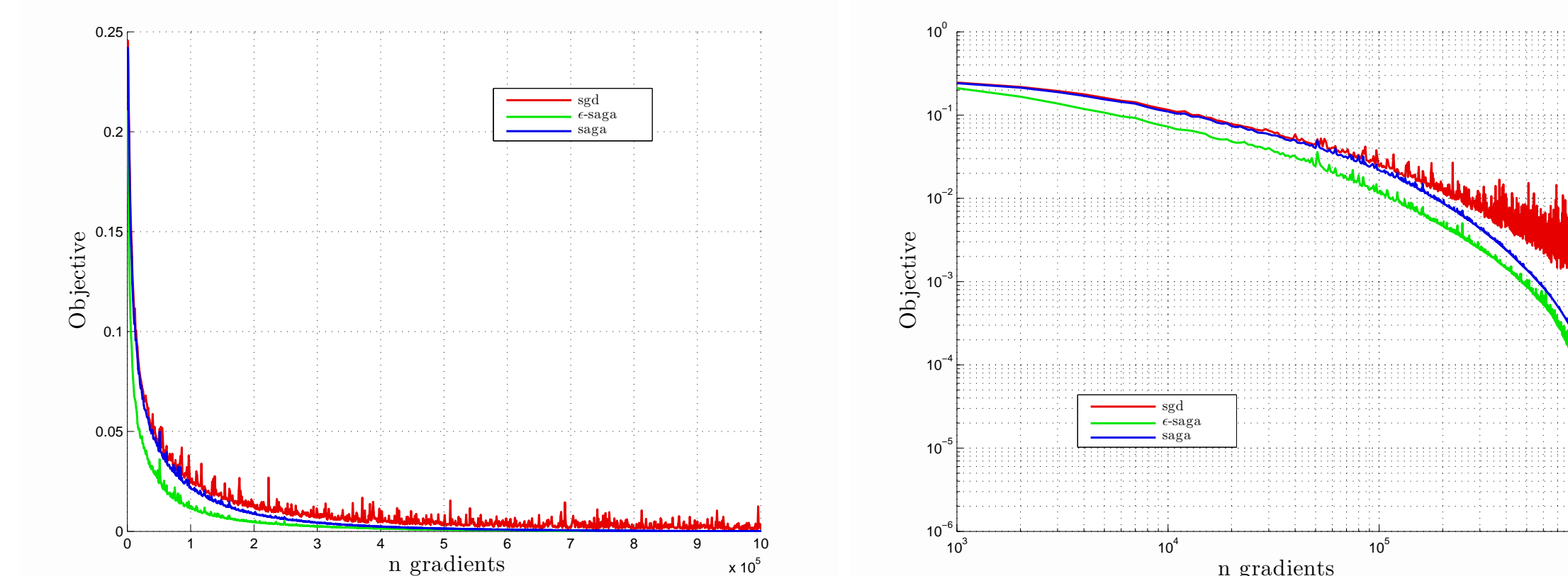
$$\mathbf{E}\|x^{t+1} - x^*\|^2 \leq \rho^t T^0 + O(\eta^2 \epsilon), \quad \rho < 1$$

- In general, there is a **constant non vanishing error** of $O(\eta^2 \epsilon)$
- Open question: How can we get $\epsilon \rightarrow 0$?

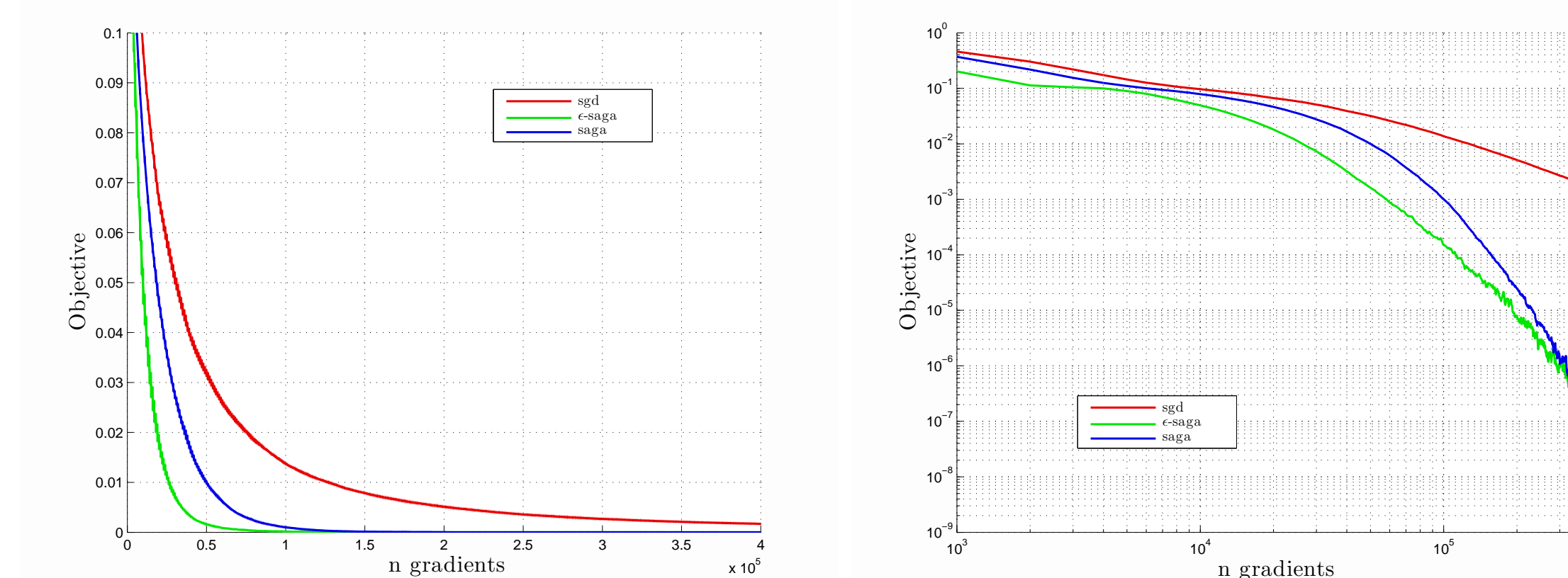
Results

- For SGD, we used $\eta_t = \frac{\eta_0 T_0}{T_0 + t}$
- Both SAGA and ϵ -SAGA (our method) use a constant step-size
- We added an ℓ_2 -regularizer with parameter $\lambda = 10^{-4}$

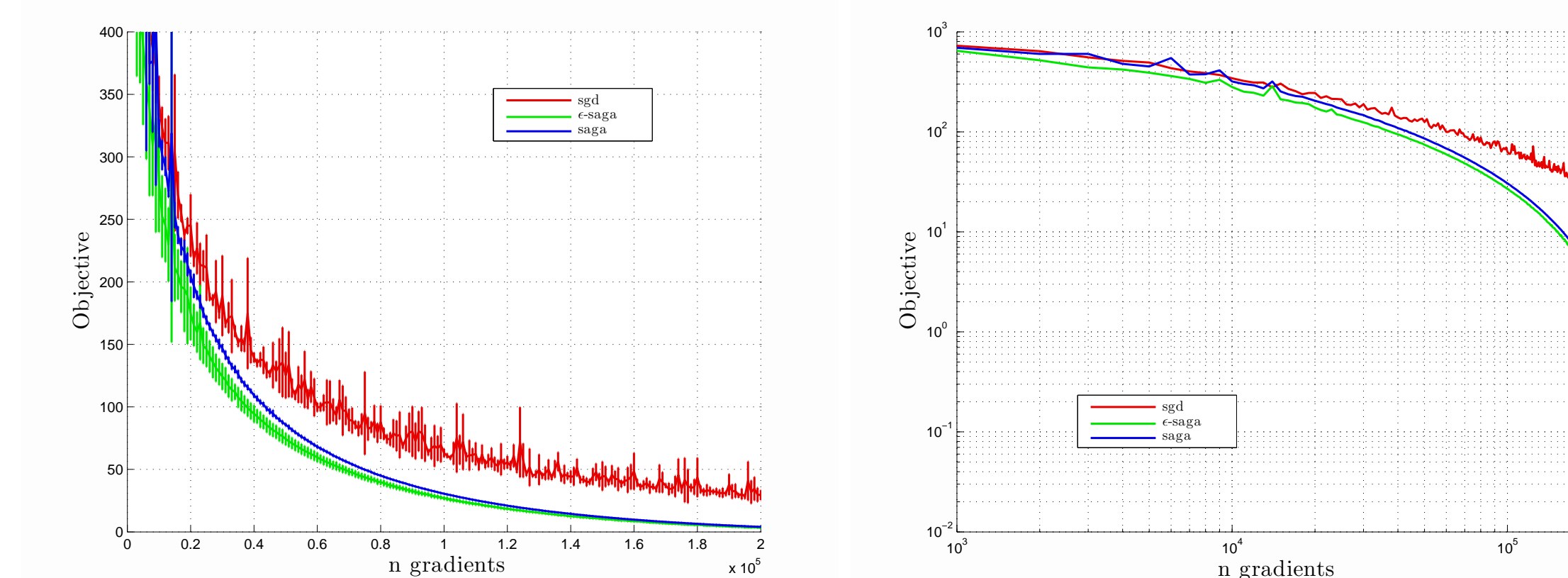
COV dataset (logistic regression)



IJCN dataset (logistic regression)



Year dataset (least-squares regression)



Future work

- Strengthen theoretical guarantees
- Study conditions required for $\epsilon \rightarrow 0$
- Generalization bound

References

- A. Defazio, F. Bach, and S. Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems*, pages 1646–1654, 2014.
- R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems*, pages 315–323, 2013.
- J. Konečný and P. Richtárik. Semi-stochastic gradient descent methods. *arXiv preprint arXiv:1312.1666*, 2013.