# Achieving the Optimal Convergence Rate in Stochastic Optimisation

Derek Driggs
CCIMI

Advisor: Dr Carola Schönlieb
Collaborators: Dr Jingwei Liang, Dr Clarice Poon

Many machine learning problems can be formulated as

$$\min_x \quad \frac{1}{n}\sum_{i=1}^{n} f_i(x) + g(x)$$

Convex, smooth functions with $L$-Lipschitz continuous gradients

Convex, non-smooth function

$$\min_x \quad \boxed{\frac{1}{n}\sum_{i=1}^{n} f_i(x)} + \boxed{g(x)}$$

$\|\mathcal{A}(\cdot)\|_2 - \ell_2$-norm with linear operator $\mathcal{A}$

$\|\cdot\|_1 - \ell_1$-norm
$\|\cdot\|_* - $ nuclear norm

Applications: LASSO, Robust PCA, Logistic Regression

$$\min_x \quad \boxed{\frac{1}{n}\sum_{i=1}^{n} f_i(x)} + \boxed{g(x)}$$

This problem can be solved using *proximal gradient descent*:

$$x_{k+1} = \mathsf{prox}_{\gamma g}\left(x_k - \gamma \nabla f(x_k)\right)$$

where the *proximal operator* is defined as

$$\mathsf{prox}_h(y) := \mathsf{argmin}_x \quad \frac{1}{2}\|x - y\|^2 + h(x)$$

This algorithm requires the evaluation of $n$ gradients per iteration.

Computing the full gradient is expensive for large $n$, so we can replace $\nabla f(x_k)$ with an estimate of the gradient, $\widetilde{\nabla} f(x_k)$, where

$$\widetilde{\nabla}_{\text{SGD}} f(x_k) = \nabla f_j(x_k), \hspace{2cm} \text{(SGD)}$$

$$\widetilde{\nabla}_{\text{SAGA}} f(x_k) = \nabla f_j(x_k) - \nabla f_j(\varphi_k^j) + \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(\varphi_k^i), \hspace{1cm} \text{(SAGA)}$$

$$\widetilde{\nabla}_{\text{SVRG}} f(x_k) = \nabla f_j(x_k) - \nabla f_j(\tilde{x}) + \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(\tilde{x}). \hspace{1cm} \text{(SVRG)}$$

The index $j$ is chosen uniformly at random.

$\varphi_k^j$ — The gradient $\nabla f_j(\varphi_k^j)$ is stored for future iterates.

$\tilde{x}$ — The full gradient $\frac{1}{n} \sum_{i=1}^{n} \nabla f(\tilde{x})$ is computed every $2n$ iterations and stored for future iterates.

With $x^*$ the minimiser of $F(x) := \frac{1}{n}\sum_{i=1}^{n} f_i(x) + g(x)$, the suboptimality at iteration $k$ is $F(x_k) - F(x^*)$.

For proximal gradient descent on convex objectives,

$$F(x_k) - F(x^*) \leq \mathcal{O}\left(\frac{1}{k}\right)$$

For SVRG and SAGA,

$$\mathbb{E}\left[F(x_k) - F(x^*)\right] \leq \mathcal{O}\left(\frac{1}{k}\right)$$

Nesterov's momentum trick is a slight modification that offers enormous acceleration:

$$y_{k+1} = x_k + \alpha(x_k - x_{k-1})$$
$$x_{k+1} = \mathsf{prox}_{\gamma g}\left(y_{k+1} - \gamma \nabla f(y_{k+1})\right)$$

With $\alpha$ and $\gamma$ chosen appropriately,

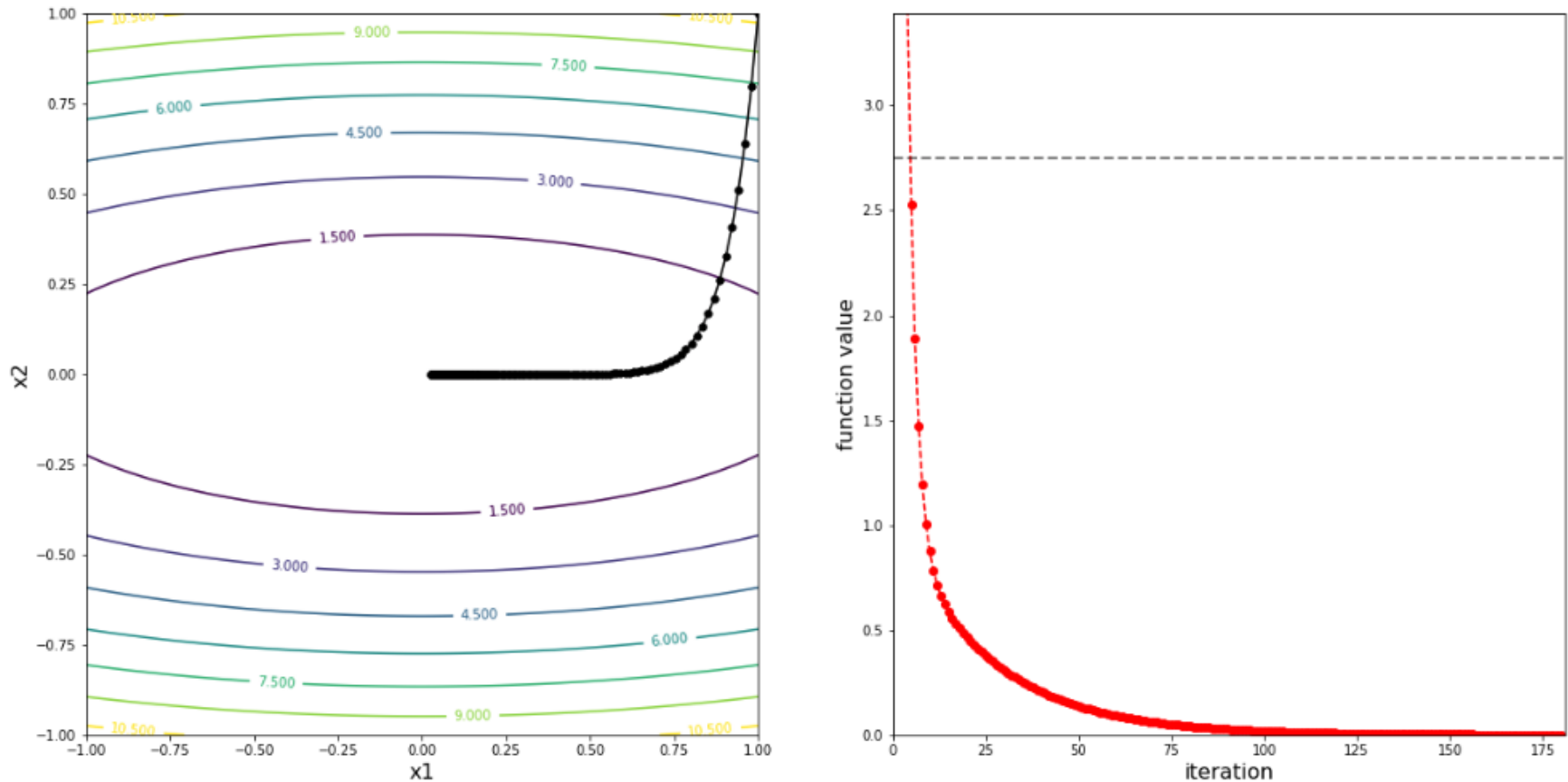$$F(x_k) - F(x^*) \leq \mathcal{O}\left(\frac{1}{k^2}\right),$$

and this rate is optimal.

**Figure:** Using gradient descent to minimize $x_1^2 + 10x_2^2$ without momentum.

**Figure:** Using gradient descent to minimize $x_1^2 + 10x_2^2$ with momentum.

In the stochastic setting, momentum propagates "bad" gradient evaluations, so it is unclear whether momentum-based methods improve performance.

"Katyusha" (Allen-Zhu, 2017) achieves the optimal $\mathcal{O}\left(\frac{1}{k^2}\right)$-rate using **negative momentum** and **linear coupling**:

$$x_{k+1} = \alpha_1 z_k + \boxed{\alpha_2 \tilde{x}} + (1 - \alpha_1 - \alpha_2)y_k$$
$$y_{k+1} = \mathsf{prox}\left(x_{k+1} - \gamma\alpha_1 \widetilde{\nabla}_{\mathsf{SVRG}} f(x_{k+1})\right)$$
$$z_{k+1} = \mathsf{prox}\left(z_k - \gamma \widetilde{\nabla}_{\mathsf{SVRG}} f(x_{k+1})\right)$$

The term $\alpha_2 \tilde{x}$ "attracts" $x_{k+1}$, supposedly limiting the effects of detrimental gradient evaluations.

Using linear-coupling analysis, we prove the following:

**Theorem**
**Consider the algorithm**

$$x_{k+1} = \alpha z_k + (1-\alpha)y_k$$
$$y_{k+1} = \mathsf{prox}\left(x_{k+1} - \gamma\alpha\widetilde{\nabla}f(x_{k+1})\right)$$
$$z_{k+1} = \mathsf{prox}\left(z_k - \gamma\widetilde{\nabla}f(x_{k+1})\right)$$

**with $\widetilde{\nabla} = \widetilde{\nabla}_{\mathsf{SAGA}}$ or $\widetilde{\nabla}_{\mathsf{SVRG}}$. Choosing $\alpha$ and $\gamma$ appropriately,**

$$F(x_k) - F(x^*) \leq \mathcal{O}\left(\frac{1}{k^2}\right)$$

Our analysis can be easily extended to many other (variance-reduced) stochastic gradient estimators.
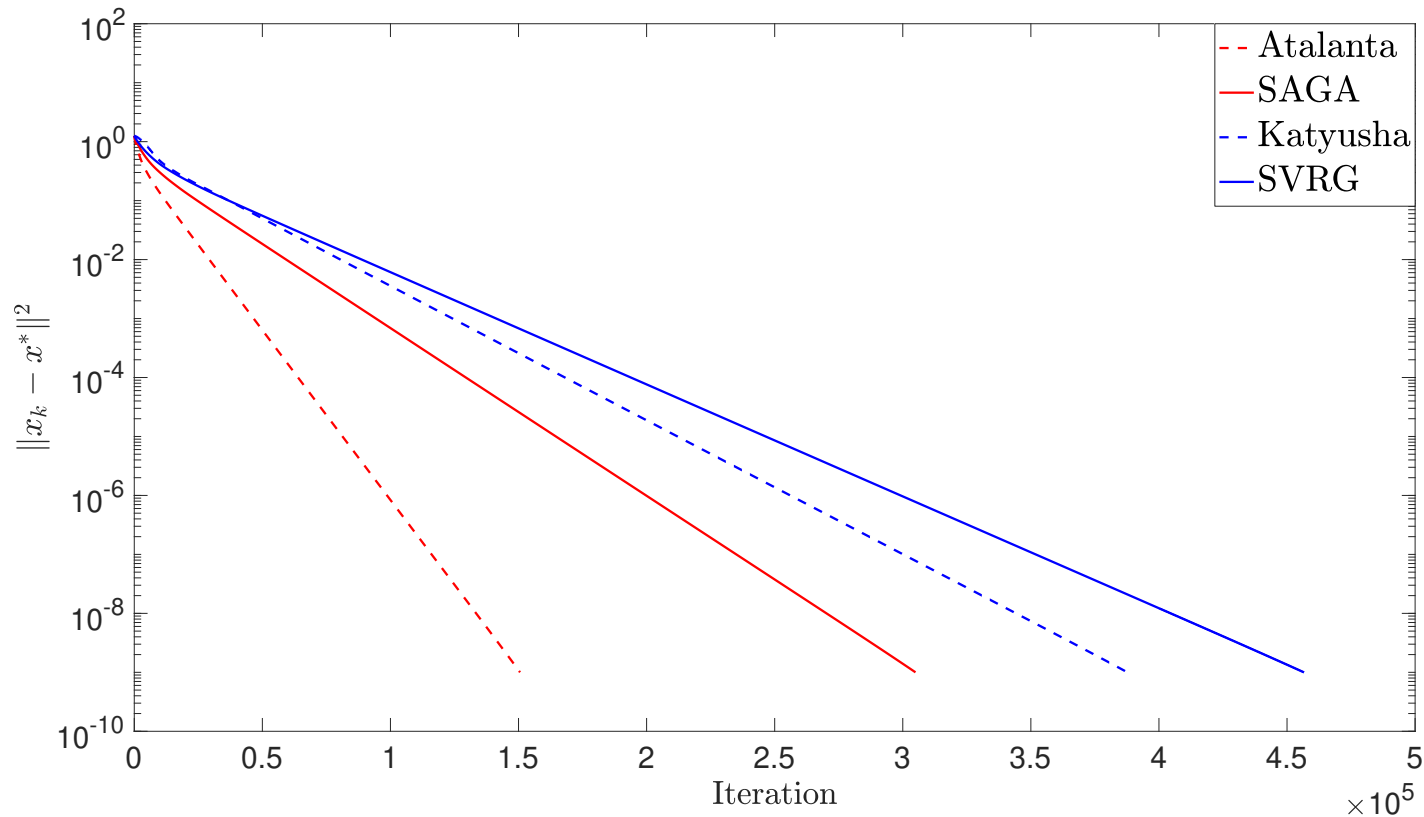
**Figure:** Comparing stochastic gradient methods on a sparse logistic regression problem (see poster for details).

"Atalanta" uses our acceleration framework with the SAGA gradient estimate, and is extremely fast.