# A Generic Acceleration Framework for Stochastic Composite Optimization

# Andrei Kulunchakov and Julien Mairal - Inria Grenoble



#### Goal

**Problem: acceleration** of convex optimization methods for minimization of the **expected risk**.

Developed tool: generic framework using iterative construction and minimization of surrogate functions.

#### Motivation and context

- Common approach in ML: approximate the expected risk by a finite sum. But this creates an additional bias of approximation.
- If infinite amount of data is available, one can minimize the **expected risk**, requiring stochastic optimization. Standard fast methods may be unstable or slow in this case.
- Having access only to approximate gradients, we substantially accelerate/stabilize these methods by the framework.

## Detailed statement of the problem

$$\min_{\mathbf{x}\in\mathbb{R}^p}\left\{F(\mathbf{x})=\frac{1}{n}\sum_{i=1}^n f_i(\mathbf{x})+\psi(\mathbf{x})\right\}\quad\text{with stochastic terms }f_i(\mathbf{x}),$$

which are  $\mu$ -strongly convex and L-smooth and  $\psi$  is convex. When n=1, we recover classical composite optimization problems. For each i, the gradient is noisy  $\tilde{\nabla} f_i(x) = \nabla f_i(x) + \xi_i$ .

#### Contributions in examples

There are mild requirements on the convergence of methods that we want to accelerate: they even can be biased!

Examples of acceleration are presented below with  $\Delta_0 = F(x_0) - F^*$ ,  $\Omega^2 = ||x_0 - x^*||^2$  and a targeted accuracy  $\varepsilon$ .

#### **Deterministic cases:**

$$(L/\mu)\log\left(\Delta_0/\varepsilon\right)$$
 becomes  $\sqrt{L/\mu}\log\left(\Delta_0/\varepsilon\right)$ ,  $(n+L/\mu)\log\left(\Delta_0/\varepsilon\right)$  becomes  $\left(n+\sqrt{nL/\mu}\right)\log\left(\Delta_0/\varepsilon\right)$ ;

#### and stochastic cases:

$$(L/\mu)\log(\Delta_0/\varepsilon)$$
 biased as  $\varepsilon \geq \sigma^2/\mu$  becomes convergent  $\sqrt{L/\mu}\log(\Delta_0/\varepsilon) + \sigma^2/\mu\varepsilon$ , (1)

#### and similarly biased

$$(n+L/\mu)\log(\Delta_0/\varepsilon)$$
 becomes convergent (2) 
$$\left(n+\sqrt{nL/\mu}\right)\log(\Delta_0/\varepsilon)+\sigma^2/\mu\varepsilon.$$

#### Original (deterministic) approach called Catalyst

- Given  $\kappa > 0$ , for each k we build a surrogate function  $h_k(x) \triangleq F(x) + (\kappa/2) \|x y_{k-1}\|^2$ .
- Enjoying better properties than F(x), it is effectively minimized.
- Iterative minimization of  $h_1, h_2, ..., h_k$  by some method  $\mathcal{M}$ , allows to accelerate  $\mathcal{M}$  on minimization of  $F(\mathbf{x})$ .

But, this is only for deterministic cases.

#### Our approach. Stochastic surrogates

We allow greater flexibility to surrogate functions  $h_k$ :

- $h_k$  is  $(\kappa + \mu)$ -strongly convex;
- $\mathbb{E}[h_k(x)|\mathcal{F}_{k-1}] \leq F(x) + (\kappa/2) ||x y_{k-1}||^2$  for some  $\alpha_{k-1}$ ;
- $\mathcal{M}$  "knows" the exact minimizer  $x_k^*$  of  $h_k$  and a point  $x_k$  such that  $\mathbb{E}\left[F(x_k)\right] \leq \mathbb{E}\left[h_k^*\right] + \delta_k$  with  $\delta_k > 0$ .

## First example of a surrogate

Given  $g_k$  as a stochastic realization of  $\nabla f(y_{k-1})$ , we consider

$$h_k(x) := f(y_{k-1}) + g_k^{\top}(x - y_{k-1}) + \frac{\mu + \kappa}{2} \|x - y_{k-1}\|^2 + \psi(x),$$
 with the exact minimizer  $x_k^* = \text{Prox}_{\psi/(\mu + \kappa)} [y_{k-1} - g_k/(\mu + \kappa)].$ 

# Algorithm [1] (exact minimization)

#### **FOR** k = 1, ..., K **DO**

- using a fixed improvable  $\mathcal{M}$ , obtain  $x_k$  and  $x_k^*$ ;
- update the extrapolated sequence

$$y_k = x_k^* + \beta_k(x_k^* - x_{k-1}) + \frac{(\kappa + \mu)(1 - \alpha_k)}{\kappa}(x_k - x_k^*),$$
 (3)

where  $\alpha_k$ ,  $\beta_k$  are from standard Nesterov extrapolation technique. **OUTPUT:**  $x_k$  (final estimate).

As a result, new accelerated SGD algorithm, converging as (1).

## Algorithm [2] (inexact minimization)

In some situations, the surrogate function  $h_k$  is such that  $x_k^*$  is not available, for example when

$$h_k(x) := F(x) + (\kappa/2) ||x - y_{k-1}||^2,$$

Then, in Equation (3) we use  $x_k$  instead of  $x_k^*$ . This results in **new** multi-stage algorithms with the convergence (2).

#### Advertisement for post-doc positions in Grenoble

in machine learning, optimization, computer vision, and ski.

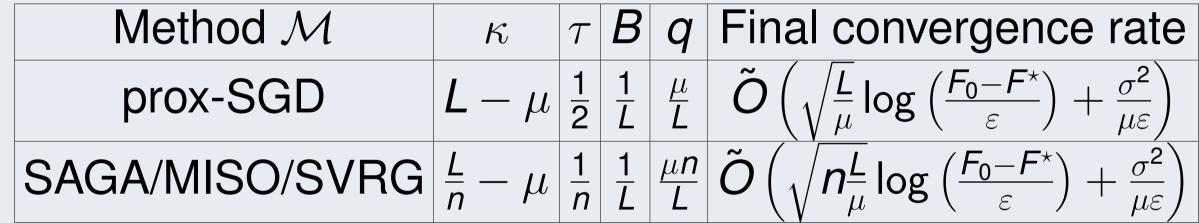
#### Basic restart schemes

In order to get converging algorithms out of Alg. [1] and [2], we sometimes address to restart procedure with mini-batching:

- at stage k, choose a target accuracy  $\varepsilon_k = \varepsilon_{k-1}/2$
- set up a mini-batch of size  $b_k = 2b_{k-1}$  to sample gradients, so that  $\sigma^2$  becomes  $\sigma^2/b_k$  at the stage k;
- minimize the objective up to accuracy  $\varepsilon_k$  using  $O(b_k/\tau)$  steps of  $\mathcal{M}$  and using previous solution as a "warm start".

#### Examples of final improvements

We provide practical choices for  $\kappa$  for different algorithms dealing with stochastic perturbations.



We see that we breach the optimal linear bias part, while preserving optimal robustness to noise.

# Experiments on logistic regression ( $\sigma^2 = 0$ on top)

- (left) CIFAR-10 represented by using a two-layer unsupervised convolutional neural network (n = 50000).
- (center) dataset with gene expressions data and the binary labels (n = 295);
- (right) Pascal Large Scale Learning Challenge (n = 250000);

