

Generic Acceleration Schemes for Gradient-Based Optimization

Hongzhou Lin¹, Julien Mairal¹, Zaid Harchaoui²

¹Inria, Grenoble

²University of Washington

Séminaire d'optimisation
Paris, 2018



An alternate title: Acceleration by Smoothing

**An alternate title:
Acceleration by Smoothing**

**Another one:
Recent Variants
of the Inexact Proximal Point Algorithm**

Collaborators



Hongzhou
Lin



Zaid
Harchaoui



Dima
Drusvyatskiy



Courtney
Paquette

Publications and pre-prints

H. Lin, J. Mairal and Z. Harchaoui. Catalyst Acceleration for First-order Convex Optimization: from Theory to Practice. *arXiv:1712.05654*. 2017.

H. Lin, J. Mairal and Z. Harchaoui. A Generic Quasi-Newton Algorithm for Faster Gradient-Based Optimization. *arXiv:1610.00960*. 2017

C. Paquette, H. Lin, D. Drusvyatskiy, J. Mairal, Z. Harchaoui. Catalyst for Gradient-Based Non-Convex Optimization. *AISTATS*. 2018

H. Lin, J. Mairal and Z. Harchaoui. A Universal Catalyst for First-Order Optimization. *Adv. NIPS* 2015.

Main motivation

Minimizing large finite sums

Consider the minimization of a large sum of convex functions

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(x) + \psi(x) \right\},$$

where each f_i is **L-smooth and convex** and ψ is a convex regularization penalty but not necessarily differentiable.

Motivation

Our goal is to accelerate existing algorithms

- with **Nesterov's principles**;
- with **Quasi-Newton** heuristics;

Why do large finite sums matter?

Empirical risk minimization

$$\min_{x \in \mathbb{R}^p} \left\{ F(x) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(x) + \psi(x) \right\},$$

- Typically, x represents **model parameters**.
- Each function f_i measures the **fidelity** of x to a data point.
- ψ is a **regularization function** to prevent overfitting.

For instance, given training data $(y_i, z_i)_{i=1, \dots, n}$ with features z_i in \mathbb{R}^p and labels y_i in $\{-1, +1\}$, we may want to predict y_i by $\text{sign}(\langle z_i, x \rangle)$. The functions f_i measure how far the prediction is from the true label.

This would be a **classification problem with a linear model**.

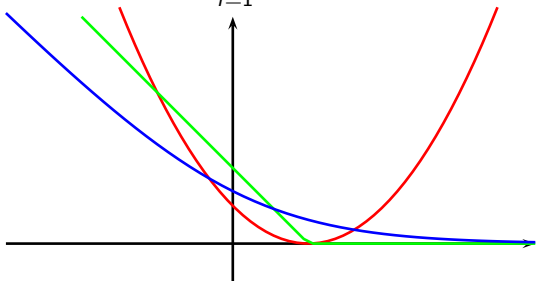
Why large finite sums matter?

A few examples

Ridge regression:
$$\min_{x \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (y_i - \langle x, z_i \rangle)^2 + \frac{\lambda}{2} \|x\|_2^2.$$

Linear SVM:
$$\min_{x \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i \langle x, z_i \rangle) + \frac{\lambda}{2} \|x\|_2^2.$$

Logistic regression:
$$\min_{x \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i \langle x, z_i \rangle}) + \frac{\lambda}{2} \|x\|_2^2.$$



Why does the composite problem matter?

A few examples

Ridge regression:
$$\min_{x \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (y_i - \langle x, z_i \rangle)^2 + \frac{\lambda}{2} \|x\|_2^2.$$

Linear SVM:
$$\min_{x \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i \langle x, z_i \rangle) + \frac{\lambda}{2} \|x\|_2^2.$$

Logistic regression:
$$\min_{x \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y_i \langle x, z_i \rangle} \right) + \frac{\lambda}{2} \|x\|_2^2.$$

The **squared l_2 -norm** penalizes large entries in x .

Why does the composite problem matter?

A few examples

Ridge regression:
$$\min_{x \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (y_i - \langle x, z_i \rangle)^2 + \lambda \|x\|_1.$$

Linear SVM:
$$\min_{x \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i \langle x, z_i \rangle)^2 + \lambda \|x\|_1.$$

Logistic regression:
$$\min_{x \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y_i \langle x, z_i \rangle} \right) + \lambda \|x\|_1.$$

When one knows in advance that x should be sparse, one should use a **sparsity-inducing** regularization such as the ℓ_1 -norm.

[Chen et al., 1999, Tibshirani, 1996].

Part I: How to address finite-sum problems?

How to minimize a large finite sum of functions?

$$\min_{x \in \mathbb{R}^p} \left\{ f(x) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(x) + \psi(x) \right\},$$

assuming here that the problem is μ -strongly convex.

We consider several alternatives

- Batch first-order methods (ISTA, FISTA).
- Stochastic first-order methods (SGD, mirror descent).
- Incremental first-order methods (SAG, SAGA, SDCA, MISO, ...).
- Quasi-Newton approaches (L-BFGS).

(Batch) gradient descent methods

Let us consider the composite problem

$$\min_{x \in \mathbb{R}^p} \{f(x) = f_0(x) + \psi(x)\},$$

where f_0 is convex, differentiable with L -Lipschitz continuous gradient and ψ is convex, but not necessarily differentiable.

The classical forward-backward/ISTA algorithm

$$x_k \leftarrow \arg \min_{x \in \mathbb{R}^p} \frac{1}{2} \left\| x - \left(x_{k-1} - \frac{1}{L} \nabla f_0(x_{k-1}) \right) \right\|_2^2 + \frac{1}{L} \psi(x).$$

- $f(x_k) - f^* = O(1/k)$ for **convex** problems;
- $f(x_k) - f^* = O((1 - \mu/L)^k)$ for **μ -strongly convex** problems;

[Nowak and Figueiredo, 2001, Daubechies et al., 2004, Combettes and Wajs, 2006, Beck and Teboulle, 2009a, Wright et al., 2009, Nesterov, 2013]...

Accelerated gradient descent methods

Nesterov introduced in the 80's an acceleration scheme for the gradient descent algorithm. It was generalized later to the composite setting.

FISTA

$$x_k \leftarrow \arg \min_{x \in \mathbb{R}^p} \frac{1}{2} \left\| x - \left(y_{k-1} - \frac{1}{L} \nabla f_0(y_{k-1}) \right) \right\|_2^2 + \frac{1}{L} \psi(x);$$

$$\text{Find } \alpha_k > 0 \text{ s.t. } \alpha_k^2 = (1 - \alpha_k) \alpha_{k-1}^2 + \frac{\mu}{L} \alpha_k;$$

$$y_k \leftarrow x_k + \beta_k (x_k - x_{k-1}) \quad \text{with} \quad \beta_k = \frac{\alpha_{k-1} (1 - \alpha_{k-1})}{\alpha_{k-1}^2 + \alpha_k}.$$

- $f(x_k) - f^* = O(1/k^2)$ for **convex** problems;
- $f(x_k) - f^* = O((1 - \sqrt{\mu/L})^k)$ for **μ -strongly convex** problems;
- Acceleration works in many practical cases.

see also [Nesterov, 1983, 2004, 2013]

Stochastic gradient descent methods

... or the recent return of Robins and Monroe, 1951. Consider

$$\min_{x \in \mathbb{R}^p} \left\{ f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}.$$

At iteration k , select at random an index i_k , and perform the update

$$x_k \leftarrow x_{k-1} - \eta_k \nabla f_{i_k}(x_{k-1})$$

Stochastic gradient descent methods

... or the recent return of Robins and Monroe, 1951. Consider

$$\min_{x \in \mathbb{R}^p} \left\{ f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}.$$

At iteration k , select at random an index i_k , and perform the update

$$x_k \leftarrow x_{k-1} - \eta_k \nabla f_{i_k}(x_{k-1}) \quad (\text{note that } \mathbb{E}[\nabla f_{i_k}(x_{k-1})] = \nabla f(x_{k-1})).$$

Stochastic gradient descent methods

... or the recent return of Robins and Monroe, 1951. Consider

$$\min_{x \in \mathbb{R}^p} \left\{ f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}.$$

At iteration k , select at random an index i_k , and perform the update

$$x_k \leftarrow x_{k-1} - \eta_k \nabla f_{i_k}(x_{k-1}) \quad (\text{note that } \mathbb{E}[\nabla f_{i_k}(x_{k-1})] = \nabla f(x_{k-1})).$$

Main features vs. batch

- **Complexity per-iteration is n times smaller;**

Stochastic gradient descent methods

... or the recent return of Robins and Monroe, 1951. Consider

$$\min_{x \in \mathbb{R}^p} \left\{ f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}.$$

At iteration k , select at random an index i_k , and perform the update

$$x_k \leftarrow x_{k-1} - \eta_k \nabla f_{i_k}(x_{k-1}) \quad (\text{note that } \mathbb{E}[\nabla f_{i_k}(x_{k-1})] = \nabla f(x_{k-1})).$$

Main features vs. batch

- **Complexity per-iteration is n times smaller;**
- Convergence rate is slower: at most $O(1/k)$ for strongly-convex problems and $O(1/\sqrt{k})$ for convex ones, see [Nemirovski et al., 2009];

Stochastic gradient descent methods

... or the recent return of Robins and Monroe, 1951. Consider

$$\min_{x \in \mathbb{R}^p} \left\{ f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}.$$

At iteration k , select at random an index i_k , and perform the update

$$x_k \leftarrow x_{k-1} - \eta_k \nabla f_{i_k}(x_{k-1}) \quad (\text{note that } \mathbb{E}[\nabla f_{i_k}(x_{k-1})] = \nabla f(x_{k-1})).$$

Main features vs. batch

- **Complexity per-iteration is n times smaller;**
- Convergence rate is slower: at most $O(1/k)$ for strongly-convex problems and $O(1/\sqrt{k})$ for convex ones, see [Nemirovski et al., 2009];
- variants are **compatible with prox** ψ , e.g., [Duchi et al., 2011].

Stochastic gradient descent methods

... or the recent return of Robins and Monroe, 1951. Consider

$$\min_{x \in \mathbb{R}^p} \left\{ f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}.$$

At iteration k , select at random an index i_k , and perform the update

$$x_k \leftarrow x_{k-1} - \eta_k \nabla f_{i_k}(x_{k-1}) \quad (\text{note that } \mathbb{E}[\nabla f_{i_k}(x_{k-1})] = \nabla f(x_{k-1})).$$

Main features vs. batch

- **Complexity per-iteration is n times smaller;**
- Convergence rate is slower: at most $O(1/k)$ for strongly-convex problems and $O(1/\sqrt{k})$ for convex ones, see [Nemirovski et al., 2009];
- variants are **compatible with prox** ψ , e.g., [Duchi et al., 2011].
- Sometimes a bit difficult to tune. When well tuned, the speed-up to obtain a solution with moderate accuracy may be huge.

Incremental gradient descent methods

$$\min_{x \in \mathbb{R}^p} \left\{ f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}.$$

Several **randomized** algorithms are designed with one ∇f_i computed per iteration, with **fast convergence rates**, e.g., SAG [Schmidt et al., 2017]:

$$x_k \leftarrow x_{k-1} - \frac{\gamma}{Ln} \sum_{i=1}^n y_i^k \quad \text{with} \quad y_i^k = \begin{cases} \nabla f_i(x_{k-1}) & \text{if } i = i_k \\ y_i^{k-1} & \text{otherwise} \end{cases}.$$

Incremental gradient descent methods

$$\min_{x \in \mathbb{R}^p} \left\{ f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}.$$

Several **randomized** algorithms are designed with one ∇f_i computed per iteration, with **fast convergence rates**, e.g., SAG [Schmidt et al., 2017]:

$$x_k \leftarrow x_{k-1} - \frac{\gamma}{Ln} \sum_{i=1}^n y_i^k \quad \text{with} \quad y_i^k = \begin{cases} \nabla f_i(x_{k-1}) & \text{if } i = i_k \\ y_i^{k-1} & \text{otherwise} \end{cases}.$$

See also SVRG, SAGA, SDCA, MISO, Finito...

Some of these algorithms perform updates of the form

$$x_k \leftarrow x_{k-1} - \eta_k g_k \quad \text{with} \quad \mathbb{E}[g_k] = \nabla f(x_{k-1}),$$

but g_k has **lower variance** than in SGD.

[Schmidt et al., 2017, Xiao and Zhang, 2014, Defazio et al., 2014a,b, Shalev-Shwartz and Zhang, 2012, Mairal, 2015, Zhang and Xiao, 2015]

Incremental gradient descent methods

These methods achieve low (**worst-case**) complexity in expectation.
The number of gradients evaluations to ensure $f(x_k) - f^* \leq \varepsilon$ is

	$\mu > 0$
FISTA	$O\left(n\sqrt{\frac{L}{\mu}} \log\left(\frac{1}{\varepsilon}\right)\right)$
SVRG, SAG, SAGA, SDCA, MISO, Finito	$O\left(\max\left(n, \frac{L}{\mu}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$

Incremental gradient descent methods

These methods achieve low (**worst-case**) complexity in expectation.
The number of gradients evaluations to ensure $f(x_k) - f^* \leq \varepsilon$ is

	$\mu > 0$
FISTA	$O\left(n\sqrt{\frac{L}{\mu}} \log\left(\frac{1}{\varepsilon}\right)\right)$
SVRG, SAG, SAGA, SDCA, MISO, Finito	$O\left(\max\left(n, \frac{L}{\mu}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$

Main features vs. stochastic gradient descent

- Same complexity per-iteration (but higher memory footprint).

Incremental gradient descent methods

These methods achieve low (**worst-case**) complexity in expectation. The number of gradients evaluations to ensure $f(x_k) - f^* \leq \varepsilon$ is

	$\mu > 0$
FISTA	$O\left(n\sqrt{\frac{L}{\mu}} \log\left(\frac{1}{\varepsilon}\right)\right)$
SVRG, SAG, SAGA, SDCA, MISO, Finito	$O\left(\max\left(n, \frac{L}{\mu}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$

Main features vs. stochastic gradient descent

- Same complexity per-iteration (but higher memory footprint).
- **Faster convergence** (exploit the finite-sum structure).

Incremental gradient descent methods

These methods achieve low (**worst-case**) complexity in expectation. The number of gradients evaluations to ensure $f(x_k) - f^* \leq \varepsilon$ is

	$\mu > 0$
FISTA	$O\left(n\sqrt{\frac{L}{\mu}} \log\left(\frac{1}{\varepsilon}\right)\right)$
SVRG, SAG, SAGA, SDCA, MISO, Finito	$O\left(\max\left(n, \frac{L}{\mu}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$

Main features vs. stochastic gradient descent

- Same complexity per-iteration (but higher memory footprint).
- **Faster convergence** (exploit the finite-sum structure).
- **Less parameter tuning** than SGD.

Incremental gradient descent methods

These methods achieve low (**worst-case**) complexity in expectation. The number of gradients evaluations to ensure $f(x_k) - f^* \leq \varepsilon$ is

	$\mu > 0$
FISTA	$O\left(n\sqrt{\frac{L}{\mu}} \log\left(\frac{1}{\varepsilon}\right)\right)$
SVRG, SAG, SAGA, SDCA, MISO, Finito	$O\left(\max\left(n, \frac{L}{\mu}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$

Main features vs. stochastic gradient descent

- Same complexity per-iteration (but higher memory footprint).
- **Faster convergence** (exploit the finite-sum structure).
- **Less parameter tuning** than SGD.
- Some variants are **compatible with composite term** ψ .

Part II: Catalyst

An old idea

Old idea: Smooth the function and then optimize.

The Moreau-Yosida envelope

Given $f : \mathbb{R}^d \rightarrow \mathbb{R}$ a convex function, the Moreau-Yosida envelope of f is the function $F : \mathbb{R}^d \rightarrow \mathbb{R}$ defined as

$$F(x) = \min_{w \in \mathbb{R}^d} \left\{ f(w) + \frac{\kappa}{2} \|w - x\|^2 \right\}.$$

The **proximal operator** $p(x)$ is the unique minimizer of the problem.

The Moreau-Yosida regularization

$$F(x) = \min_{w \in \mathbb{R}^d} \left\{ f(w) + \frac{\kappa}{2} \|w - x\|^2 \right\}.$$

Basic properties [see Lemaréchal and Sagastizábal, 1997]

- Minimizing f and F is equivalent in the sense that

$$\min_{x \in \mathbb{R}^d} F(x) = \min_{x \in \mathbb{R}^d} f(x),$$

and the solution set of the two problems coincide with each other.

- F is continuously differentiable even when f is not and

$$\nabla F(x) = \kappa(x - p(x)).$$

In addition, ∇F is Lipschitz continuous with parameter $L_F = \kappa$.

- If f is μ -strongly convex then F is also strongly convex with parameter $\mu_F = \frac{\mu\kappa}{\mu + \kappa}$.

The Moreau-Yosida regularization

$$F(x) = \min_{w \in \mathbb{R}^d} \left\{ f(w) + \frac{\kappa}{2} \|w - x\|^2 \right\}.$$

Basic properties [see Lemaréchal and Sagastizábal, 1997]

- Minimizing f and F is equivalent in the sense that

$$\min_{x \in \mathbb{R}^d} F(x) = \min_{x \in \mathbb{R}^d} f(x),$$

and the solution set of the two problems coincide with each other.

- F is continuously differentiable even when f is not and

$$\nabla F(x) = \kappa(x - p(x)).$$

In addition, ∇F is Lipschitz continuous with parameter $L_F = \kappa$.

F enjoys nice properties: smoothness, (strong) convexity and we can control its condition number $1/q = 1 + \kappa/\mu$.

The proximal point algorithm

A naive approach consists of **minimizing the smoothed objective F instead of f** with a method designed for smooth optimization.

Consider indeed

$$x_{k+1} = x_k - \frac{1}{\kappa} \nabla F(x_k).$$

By rewriting the gradient $\nabla F(x_k)$ as $\kappa(x_k - p(x_k))$, we obtain

$$x_{k+1} = p(x_k) = \arg \min_{w \in \mathbb{R}^p} \left\{ f(w) + \frac{\kappa}{2} \|w - x_k\|^2 \right\}.$$

This is exactly the **proximal point algorithm** [Martinet, 1970, Rockafellar, 1976].

The accelerated proximal point algorithm

Consider now

$$x_{k+1} = y_k - \frac{1}{\kappa} \nabla F(y_k) \quad \text{and} \quad y_{k+1} = x_{k+1} + \beta_{k+1}(x_{k+1} - x_k),$$

where β_{k+1} is a Nesterov-like extrapolation parameter. We may now rewrite the update using the value of ∇F , which gives:

$$x_{k+1} = p(y_k) \quad \text{and} \quad y_{k+1} = x_{k+1} + \beta_{k+1}(x_{k+1} - x_k)$$

This is the **accelerated proximal point algorithm** of Güler [1992].

The accelerated proximal point algorithm

Consider now

$$x_{k+1} = y_k - \frac{1}{\kappa} \nabla F(y_k) \quad \text{and} \quad y_{k+1} = x_{k+1} + \beta_{k+1}(x_{k+1} - x_k),$$

where β_{k+1} is a Nesterov-like extrapolation parameter. We may now rewrite the update using the value of ∇F , which gives:

$$x_{k+1} = p(y_k) \quad \text{and} \quad y_{k+1} = x_{k+1} + \beta_{k+1}(x_{k+1} - x_k)$$

This is the **accelerated proximal point algorithm** of Güler [1992].

Remarks

- F may be **better conditioned** than f when $1 + \kappa/\mu \leq L/\mu$;
- Computing $p(y_k)$ has a cost!

A fresh look at Catalyst [Lin, Mairal, and Harchaoui, 2015]

Catalyst is a particular **accelerated proximal point algorithm with inexact gradients** [Güler, 1992].

$$x_{k+1} \approx p(y_k) \quad \text{and} \quad y_{k+1} = x_{k+1} + \beta_{k+1}(x_{k+1} - x_k)$$

The quantity x_{k+1} is obtained by using an optimization method \mathcal{M} for approximately solving:

$$x_{k+1} \approx \arg \min_{w \in \mathbb{R}^p} \left\{ f(w) + \frac{\kappa}{2} \|w - y_k\|^2 \right\},$$

Catalyst provides Nesterov's acceleration to \mathcal{M} with...

- **restart strategies** for solving the sub-problems;
- **global complexity analysis** resulting in theoretical acceleration;
- **optimal balancing between outer and inner computations.**

see also [Frostig et al., 2015, Schmidt et al., 2011, Salzo and Villa, 2012, Devolder et al., 2014, Shalev-Shwartz and Zhang, 2016]

This work

Contributions

- **Generic acceleration scheme**, which applies to algorithms \mathcal{M} that have **linear convergence rates** for strongly convex problems..
- Provides explicit **support to non-strongly convex objectives**.
- Complexity analysis for μ -strongly convex objectives.
- Complexity analysis for non-strongly convex objectives.
- Extension to **non-convex optimization** by Paquette, Lin, Drusvyatskiy, Mairal, and Harchaoui [2017].

Requirements on \mathcal{M}

Linear convergence

- Say a sub-problem consists of minimizing h ; we want \mathcal{M} to produce a sequence of iterates $(z_t)_{t \geq 0}$ with **linear convergence rate**

$$h(z_t) - h^* \leq C_{\mathcal{M}}(1 - \tau_{\mathcal{M}})^t (h(z_0) - h^*),$$

which may possibly hold only **in expectation** if \mathcal{M} is randomized.

- **No assumption** is made on the behavior of \mathcal{M} for **non-strongly convex problems**.
- Variants may be allowed when linear convergence is stated in terms of dual certificate.

When do we stop the method \mathcal{M} ?

Three strategies to balance outer and inner computations

- (a) use a **pre-defined sequence** $(\varepsilon_k)_{k \geq 0}$ and stop the optimization method \mathcal{M} when the sub-problems $\min h_k$ satisfies

$$h_k(z_t) - h_k^* \leq \varepsilon_k.$$

- (b) use a **pre-defined sequence** $(\delta_k)_{k \geq 0}$ and stop the optimization method \mathcal{M} when the sub-problems $\min h_k$ satisfies

$$h_k(z_t) - h_k^* \leq \frac{\delta_k}{2} \|z_t - y_k\|^2.$$

- (c) use a **pre-defined budget** $T_{\mathcal{M}}$ of iterations of the method \mathcal{M} .

When do we stop the method \mathcal{M} ?

Three strategies to balance outer and inner computations

- (a) use a **pre-defined sequence** $(\varepsilon_k)_{k \geq 0}$ and stop the optimization method \mathcal{M} when the sub-problems $\min h_k$ satisfies

$$h_k(z_t) - h_k^* \leq \varepsilon_k.$$

- (b) use a **pre-defined sequence** $(\delta_k)_{k \geq 0}$ and stop the optimization method \mathcal{M} when the sub-problems $\min h_k$ satisfies

$$h_k(z_t) - h_k^* \leq \frac{\delta_k}{2} \|z_t - y_k\|^2.$$

- (c) use a **pre-defined budget** $T_{\mathcal{M}}$ of iterations of the method \mathcal{M} .

Remark

- (c) implies (a) and requires $T_{\mathcal{M}}$ to be larger than necessary in practice; it leads to the simplest and most effective strategies.

When do we stop the method \mathcal{M} ?

Three strategies for μ -strongly convex objectives f

(a) use

$$\varepsilon_k = \frac{1}{2}C(1 - \rho)^{k+1} \quad \text{with} \quad C \geq f(x_0) - f^* \quad \text{and} \quad \rho < \sqrt{q}.$$

where q is the inverse of the condition number of F : $q = \frac{\mu}{(\mu + \kappa)}$

(b) use

$$\delta_k = \frac{\sqrt{q}}{2 - \sqrt{q}}.$$

(c) use a **pre-defined budget** $T_{\mathcal{M}}$ of iterations of the method \mathcal{M} for solving each sub-problem with

$$T_{\mathcal{M}} = \frac{1}{\tau_{\mathcal{M}}} \log \left(19C_{\mathcal{M}} \frac{L + \kappa}{\kappa} \right). \quad (\text{be more aggressive in practice})$$

When do we stop the method \mathcal{M} ?

Three strategies for $\mu = 0$

(a) use

$$\varepsilon_k = \frac{f(x_0) - f^*}{2(k+1)^{4+\gamma}} \quad \text{with } \gamma > 0.$$

(b) use

$$\delta_k = \frac{1}{(k+1)^2}.$$

(c) use a **pre-defined budget** T_k of iterations of the method \mathcal{M} for solving each sub-problem h_k with

$$T_k = O(\log(k)) \quad (\text{use a constant in practice})$$

Other implementation details

See the arXiv paper for

- Nesterov's extrapolation parameters (fairly standard).
- restart strategies for solving the sub-problems.

Other implementation details

See the arXiv paper for

- Nesterov's extrapolation parameters (fairly standard).
- restart strategies for solving the sub-problems.

Spoiler: optimal balance for inner/outer computations

To choose κ , maximize

$$\frac{\tau_{\mathcal{M}}}{\sqrt{\mu + \kappa}}.$$

Remember that $\tau_{\mathcal{M}}$ drives the convergence rate for the sub-problems

$$h(w_t) - h^* \leq C_{\mathcal{M}}(1 - \tau_{\mathcal{M}})^t(h(w_0) - h^*).$$

For the standard gradient descent method, use $\kappa = L - 2\mu$.

Outer-loop convergence analysis

With strong convexity

Using strategy (a),

$$f(x_k) - f^* \leq C(1 - \rho)^{k+1}(f(x_0) - f^*) \quad \text{with } \rho < \sqrt{q},$$

and a similar result holds for (b).

Without strong convexity

Using strategy (b),

$$f(x_k) - f^* \leq \frac{4\kappa \|x_0 - x^*\|^2}{(k+1)^2}.$$

and a similar result holds for (a).

Inner-loop convergence analysis

Using appropriate restart strategies, the inner-loop stopping criterions are satisfied after T_k iterations, where

$$T_k = \tilde{O}\left(\frac{1}{\tau_{\mathcal{M}}}\right) \quad \text{when } \mu > 0,$$

and

$$T_k = \tilde{O}\left(\frac{\log(k)}{\tau_{\mathcal{M}}}\right) \quad \text{when } \mu = 0.$$

The \tilde{O} hides logarithmic quantities in μ, κ and universal constants.

Global complexity analysis

By combining the two previous strategies, we obtain that the guarantee $f(x_k) - f^* \leq \varepsilon$ is achieved after N iterations of the method \mathcal{M} , where

$$N = \tilde{O} \left(\frac{1}{\tau_{\mathcal{M}} \sqrt{q}} \log \left(\frac{1}{\varepsilon} \right) \right) \quad \text{when } \mu > 0,$$

and

$$N = \tilde{O} \left(\frac{1}{\tau_{\mathcal{M}}} \sqrt{\frac{\kappa}{\varepsilon}} \log \left(\frac{1}{\varepsilon} \right) \right) \quad \text{when } \mu = 0.$$

Similar results hold also for randomized algorithms.

Global complexity analysis

By combining the two previous strategies, we obtain that the guarantee $f(x_k) - f^* \leq \varepsilon$ is achieved after N iterations of the method \mathcal{M} , where

$$N = \tilde{O} \left(\frac{1}{\tau_{\mathcal{M}} \sqrt{q}} \log \left(\frac{1}{\varepsilon} \right) \right) \quad \text{when } \mu > 0,$$

and

$$N = \tilde{O} \left(\frac{1}{\tau_{\mathcal{M}}} \sqrt{\frac{\kappa}{\varepsilon}} \log \left(\frac{1}{\varepsilon} \right) \right) \quad \text{when } \mu = 0.$$

Similar results hold also for randomized algorithms.

Theoretical choice of κ

maximize

$$\frac{\tau_{\mathcal{M}}}{\sqrt{\mu + \kappa}}.$$

For gradient descent, $\tau_{\mathcal{M}} = \frac{\mu + \kappa}{L + \kappa} \Rightarrow \kappa = L - 2\mu \Rightarrow \frac{1}{\tau_{\mathcal{M}} \sqrt{q}} \leq 2\sqrt{\frac{L}{\mu}}$

Applications

Expected computational complexity in the regime $n \leq L/\mu$ when $\mu > 0$,

	$\mu > 0$	$\mu = 0$	Catalyst $\mu > 0$	Cat. $\mu = 0$
FG	$O\left(n\left(\frac{L}{\mu}\right)\log\left(\frac{1}{\varepsilon}\right)\right)$	$O\left(n\frac{L}{\varepsilon}\right)$	$\tilde{O}\left(n\sqrt{\frac{L}{\mu}}\log\left(\frac{1}{\varepsilon}\right)\right)$	$\tilde{O}\left(n\sqrt{\frac{L}{\varepsilon}}\right)$
SAG	$O\left(\frac{L}{\mu}\log\left(\frac{1}{\varepsilon}\right)\right)$		NA	$\tilde{O}\left(\sqrt{\frac{nL}{\mu}}\log\left(\frac{1}{\varepsilon}\right)\right)$
SAGA				
Finito/MISO				
SDCA				
SVRG				
Acc-FG	$O\left(n\sqrt{\frac{L}{\mu}}\log\left(\frac{1}{\varepsilon}\right)\right)$	$O\left(n\sqrt{\frac{L}{\varepsilon}}\right)$	no acceleration	
Acc-SDCA	$\tilde{O}\left(\sqrt{\frac{nL}{\mu}}\log\left(\frac{1}{\varepsilon}\right)\right)$	NA		

Part III: QNing

Limited-Memory BFGS (L-BFGS)

Pros

- **one of the largest practical success of smooth optimization.**

Limited-Memory BFGS (L-BFGS)

Pros

- **one of the largest practical success of smooth optimization.**

Cons

- worst-case convergence rates for strongly-convex functions are linear, but **much worse than the gradient descent method.**
- proximal variants typically requires solving many times

$$\min_{x \in \mathbb{R}^d} \frac{1}{2}(x - z)B_k(z - z) + \psi(x).$$

- no guarantee of approximating the Hessian.

An old idea (again)

Old idea: Smooth the function and then optimize.

- The strategy appears in early work about variable metric bundle methods. [Chen and Fukushima, 1999, Fukushima and Qi, 1996, Mifflin, 1996, Fuentes, Malick, and Lemaréchal, 2012, Burke and Qian, 2000] ...

The Moreau-Yosida envelope

Given $f : \mathbb{R}^d \rightarrow \mathbb{R}$ a convex function, the Moreau-Yosida envelope of f is the function $F : \mathbb{R}^d \rightarrow \mathbb{R}$ defined as

$$F(x) = \min_{w \in \mathbb{R}^d} \left\{ f(w) + \frac{\kappa}{2} \|w - x\|^2 \right\}.$$

The **proximal operator** $p(x)$ is the unique minimizer of the problem.

Main recipe

- L-BFGS applied to the **smoothed objective** F with **inexact gradients** [see Friedlander and Schmidt, 2012].
- inexact gradients are obtained by **solving sub-problems** using a first-order optimization method \mathcal{M} ;
- ideally, \mathcal{M} is **able to adapt to the problem structure** (finite sum, composite regularization).
- replace L-BFGS steps by proximal point steps if no sufficient decrease is estimated \Rightarrow **no line search on F** ;

Obtaining inexact gradients

Algorithm Procedure ApproxGradient

input Current point x in \mathbb{R}^d ; smoothing parameter $\kappa > 0$.

- 1: Compute the approximate mapping using an optimization method \mathcal{M} :

$$z \approx \arg \min_{w \in \mathbb{R}^d} \left\{ h(w) \triangleq f(w) + \frac{\kappa}{2} \|w - x\|^2 \right\},$$

- 2: Estimate the gradient $\nabla F(x)$

$$g = \kappa(x - z).$$

output approximate gradient estimate g , objective value $F_a \triangleq h(z)$, proximal mapping z .

Algorithm QuickeNing

input x_0 in \mathbb{R}^p ; number of iterations K ; $\kappa > 0$; minimization algorithm \mathcal{M} .

1: Initialization: $(g_0, F_0, z_0) = \text{ApproxGradient}(x_0, \mathcal{M})$; $B_0 = \kappa I$.

2: **for** $k = 0, \dots, K - 1$ **do**

3: Perform the Quasi-Newton step

$$x_{\text{test}} = x_k - B_k^{-1} g_k$$

$$(g_{\text{test}}, F_{\text{test}}, z_{\text{test}}) = \text{ApproxGradient}(x_{\text{test}}, \mathcal{M}) .$$

4: **if** $F_{\text{test}} \leq F_k - \frac{1}{2\kappa} \|g_k\|^2$, **then**

5: $(x_{k+1}, g_{k+1}, F_{k+1}, z_{k+1}) = (x_{\text{test}}, g_{\text{test}}, F_{\text{test}}, z_{\text{test}})$.

6: **else**

7: Update the current iterate with the last proximal mapping:

$$x_{k+1} = z_k = x_k - (1/\kappa)g_k$$

$$(g_{k+1}, F_{k+1}, z_{k+1}) = \text{ApproxGradient}(x_{k+1}, \mathcal{M}) .$$

8: **end if**

9: update $B_{k+1} = \text{L-BFGS}(B_k, x_{k+1} - x_k, g_{k+1} - g_k)$.

10: **end for**

output last proximal mapping z_K (solution).

Algorithm QuickeNing

input x_0 in \mathbb{R}^p ; number of iterations K ; $\kappa > 0$; minimization algorithm \mathcal{M} .

1: Initialization: $(g_0, F_0, z_0) = \text{ApproxGradient}(x_0, \mathcal{M})$; $B_0 = \kappa I$.

2: **for** $k = 0, \dots, K - 1$ **do**

3: Perform the Quasi-Newton step

$$x_{\text{test}} = x_k - B_k^{-1} g_k$$
$$(g_{\text{test}}, F_{\text{test}}, z_{\text{test}}) = \text{ApproxGradient}(x_{\text{test}}, \mathcal{M}).$$

The main characters:

- the sequence $(x_k)_{k \geq 0}$ that minimizes F ;
- the sequence $(z_k)_{k \geq 0}$ produced by \mathcal{M} that minimizes f ;
- the gradient approximations $g_k \approx \nabla F(x_k)$;
- the function value approximations $F_k \approx F(x_k)$;
- an L-BFGS update with inexact gradients;
- an approximate sufficient descent condition.

10: **end for**

output last proximal mapping z_K (solution).

Requirements on \mathcal{M} and restarts

Method \mathcal{M}

- Say a sub-problem consists of minimizing h ; we want \mathcal{M} to produce a sequence of iterates $(w_t)_{t \geq 0}$ with **linear convergence rate**

$$h(w_t) - h^* \leq C_{\mathcal{M}}(1 - \tau_{\mathcal{M}})^t (h(w_0) - h^*).$$

Restarts

- When f is smooth, we **initialize** $w_0 = x$ when solving

$$\min_{w \in \mathbb{R}^d} \left\{ f(w) + \frac{\kappa}{2} \|w - x\|^2 \right\}.$$

- When $f = f_0 + \psi$ is composite, we use the initialization

$$w_0 = \arg \min_{w \in \mathbb{R}^d} \left\{ f_0(x) + \langle \nabla f_0(x), w - x \rangle + \frac{L + \kappa}{2} \|w - x\|^2 + \psi(w) \right\}.$$

When do we stop the method \mathcal{M} ?

Three strategies to balance outer and inner computations

- (a) use a **pre-defined sequence** $(\varepsilon_k)_{k \geq 0}$ and stop the optimization method \mathcal{M} when the approximate proximal mapping is ε_k -accurate.
- (b) define an **adaptive stopping criterion** that depends on quantities that are available at iteration k .
- (c) use a **pre-defined budget** $T_{\mathcal{M}}$ of iterations of the method \mathcal{M} for solving each sub-problem.

When do we stop the method \mathcal{M} ?

Three strategies to balance outer and inner computations

- (a) use a **pre-defined sequence** $(\varepsilon_k)_{k \geq 0}$ and stop the optimization method \mathcal{M} when the approximate proximal mapping is ε_k -accurate.
- (b) define an **adaptive stopping criterion** that depends on quantities that are available at iteration k .
- (c) use a **pre-defined budget** $T_{\mathcal{M}}$ of iterations of the method \mathcal{M} for solving each sub-problem.

Remarks

- We have already seen all of this for Catalyst We have already seen all of this for Catalyst..

When do we stop the method \mathcal{M} ?

Three strategies for μ -strongly convex objectives f

- (a) use a **pre-defined sequence** $(\varepsilon_k)_{k \geq 0}$ and stop the optimization method \mathcal{M} when the approximate proximal mapping is ε_k -accurate.

$$\varepsilon_k = \frac{1}{2}C(1 - \rho)^{k+1} \quad \text{with} \quad C \geq f(x_0) - f^* \quad \text{and} \quad \rho = \frac{\mu}{4(\mu + \kappa)}.$$

- (b) For minimizing $h(w) = f(w) + (\kappa/2)\|w - x\|^2$, stop when

$$h(w_t) - h^* \leq \frac{\kappa}{36}\|w_t - x\|^2.$$

- (c) use a **pre-defined budget** $T_{\mathcal{M}}$ of iterations of the method \mathcal{M} for solving each sub-problem with

$$T_{\mathcal{M}} = \frac{1}{\tau_{\mathcal{M}}} \log \left(19C_{\mathcal{M}} \frac{L + \kappa}{\kappa} \right). \quad (\text{be more aggressive in practice})$$

Remarks and worst-case global complexity

Composite objectives and sparsity

Consider a composite problem with a sparse solution (e.g., $\psi = \ell_1$). The method produces two sequences $(x_k)_{k \geq 0}$ and $(z_k)_{k \geq 0}$;

- $F(x_k) \rightarrow F^*$, minimizes the **smoothed objective** \Rightarrow no sparsity;
- $f(z_k) \rightarrow f^*$, minimizes the **true objective** \Rightarrow the iterates may be sparse if \mathcal{M} handles composite optimization problems;

Global complexity

The number of iterations of \mathcal{M} to guarantee $f(z_k) - f^* \leq \varepsilon$ is at most

- $\tilde{O}\left(\frac{\mu + \kappa}{\tau_{\mathcal{M}} \mu} \log(1/\varepsilon)\right)$ for μ -strongly convex problems.
- $\tilde{O}\left(\frac{\kappa R^2}{\tau_{\mathcal{M}} \varepsilon}\right)$ for convex problems.

Global Complexity and choice of κ

Example for gradient descent

With the right step-size, we have $\tau_{\mathcal{M}} = (\mu + \kappa)/(L + \kappa)$ and the complexity for $\mu > 0$ becomes

$$\tilde{O}\left(\frac{L + \kappa}{\mu} \log(1/\varepsilon)\right).$$

Example for SVRG for minimizing the sum of n functions

$\tau_{\mathcal{M}} = \min(1/n, (\mu + \kappa)/(L + \kappa))$ and the complexity for $\mu > 0$ is

$$\tilde{O}\left(\max\left(\frac{\mu + \kappa}{\mu} n, \frac{L + \kappa}{\mu}\right) \log(1/\varepsilon)\right).$$

Global Complexity and choice of κ

Example for gradient descent

With the right step-size, we have $\tau_{\mathcal{M}} = (\mu + \kappa)/(L + \kappa)$ and the complexity for $\mu > 0$ becomes

$$\tilde{O}\left(\frac{L + \kappa}{\mu} \log(1/\varepsilon)\right).$$

Example for SVRG for minimizing the sum of n functions

$\tau_{\mathcal{M}} = \min(1/n, (\mu + \kappa)/(L + \kappa))$ and the complexity for $\mu > 0$ is

$$\tilde{O}\left(\max\left(\frac{\mu + \kappa}{\mu} n, \frac{L + \kappa}{\mu}\right) \log(1/\varepsilon)\right).$$

QuickeNing does not provide any theoretical acceleration, but it does not degrade significantly the worst-case performance of \mathcal{M} (unlike L-BFGS vs gradient descent).

Global Complexity and choice of κ

Example for gradient descent

With the right step-size, we have $\tau_{\mathcal{M}} = (\mu + \kappa)/(L + \kappa)$ and the complexity for $\mu > 0$ becomes

$$\tilde{O}\left(\frac{L + \kappa}{\mu} \log(1/\varepsilon)\right).$$

Example for SVRG for minimizing the sum of n functions

$\tau_{\mathcal{M}} = \min(1/n, (\mu + \kappa)/(L + \kappa))$ and the complexity for $\mu > 0$ is

$$\tilde{O}\left(\max\left(\frac{\mu + \kappa}{\mu} n, \frac{L + \kappa}{\mu}\right) \log(1/\varepsilon)\right).$$

Then, how to choose κ ?

- (i) assume that L-BFGS steps do as well as Nesterov.
- (ii) **choose κ as in Catalyst.**

Experiments: formulations

- ℓ_2 -regularized Logistic Regression:

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \log \left(1 + \exp(-b_i a_i^T x) \right) + \frac{\mu}{2} \|x\|^2,$$

- ℓ_1 -regularized Linear Regression (LASSO):

$$\min_{x \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (b_i - a_i^T x)^2 + \lambda \|x\|_1,$$

- $\ell_1 - \ell_2^2$ -regularized Linear Regression (Elastic-Net):

$$\min_{x \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (b_i - a_i^T x)^2 + \lambda \|x\|_1 + \frac{\mu}{2} \|x\|^2,$$

Experiments: Datasets

We consider four standard machine learning datasets with different characteristics in terms of size and dimension

name	covtype	alpha	real-sim	rcv1
n	581 012	250 000	72 309	781 265
d	54	500	20 958	47 152

- we simulate the ill-conditioned regime $\mu = 1/(100n)$;
- λ for the Lasso leads to about 10% non-zero coefficients.

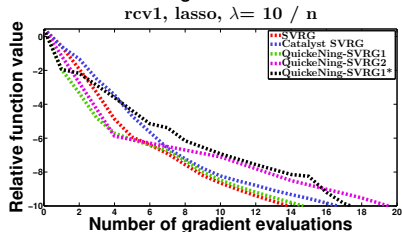
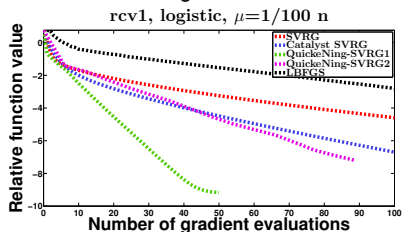
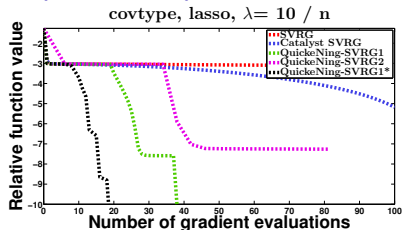
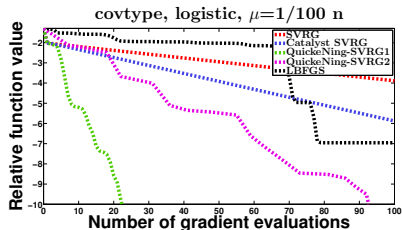
Experiments: QuickeNing-SVRG

We consider the methods

- **SVRG**: the Prox-SVRG algorithm of Xiao and Zhang [2014].
- **Catalyst-SVRG**: Catalyst applied to SVRG;
- **L-BFGS** (for smooth objectives): Mark Schmidt's implementation.
- **QuickeNing-SVRG1**: QuickeNing with aggressive strategy (c): one pass over the data in the inner loop.
- **QuickeNing-SVRG2**: strategy (b), compatible with theory.

We produce 12 figures (3 formulations, 4 datasets).

Experiments: QuickeNing-SVRG (log scale)



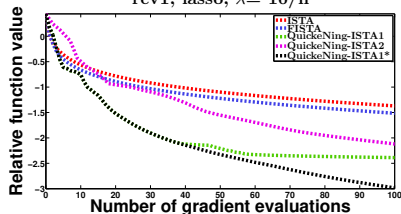
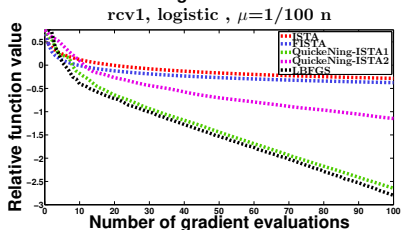
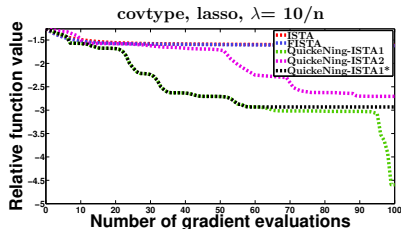
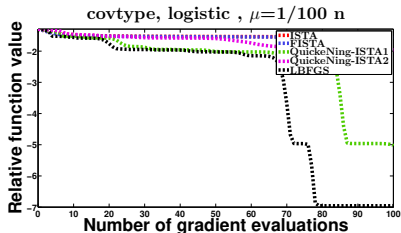
- QuickeNing-SVRG1 \geq SVRG, QuickeNing-SVRG2;
- QuickeNing-SVRG2 \geq SVRG;
- QuickeNing-SVRG1 \geq Catalyst-SVRG in 10/12 cases.

Experiments: QuickeNing-ISTA

We consider the methods

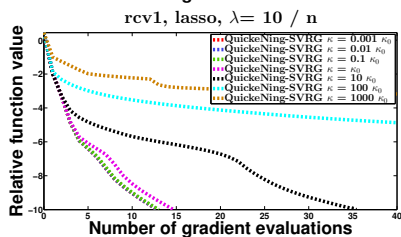
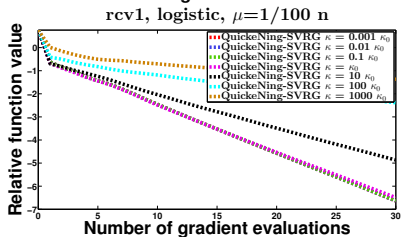
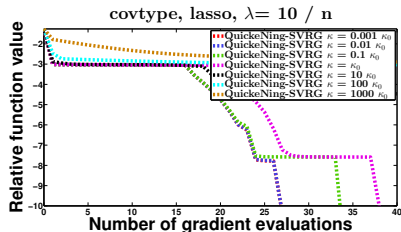
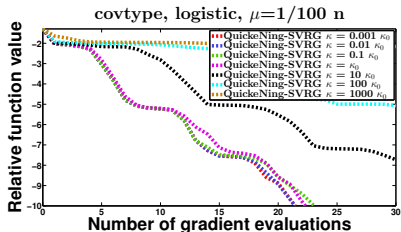
- **ISTA**: the proximal gradient descent method with line search.
- **FISTA**: the accelerated ISTA of Beck and Teboulle [2009b].
- **L-BFGS** (for smooth objectives): Mark Schmidt's implementation.
- **QuickeNing-ISTA1**: QuickeNing with aggressive strategy (c): one pass over the data in the inner loop.
- **QuickeNing-ISTA2**: strategy (b), compatible with theory.

Experiments: QuickeNing-ISTA (log scale)



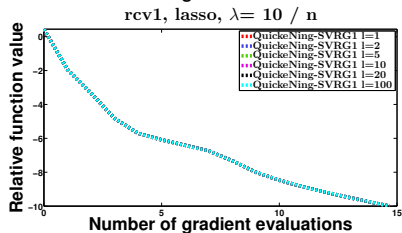
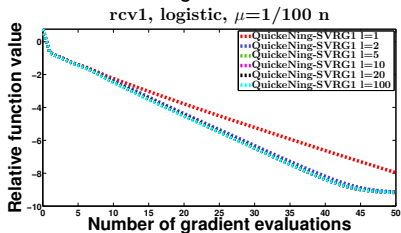
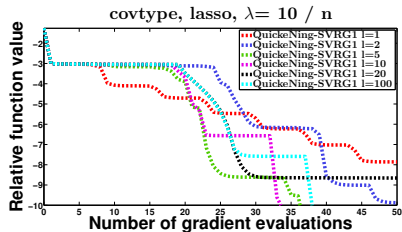
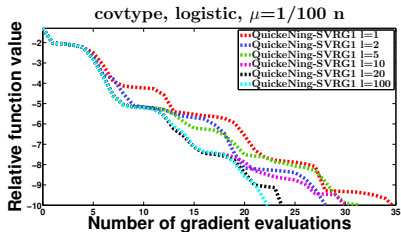
- L-BFGS (for smooth f) is slightly better than QuickeNing-ISTA1;
- QuickeNing-ISTA \geq or \gg FISTA in 11/12 cases.
- QuickeNing-ISTA1 \geq QuickeNing-ISTA2.

Experiments: Influence of κ



- κ_0 is the parameter (same as in Catalyst) used in all experiments;
- QuickeNing slows down when using $\kappa > \kappa_0$;
- here, for SVRG, QuickeNing is robust to small values of κ !

Experiments: Influence of l



- $l = 100$ in all previous experiments;
- $l = 5$ seems to be a reasonable choice in many cases, especially for sparse problems.

Conclusions and perspectives

- A simple generic Quasi-Newton method for composite functions, with simple sub-problems, and complexity guarantees.
- We also have a variant for dual approaches.
- Does not solve the gap between theory and practice for L-BFGS.

Perspectives

- QuickeNing-BCD, QuickeNing-SAG, SAGA, SDCA...
- Other types of smoothing? \Rightarrow Links with recent Quasi-Newton methods applied to other envelopes [Stella et al., 2016].
- Simple line search improves slightly the performance.

Outer-loop convergence analysis

Lemma: approximate descent property

$$F(x_{k+1}) \leq f(z_k) \leq F(x_k) - \frac{1}{4\kappa} \|\nabla F(x_k)\|_2^2 + 2\varepsilon_k.$$

Then, ε_k should be smaller than $\frac{1}{4\kappa} \|\nabla F(x_k)\|_2^2$, and indeed

Outer-loop convergence analysis

Lemma: approximate descent property

$$F(x_{k+1}) \leq f(z_k) \leq F(x_k) - \frac{1}{4\kappa} \|\nabla F(x_k)\|_2^2 + 2\varepsilon_k.$$

Then, ε_k should be smaller than $\frac{1}{4\kappa} \|\nabla F(x_k)\|_2^2$, and indeed

Proposition: convergence with impractical ε_k and $\mu > 0$

If $\varepsilon_k \leq \frac{1}{16\kappa} \|\nabla F(x_k)\|_2^2$, define $\rho = \frac{\mu}{4(\mu+\kappa)}$, then

$$F(x_{k+1}) - F^* \leq f(z_k) - f^* \leq (1 - \rho)^{k+1} (f(x_0) - f^*).$$

Unfortunately, $\|\nabla F(x_k)\|$ is unknown.

Outer-loop convergence analysis

Lemma: approximate descent property

$$F(x_{k+1}) \leq f(z_k) \leq F(x_k) - \frac{1}{4\kappa} \|\nabla F(x_k)\|_2^2 + 2\varepsilon_k.$$

Then, ε_k should be smaller than $\frac{1}{4\kappa} \|\nabla F(x_k)\|_2^2$, and indeed

Proposition: convergence with impractical ε_k and $\mu > 0$

If $\varepsilon_k \leq \frac{1}{16\kappa} \|\nabla F(x_k)\|_2^2$, define $\rho = \frac{\mu}{4(\mu+\kappa)}$, then

$$F(x_{k+1}) - F^* \leq f(z_k) - f^* \leq (1 - \rho)^{k+1} (f(x_0) - f^*).$$

Unfortunately, $\|\nabla F(x_k)\|$ is unknown.

Lemma: convergence with adaptive ε_k and $\mu > 0$

If $\varepsilon_k \leq \frac{1}{36\kappa} \|g_k\|^2$, then $\varepsilon_k \leq \frac{1}{16} \|\nabla F(x_k)\|_2^2$.

This is strategy (b). g_k is known and easy to compute.

Inner-loop complexity analysis

Restart for L -smooth functions

For minimizing h , initialize the method \mathcal{M} with $w_0 = x$. Then,

$$h(w_0) - h^* \leq \frac{L + \kappa}{2\kappa^2} \|\nabla F(x)\|^2. \quad (1)$$

Proof.

We have the optimality condition $\nabla f(w^*) + \kappa(w^* - x) = 0$. As a result,

$$\begin{aligned} h(w_0) - h^* &= f(x) - \left(f(w^*) + \frac{\kappa}{2} \|w^* - x\|^2 \right) \\ &\leq f(w^*) + \langle \nabla f(w^*), x - w^* \rangle + \frac{L}{2} \|x - w^*\|^2 - \left(f(w^*) + \frac{\kappa}{2} \|w^* - x\|^2 \right) \\ &= \frac{L + \kappa}{2} \|w^* - x\|^2 = \frac{L + \kappa}{2\kappa^2} \|\nabla F(x)\|^2. \end{aligned}$$

References I

- A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009a.
- A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009b.
- J.V. Burke and Maijian Qian. On the superlinear convergence of the variable metric proximal point algorithm using Broyden and BFGS matrix secant updating. *Mathematical Programming*, 88(1):157–181, 2000.
- S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20:33–61, 1999.
- Xiaojun Chen and Masao Fukushima. Proximal quasi-Newton methods for nondifferentiable convex optimization. *Mathematical Programming*, 85(2): 313–334, 1999.
- P. L. Combettes and V. R. Wajs. Signal recovery by proximal forward-backward splitting. *SIAM Multiscale Modeling and Simulation*, 4(4):1168–1200, 2006.

References II

- I. Daubechies, M. Defrise, and C. De Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics*, 57(11):1413–1457, 2004.
- Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems (NIPS)*, 2014a.
- Aaron Defazio, Justin Domke, and Tibério S Caetano. Finito: A faster, permutable incremental gradient method for big data problems. In *Proceedings of the International Conferences on Machine Learning (ICML)*, 2014b.
- Olivier Devolder, F. Glineur, and Yurii Nesterov. First-order methods of smooth convex optimization with inexact oracle. *Mathematical Programming*, 146 (1-2):37–75, 2014.
- J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. 12:2121–2159, 2011.

References III

- Michael P Friedlander and Mark Schmidt. Hybrid deterministic-stochastic methods for data fitting. *SIAM Journal on Scientific Computing*, 34(3): A1380–A1405, 2012.
- Roy Frostig, Rong Ge, Sham M Kakade, and Aaron Sidford. Un-regularizing: approximate proximal point and faster stochastic algorithms for empirical risk minimization. In *Proceedings of the International Conferences on Machine Learning (ICML)*, 2015.
- Marc Fuentes, Jérôme Malick, and Claude Lemaréchal. Descentwise inexact proximal algorithms for smooth optimization. *Computational Optimization and Applications*, 53(3):755–769, 2012.
- Masao Fukushima and Liqun Qi. A globally and superlinearly convergent algorithm for nonsmooth convex minimization. *SIAM Journal on Optimization*, 6(4):1106–1120, 1996.
- O. Güler. New proximal point algorithms for convex minimization. *SIAM Journal on Optimization*, 2(4):649–664, 1992.

References IV

- Claude Lemaréchal and Claudia Sagastizábal. Practical aspects of the Moreau–Yosida regularization: Theoretical preliminaries. *SIAM Journal on Optimization*, 7(2):367–385, 1997.
- Hongzhou Lin, Julien Mairal, and Zaid Harchaoui. A universal catalyst for first-order optimization. In *Advances in Neural Information Processing Systems (NIPS)*, 2015.
- J. Mairal. Incremental majorization-minimization optimization with application to large-scale machine learning. *SIAM Journal on Optimization*, 25(2): 829–855, 2015.
- Bernard Martinet. Brève communication. régularisation d'inéquations variationnelles par approximations successives. 4(3):154–158, 1970.
- Robert Mifflin. A quasi-second-order proximal bundle algorithm. *Mathematical Programming*, 73(1):51–72, 1996.
- Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.

References V

- Y. Nesterov. *Introductory lectures on convex optimization: a basic course*. Kluwer Academic Publishers, 2004.
- Y. Nesterov. Gradient methods for minimizing composite objective function. *Mathematical Programming*, 140(1):125–161, 2013.
- Yurii Nesterov. A method for unconstrained convex minimization problem with the rate of convergence $o(1/k^2)$. In *Doklady an SSSR*, volume 269, pages 543–547, 1983.
- R. D. Nowak and M. A. T. Figueiredo. Fast wavelet-based image deconvolution using the EM algorithm. In *Conference Record of the Thirty-Fifth Asilomar Conference on Signals, Systems and Computers.*, 2001.
- Courtney Paquette, Hongzhou Lin, Dmitriy Drusvyatskiy, Julien Mairal, and Zaid Harchaoui. Catalyst acceleration for gradient-based non-convex optimization. *arXiv preprint arXiv:1703.10993*, 2017.
- R Tyrrell Rockafellar. Monotone operators and the proximal point algorithm. *SIAM Journal on Control and Optimization*, 14(5):877–898, 1976.
- Saverio Salzo and Silvia Villa. Inexact and accelerated proximal point algorithms. *Journal of Convex Analysis*, 19(4):1167–1192, 2012.

References VI

- M. Schmidt, N. Le Roux, and F. Bach. Convergence rates of inexact proximal-gradient methods for convex optimization. 2011.
- M. Schmidt, N. Le Roux, and F. Bach. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 160(1):83–112, 2017.
- S. Shalev-Shwartz and T. Zhang. Proximal stochastic dual coordinate ascent. *arXiv:1211.2717*, 2012.
- S. Shalev-Shwartz and T. Zhang. Accelerated proximal stochastic dual coordinate ascent for regularized loss minimization. *Mathematical Programming*, 155(1):105–145, 2016.
- Lorenzo Stella, Andreas Themelis, and Panagiotis Patrinos. Forward-backward quasi-newton methods for nonsmooth optimization problems. *arXiv preprint arXiv:1604.08096*, 2016.
- R. Tibshirani. Regression shrinkage and selection via the Lasso. 58(1):267–288, 1996.

References VII

- S. J. Wright, R. D. Nowak, and M. A. T. Figueiredo. Sparse reconstruction by separable approximation. *IEEE Transactions on Signal Processing*, 57(7): 2479–2493, 2009.
- L. Xiao and T. Zhang. A proximal stochastic gradient method with progressive variance reduction. *SIAM Journal on Optimization*, 24(4):2057–2075, 2014.
- Y. Zhang and L. Xiao. Stochastic primal-dual coordinate method for regularized empirical risk minimization. In *Proceedings of the International Conferences on Machine Learning (ICML)*, 2015.