

# An Inexact Variable Metric PPA for Generic Quasi-Newton Acceleration

Hongzhou Lin<sup>1,3</sup>, Julien Mairal<sup>1</sup>, Zaid Harchaoui<sup>2</sup>

<sup>1</sup>Inria, Grenoble

<sup>2</sup>University of Washington

<sup>3</sup>MIT

ISMP

Bordeaux, 2018



# 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  $\psi$  is a convex regularization penalty but not necessarily differentiable.

## Motivation

Our goal is to accelerate existing algorithms

- with Nesterov's principles (previous work Catalyst);
- with **Quasi-Newton** heuristics (this work);

# 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.
- $\psi$  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**.

## 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  $\mu$ -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).

## 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  $\nabla f_i$  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  $\nabla f_i$  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	$O\left(\max\left(n, \frac{\bar{L}}{\mu}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$
Catalyst-SVRG, SAG, SAGA, ...	$\tilde{O}\left(\max\left(n, \sqrt{\frac{n\bar{L}}{\mu}}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$

## 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	$O\left(\max\left(n, \frac{\bar{L}}{\mu}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$
Catalyst-SVRG, SAG, SAGA, ...	$\tilde{O}\left(\max\left(n, \sqrt{\frac{n\bar{L}}{\mu}}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$

### 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	$O\left(\max\left(n, \frac{\bar{L}}{\mu}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$
Catalyst-SVRG, SAG, SAGA, ...	$\tilde{O}\left(\max\left(n, \sqrt{\frac{n\bar{L}}{\mu}}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$

### 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	$O\left(\max\left(n, \frac{\bar{L}}{\mu}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$
Catalyst-SVRG, SAG, SAGA, ...	$\tilde{O}\left(\max\left(n, \sqrt{\frac{n\bar{L}}{\mu}}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$

### 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	$O\left(\max\left(n, \frac{\bar{L}}{\mu}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$
Catalyst-SVRG, SAG, SAGA, ...	$\tilde{O}\left(\max\left(n, \sqrt{\frac{n\bar{L}}{\mu}}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$

### 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**  $\psi$ .

## Catalyst and QNing: An 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,  $\nabla F$  is Lipschitz continuous with parameter  $L_F = \kappa$ .

- If  $f$  is  $\mu$ -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,  $\nabla F$  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  $\beta_{k+1}$  is a Nesterov-like extrapolation parameter. We may now rewrite the update using the value of  $\nabla F$ , 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  $\beta_{k+1}$  is a Nesterov-like extrapolation parameter. We may now rewrite the update using the value of  $\nabla F$ , 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]

# 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 **not better than the gradient descent method.**
- proximal variants typically requires solving many times

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

- no guarantee of approximating the Hessian.

## Back to the old idea

### 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  $\mathcal{M}$ ;
- the approach is useful if  $\mathcal{M}$  is **able to adapt to the problem structure** (finite sum, composite regularization).

## 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 QNing

---

**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 QNing

**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  $\varepsilon_k$ -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 for $\mu$ -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  $\varepsilon_k$ -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  $\mu$ -strongly convex problems.
- $\tilde{O}\left(\frac{\kappa R^2}{\tau_{\mathcal{M}} \varepsilon}\right)$  for convex problems.

## Global Complexity and choice of $\kappa$

### 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 $\kappa$

### 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

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

**Qning 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 $\kappa$

### 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  $\kappa$ ?

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

## Experiments: formulations

- $\ell_2$ -regularized Logistic Regression:

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

- $\ell_1$ -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  $\mu = 1/(100n)$ ;
- $\lambda$  for the Lasso leads to about 10% non-zero coefficients.

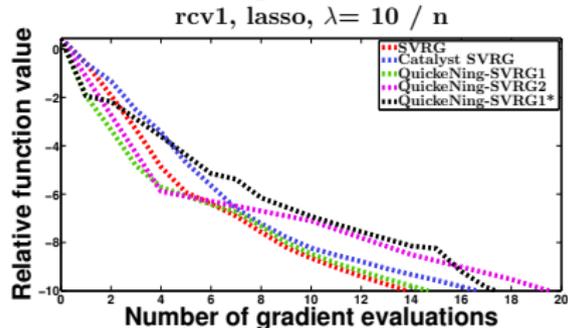
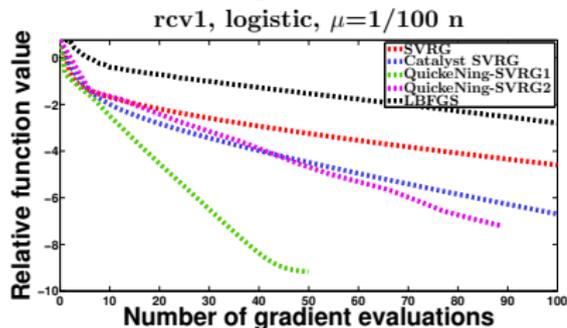
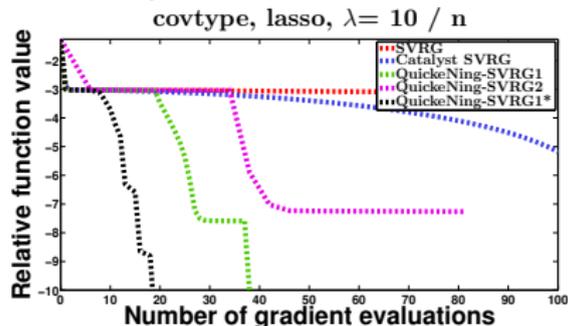
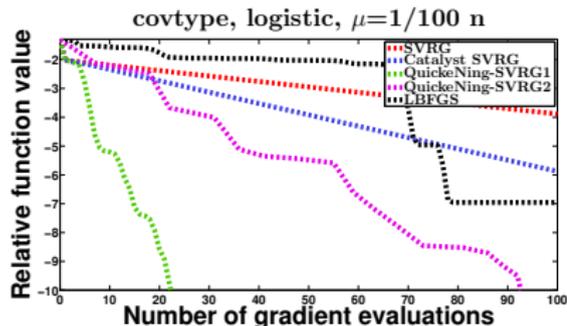
## Experiments: QNing-SVRG

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.
- **QNing-SVRG1**: QNing with aggressive strategy (c): one pass over the data in the inner loop.
- **QNing-SVRG2**: strategy (b), compatible with theory.

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

# Experiments: QNing-SVRG (log scale)



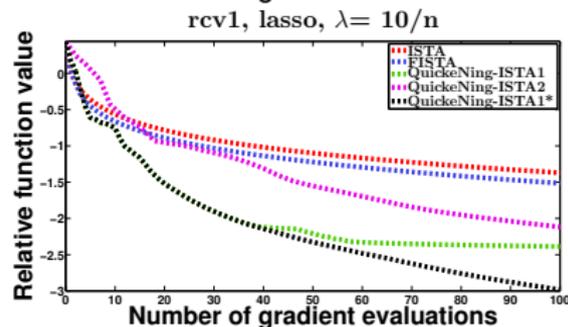
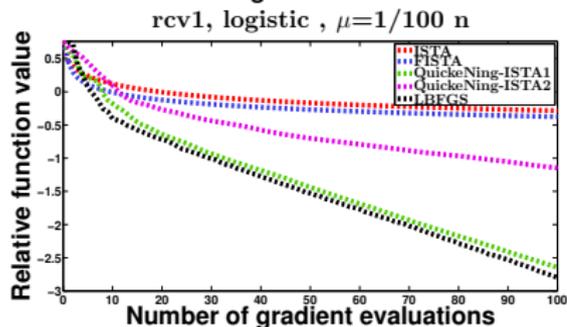
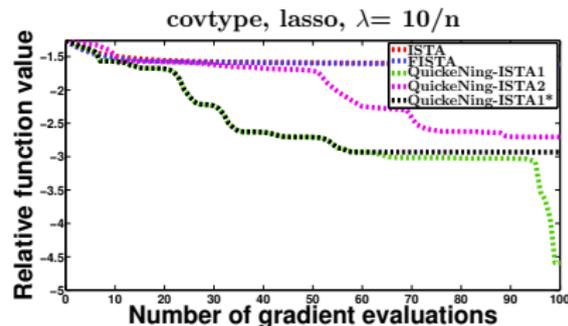
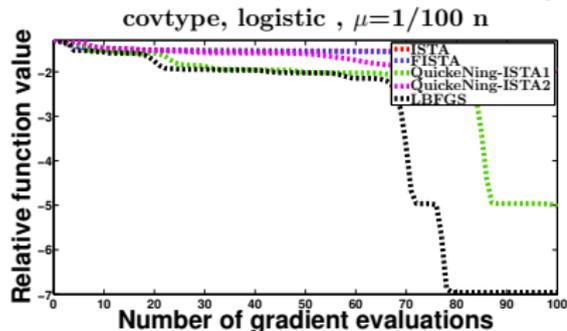
- $\text{QNing-SVRG1} \geq \text{SVRG}, \text{QNing-SVRG2}$ ;
- $\text{QNing-SVRG2} \geq \text{SVRG}$ ;
- $\text{QNing-SVRG1} \geq \text{Catalyst-SVRG}$  in 10/12 cases.

## Experiments: QNing-ISTA

We consider the methods

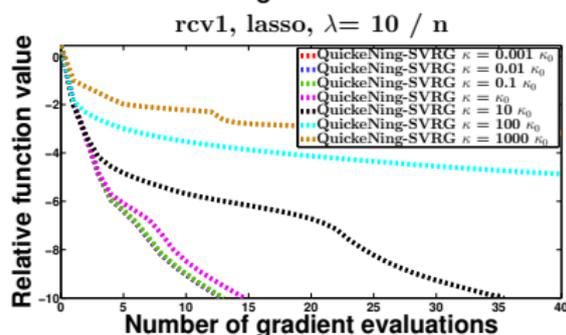
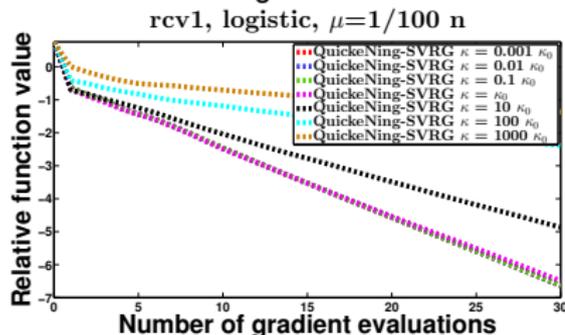
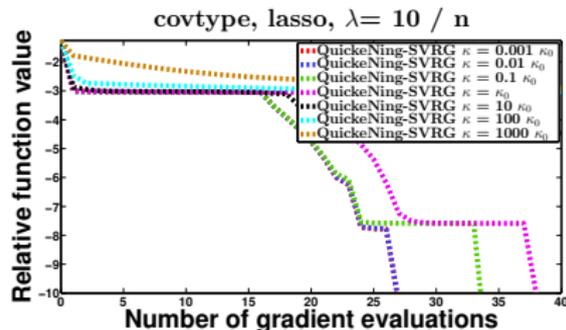
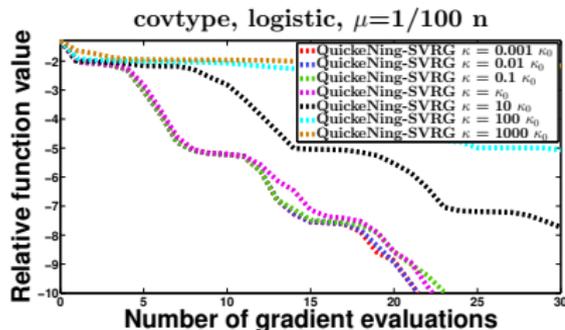
- **ISTA**: the proximal gradient descent method with line search.
- **FISTA**: the accelerated ISTA of Beck and Teboulle [2009].
- **L-BFGS** (for smooth objectives): Mark Schmidt's implementation.
- **QNing-ISTA1**: QNing with aggressive strategy (c): one pass over the data in the inner loop.
- **QNing-ISTA2**: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

# Experiments: QNing-ISTA (log scale)



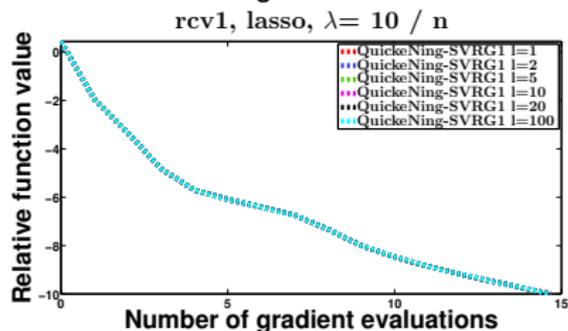
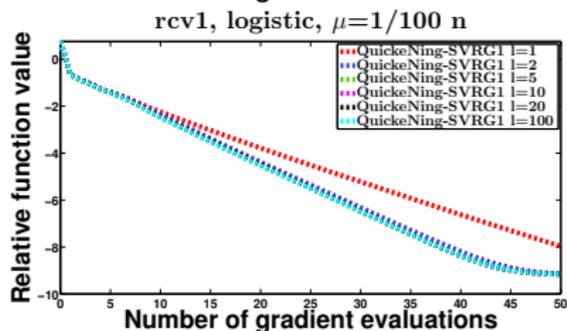
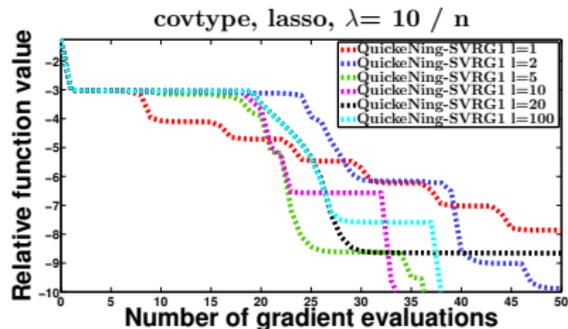
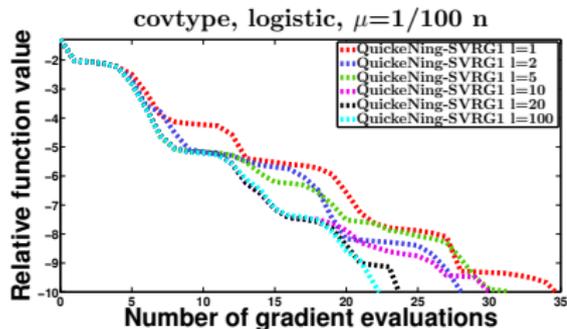
- L-BFGS (for smooth  $f$ ) is slightly better than QNing-ISTA1;
- QNing-ISTA  $\geq$  or  $\gg$  FISTA in 11/12 cases.
- QNing-ISTA1  $\geq$  QNing-ISTA2.

# Experiments: Influence of $\kappa$



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

# 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

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

**Proposition: convergence with impractical  $\varepsilon_k$  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,  $\varepsilon_k$  should be smaller than  $\frac{1}{4\kappa} \|\nabla F(x_k)\|_2^2$ , and indeed

**Proposition: convergence with impractical  $\varepsilon_k$  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  $\varepsilon_k$  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, 2009.
- 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.
- Xiaojun Chen and Masao Fukushima. Proximal quasi-Newton methods for nondifferentiable convex optimization. *Mathematical Programming*, 85(2): 313–334, 1999.
- 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.

## References II

- 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.
- 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.

## References III

- 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.
- 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.

## References IV

- Robert Mifflin. A quasi-second-order proximal bundle algorithm. *Mathematical Programming*, 73(1):51–72, 1996.
- 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.
- 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.

## References V

- Lorenzo Stella, Andreas Themelis, and Panagiotis Patrinos. Forward-backward quasi-newton methods for nonsmooth optimization problems. *arXiv preprint arXiv:1604.08096*, 2016.
- 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.