

# Acceleration methods in Optimization

Aude Rondepierre



Institut de Mathématiques de Toulouse, INSA de Toulouse & LAAS-CNRS

CIMI-ANITI School on Optimisation - September 2th, 2021

# The framework

$$\min_{x \in \mathbb{R}^N} F(x)$$

where  $F : \mathbb{R}^N \rightarrow \mathbb{R}$  is a differentiable convex function admitting at least one minimizer  $x^*$ .

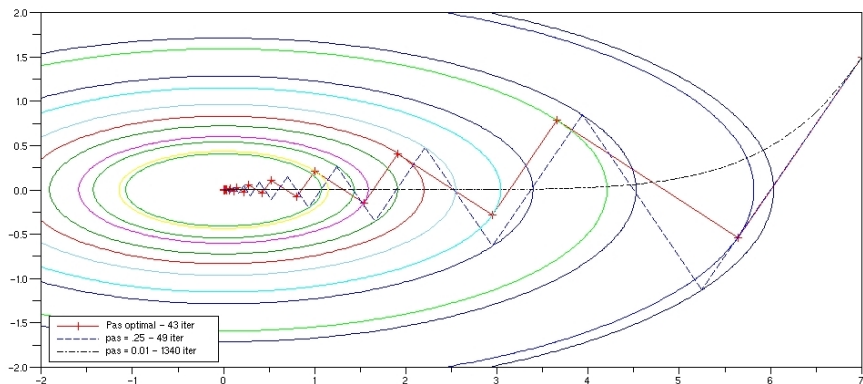
## The gradient descent (GD) method

$$\begin{array}{l|l} x_0 & \in \mathbb{R}^N \\ x_{k+1} & = x_k - s \nabla F(x_k), \quad s > 0. \end{array}$$

- A very simple algorithm, does not require second order derivative.
- Each iteration is of the order of  $N$  operations.

## A quadratic example

$$\min_{(x,y) \in \mathbb{R}^2} \mathbf{F}(x,y) = \frac{1}{2}x^2 + \frac{7}{2}y^2, \quad d_k = -\nabla f(X_k) = \begin{pmatrix} -x_k \\ -7y_k \end{pmatrix}.$$



Gradient method is often slow ; the convergence is very dependent on scaling.

# Methods with improved convergence rate

## Methods with improved convergence

- Conjugate gradient method
- Accelerated gradient method
- Quasi-Newton methods



# Methods with improved convergence rate

## Methods with improved convergence

- Conjugate gradient method
- Accelerated gradient method
- Quasi-Newton methods

# Methods with improved convergence rate

## Methods with improved convergence

- Conjugate gradient method
- Accelerated gradient method
- Quasi-Newton methods

## Natural extensions to composite optimization $F(x) = f(x) + g(x)$ where

- $f$  is a convex differentiable function with a  $L$ -Lipschitz gradient
- $g$  is a convex lsc (possibly nonsmooth but quite simple) function.

↪ Motivation: application to least square problems, LASSO:

$$\min_{x \in \mathbb{R}^N} \frac{1}{2} \|Ax - b\|^2 + \|x\|_1.$$

Applications in Image and Signal processing, machine learning,...

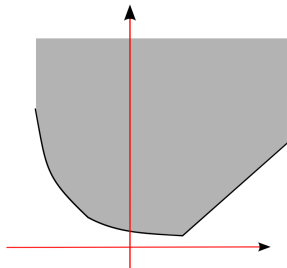
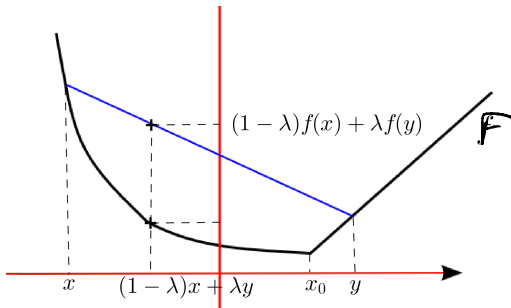
- 1 Preliminaries: local geometry of convex functions
- 2 Gradient descent methods
- 3 Accelerated gradient methods
  - The Heavy ball method
  - The Nesterov's accelerated gradient method
  - Natural extension to composite optimization
  - Some numerical experiments
- 4 Improving the state of the art results with the help of the geometry
- 5 Conclusion

# Geometry of convex functions

## Convexity (1/2)

A function  $F : \mathbb{R}^n \rightarrow \mathbb{R}$  is said convex if:

$$\forall (x, y) \in \text{dom}(F)^2, \forall \lambda \in [0, 1], F((1 - \lambda)x + \lambda y) \leq (1 - \lambda)F(x) + \lambda F(y).$$



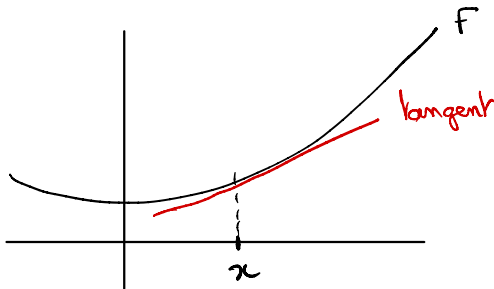
$F$  convex iff its epigraph is convex

# Geometry of convex functions

## Convexity (2/2)

Let  $F : \mathbb{R}^n \rightarrow \mathbb{R}$  be a differentiable function.  $F$  is convex if:

$$\forall (x, y) \in \mathbb{R}^n \times \mathbb{R}^n, f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle.$$



⇒ Monotonicity of the gradient:

$$\forall (x, y) \in \mathbb{R}^N \times \mathbb{R}^N, \langle \nabla F(y) - \nabla F(x), y - x \rangle \geq 0.$$

# Geometry of convex functions

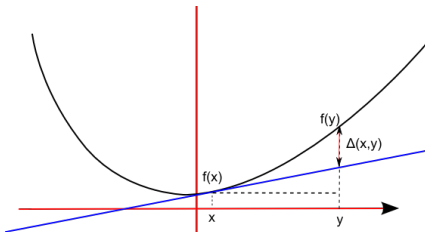
## Functions having a $L$ -Lipschitz gradient / $L$ -smooth functions (1/3)

Let  $F : \mathbb{R}^n \rightarrow \mathbb{R}$  be a continuously differentiable function and  $L > 0$ . The function  $F$  has a  $L$ -Lipschitz gradient iff:

$$\forall (x, y) \in \mathbb{R}^N \times \mathbb{R}^N, \|\nabla F(x) - \nabla F(y)\| \leq L\|x - y\|.$$

❶ Quadratic upper bound: For all  $(x, y) \in \mathbb{R}^N \times \mathbb{R}^N$ , we have:

$$F(y) \leq \underbrace{F(x) + \langle \nabla F(x), y - x \rangle}_{\text{linear approximation}} + \underbrace{\frac{L}{2}\|y - x\|^2}_{=\Delta(x,y)}$$



## Geometry of convex functions

### Proof of the quadratic upper bound (2/3)

$$\varphi(t) = F(x + t(y-x)) \quad \nabla \varphi(t) = \langle \nabla F(x + t(y-x)), y-x \rangle$$

$$\begin{aligned} F(y) - F(x) &= \varphi(1) - \varphi(0) = \int_0^1 \varphi'(t) dt \\ &= \int_0^1 \langle \nabla F(x + t(y-x)), y-x \rangle dt \end{aligned}$$

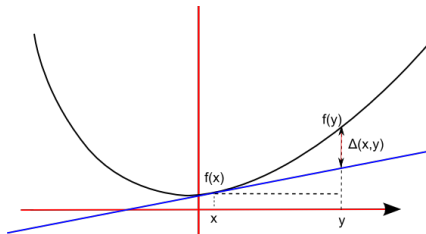
$$\begin{aligned} F(y) - F(x) - \langle \nabla F(x), y-x \rangle &= \int_0^1 \langle \nabla F(x + t(y-x)) - \nabla F(x), y-x \rangle dt \\ &\leq \int_0^1 \overbrace{\|\nabla F(x + t(y-x)) - \nabla F(x)\|}^{\leq L \|x + t(y-x) - x\|} \|y-x\| dt \\ &\leq L \int_0^1 t \|y-x\|^2 dt = \frac{L}{2} \|y-x\|^2 \end{aligned}$$

# Geometry of convex functions

## Strong convexity (1/2)

$F : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\mu$ -strongly convex i.e. that there exists  $\mu > 0$  such that:

$$\forall (x, y) \in \mathbb{R}^n \times \mathbb{R}^n, F(y) \geq F(x) + \langle \nabla F(x), y - x \rangle + \frac{\mu}{2} \|y - x\|^2.$$





# Geometry of convex functions

## Strong convexity (2/2)

This class of functions satisfies a **global quadratic growth condition**: for any minimizer  $x^*$  we have:

$$\forall x \in \mathbb{R}^n, F(x) - F(x^*) \geq \frac{\mu}{2} \|x - x^*\|^2.$$

and:

$$\forall x \in \mathbb{R}^n, \|\nabla F(x)\|^2 \geq 2\mu(F(x) - F(x^*)).$$

Logarithmic  
property with  
 $\Theta = \frac{1}{2}$

- 1 Preliminaries: local geometry of convex functions
- 2 Gradient descent methods
- 3 Accelerated gradient methods
  - The Heavy ball method
  - The Nesterov's accelerated gradient method
  - Natural extension to composite optimization
  - Some numerical experiments
- 4 Improving the state of the art results with the help of the geometry
- 5 Conclusion

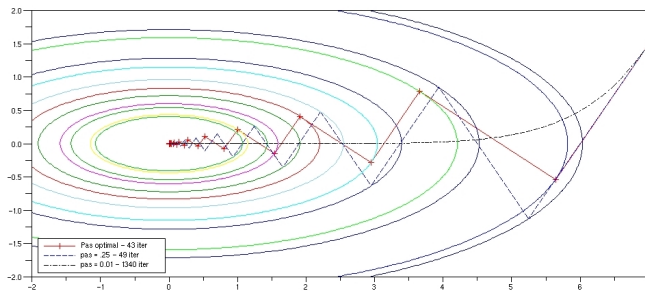
# The gradient descent method

## Algorithm and basic properties

Let  $F : \mathbb{R}^N \rightarrow \mathbb{R}$  be a continuously differentiable function having a  $L$ -Lipschitz gradient and admitting at least one minimizer.

$$\begin{aligned}x_0 &\in \mathbb{R}^N \\ x_{k+1} &= x_k - s_k \nabla F(x_k), \quad s_k > 0.\end{aligned}$$

with a fixed step  $s_k := s > 0$  or backtracking linesearch.



Step	0.325	0.25	0.125	0.05	0.01
Nb of it.	DV	49	101	263	1340

# The gradient descent method

## Properties

Properties:

- ① GD is a descent method when  $s_k < \frac{2}{L}$  for all  $k \in \mathbb{N}$ , :

$$\forall k \in \mathbb{N}, F(x_{k+1}) - F(x_k) \leq s_k \left( \frac{L}{2} s_k - 1 \right) \|\nabla F(x_k)\|^2 \leq 0.$$

- ② Assume that  $F$  is additionally convex and  $s_k \leq \frac{1}{L}$ . The distance to the optimal set decrease. Let  $x^* \in \arg \min(F)$ .

$$\forall k \in \mathbb{N}, \|x_{k+1} - x^*\| \leq \|x_k - x^*\|.$$

# The gradient descent method

Proof (1/2)

$$x_{k+1} = x_k - s_k \nabla F(x_k)$$

①

$$F(x_{k+1}) \leq F(x_k) + \underbrace{(\nabla F(x_k), x_{k+1} - x_k)}_{-s_k \|\nabla F(x_k)\|^2} + \frac{L}{2} \|x_{k+1} - x_k\|^2$$

$$\Rightarrow F(x_{k+1}) \leq F(x_k) - s_k \|\nabla F(x_k)\|^2 + \frac{L}{2} s_k^2 \|\nabla F(x_k)\|^2$$

$$\leq F(x_k) - s_k \underbrace{\left(1 - \frac{L}{2} s_k\right)}_{\geq 0} \|\nabla F(x_k)\|^2$$

$$\leq F(x_k) \quad \text{if } s_k < \frac{2}{L}$$

②  $\forall k, F(x_{k+1}) - F(x_k) \leq \frac{L}{2} (\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2)$

Using a Series argument

$$\begin{aligned} &= \sum_{k=0}^{N-1} (F(x_{k+1}) - F(x_k)) \leq \frac{L}{2} (\|x_0 - x^*\|^2 - \|x_N - x^*\|^2) \\ &\leq \frac{L}{2} \|x_0 - x^*\|^2 \end{aligned}$$

# The gradient descent method

## Proof (2/2)

$$\Rightarrow \forall m, m(F(x_m) - F^*) \leq \frac{L}{2} \|x_m - x^*\|^2$$

$$\Rightarrow \forall m, F(x_m) - F^* \leq \frac{L \|x_m - x^*\|^2}{2m}$$

# The dynamical system intuition

Link with the ODEs - A guideline to study optimization algorithms

## General methodology to analyze optimization algorithms

- Interpreting the optimization algorithm as a discretization of a given ODE:

Gradient descent iteration:  $\frac{x_{n+1} - x_n}{\Delta t} + \nabla F(x_n) = 0$

Associated ODE:  $\dot{x}(t) + \nabla F(x(t)) = 0.$

- Analysis of ODEs using a Lyapunov approach:

$$\mathcal{E}(t) = F(x(t)) - F^*$$

$$\mathcal{E}(t) = t(F(x(t)) - F^*) + \frac{1}{2}\|x(t) - x^*\|^2.$$

- Building a sequence of discrete Lyapunov energies adapted to the optimization scheme to get the same decay rates

## Gradient descent for strongly convex functions

A Lyapunov analysis of the ODE  $\dot{x}(t) + \nabla F(x(t)) = 0$  (1/3)

Let:

$$\mathcal{E}(t) = F(x(t)) - F^*.$$

①  $\mathcal{E}$  is a Lyapunov energy (i.e. non increasing along the trajectories  $x(t)$ ):

$$\begin{aligned}\mathcal{E}'(t) &= \langle \nabla F(x(t)), \dot{x}(t) \rangle = - \|\nabla F(x(t))\|^2 \\ &\leq 0\end{aligned}$$

$$\Rightarrow \forall t \geq t_0, \quad \mathcal{E}(t) \leq \mathcal{E}(t_0)$$

$$\Rightarrow \forall t \geq t_0, \quad F(x(t)) - F^* \leq F(x_0) - F^*$$



## Gradient descent for strongly convex functions

A Lyapunov analysis of the ODE  $\dot{x}(t) + \nabla F(x(t)) = 0$  (2/3)

- ② Assume now that  $F$  is additionally  $\mu$ -strongly convex. Remember that:

$$\forall y \in \mathbb{R}^N, \|\nabla F(y)\|^2 \geq 2\mu(F(y) - F^*),$$

$$\begin{aligned} \mathcal{E}'(t) &= -\|\nabla F(x(t))\|^2 \leq -2\mu(F(x(t)) - F^*) \\ &\leq -2\mu \mathcal{E}(t) \end{aligned}$$

$$\Rightarrow \forall t \geq t_0, \quad \mathcal{E}(t) \leq e^{-2\mu t} \mathcal{E}(t_0)$$

$$\Rightarrow \forall t \geq t_0, \quad F(x(t)) - F^* \leq e^{-2\mu t} (F(x_0) - F^*)$$

## Gradient descent for strongly convex functions

A Lyapunov analysis of the ODE  $\dot{x}(t) + \nabla F(x(t)) = 0$ ,  $x(0) = x_0$

③ Assume that  $F$  is only convex. Let  $x^* \in \arg \min(F)$ .

$$\mathcal{E}(t) = t(F(x(t)) - F^*) + \frac{1}{2} \|x(t) - x^*\|^2.$$

$\mathcal{E}$  is also Lyapunov energy:  $-\nabla F(x(t))$

$$\begin{aligned} \mathcal{E}'(t) &= t \langle \nabla F(x(t)), \dot{x}(t) \rangle + F(x(t)) - F^* \\ &\quad + \langle x(t) - x^*, \dot{x}(t) \rangle \end{aligned}$$

$$= -t \|\nabla F(x(t))\|^2 + \underbrace{F(x(t)) - F^* \langle x(t) - x^*, \nabla F(x(t)) \rangle}_{\leq 0}$$

by convexity.

$$\leq -t \|\nabla F(x(t))\|^2 \leq 0$$

$$\begin{aligned} \Rightarrow \forall t \geq t_0, \quad t(F(x(t)) - F^*) &\leq \mathcal{E}(t) \leq \mathcal{E}(t_0) \\ &\hookrightarrow \text{cv in } \mathcal{O}\left(\frac{1}{t}\right) \end{aligned}$$

# Gradient descent for strongly convex functions

## From the continuous to the discrete (1/3)

$$\mathcal{E}_n = F(x_n) - F^* \quad \text{with:} \quad x_{n+1} = x_n - s \nabla F(x_n).$$

①  $\mathcal{E}_n$  is a discretization of the Lyapunov energy  $\mathcal{E}(t)$ . We have:

$$\begin{aligned} \mathcal{E}_{n+1} - \mathcal{E}_n &= F(x_{n+1}) - F(x_n) \leq \langle \nabla F(x_n), x_{n+1} - x_n \rangle + \frac{L}{2} \|x_{n+1} - x_n\|^2 \\ &\leq -s \left(1 - \frac{L}{2}s\right) \|\nabla F(x_n)\|^2 \end{aligned}$$

If the step  $s$  satisfies:

$$s < \frac{2}{L}$$

then the GD is a descent algorithm ( $\forall n, F(x_{n+1}) < F(x_n)$ ) and the values  $F(x_n) - F^*$  remain bounded.

# Gradient descent for strongly convex functions

## From the continuous to the discrete (1/3)

$$\mathcal{E}_n = F(x_n) - F^* \quad \text{with:} \quad x_{n+1} = x_n - s \nabla F(x_n).$$

① Assume now that  $F$  is additionally  $\mu$ -strongly convex and  $h < \frac{2}{L}$ . We have:

$$\forall n, \|\nabla F(x_n)\|^2 \geq 2\mu(F(x_n) - F^*) = 2\mu\mathcal{E}_n,$$

and

$$\mathcal{E}_{n+1} - \mathcal{E}_n \leq -s \left(1 - \frac{L}{2}s\right) \|\nabla F(x_n)\|^2.$$

Hence:

$$\mathcal{E}_{n+1} - \mathcal{E}_n \leq -2\mu s \left(1 - \frac{L}{2}s\right) \mathcal{E}_n$$

For example if  $s = \frac{1}{L}$  we get:

$$\forall n, \mathcal{E}_{n+1} - \mathcal{E}_n \leq -\kappa \mathcal{E}_n \Rightarrow \mathcal{E}_n \leq (1 - \kappa)^n \mathcal{E}_0$$

hence:

$$F(x_n) - F^* \leq (F(x_0) - F^*)(1 - \kappa)^n.$$

# The gradient descent method

## Convergence rates for convex and strongly convex functions

Let  $F : \mathbb{R}^N \rightarrow \mathbb{R}$  be a convex function having a  $L$ -Lipschitz gradient. Choosing  $s = \frac{1}{L}$  in (GD), we get:

### Theorem

$$\forall k \in \mathbb{N}, F(x_k) - F(x^*) \leq \frac{L \|x_0 - x^*\|^2}{2k}.$$

*The number of iterations to reach  $F(x_k) - F^* \leq \varepsilon$  is in  $\mathcal{O}(\frac{1}{\varepsilon})$ .*

Additionally assume  $F$   $\mu$ -strongly convex for some  $\mu > 0$ . Then:

### Theorem

$$\forall n \in \mathbb{N}, F(x_n) - F^* \leq (1 - \kappa)^n (F(x_0) - F^*) \quad \text{where } \kappa = \frac{\mu}{L}.$$

*The number of iterations to reach  $F(x_n) - F^* \leq \varepsilon$  is in  $\mathcal{O}(\log(\frac{1}{\varepsilon}))$ .*

**Rmq:** the optimal convergence factor is:  $F(x_n) - F^* \leq \left(\frac{L-\mu}{L+\mu}\right)^n (F(x_0) - F^*)$   
attained for  $s = \frac{2}{L+\mu}$ .

# The gradient method

## Limits on the convergence rates of first order methods

**First order method:** any iterative algorithm that selects  $x_{k+1}$  in the set

$$x_0 + \text{span} \{ \nabla F(x_0), \dots, \nabla F(x_k) \}.$$

### Theorem (Nemirovski Yudin 1983, Nesterov 2003)

Let  $k \leq \frac{N-1}{2}$  and  $L > 0$ . There exists a convex function  $F$  having a  $L$ -Lipschitz gradient over  $\mathbb{R}^N$  such that for any first order method

$$F(x_k) - F^* \geq \frac{3L\|x_0 - x^*\|^2}{32(k+1)^2}.$$

- ↪ Suggests that the rate  $\frac{1}{k}$  for GD is not optimal !
- ↪ We will see that recent accelerated gradient methods have a  $\frac{1}{k^2}$  convergence rate.

- 1 Preliminaries: local geometry of convex functions
- 2 Gradient descent methods
- 3 Accelerated gradient methods
  - The Heavy ball method
  - The Nesterov's accelerated gradient method
  - Natural extension to composite optimization
  - Some numerical experiments
- 4 Improving the state of the art results with the help of the geometry
- 5 Conclusion

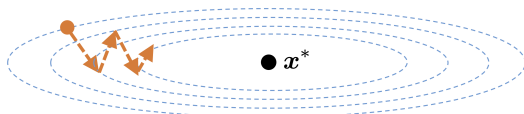
# The Heavy Ball method

A first inertial method (Polyak 1964)

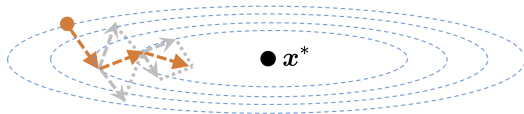
## The Heavy ball method

$$\begin{aligned}y_k &= x_k + a(x_k - x_{k-1}) \\ x_{k+1} &= y_k - s \nabla F(x_k)\end{aligned}, \quad \alpha \in [0, 1], \quad s > 0.$$

where  $a \in [0, 1]$  is an *fixed* inertial coefficient added to mitigate zigzagging.



gradient descent



heavy-ball method



# The Heavy Ball method

## The dynamical system intuition

Let us consider:

$$\ddot{x}(t) + \alpha \dot{x}(t) + \nabla F(x(t)) = 0.$$

- Describe the motion of a body (a heavy ball) in a potential field  $F$  subject to a friction proportional to its velocity.
- Natural intuition: the body reaches a minimum of the potential  $F$ .

# The Heavy Ball method

## Link between the continuous ODE and the discrete scheme

The HB algorithm:

$$\begin{aligned} y_k &= x_k + a(x_k - x_{k-1}) \\ x_{k+1} &= y_k - s \nabla F(x_k) \end{aligned}, \quad \alpha \in [0, 1], \quad s > 0.$$

can be seen as a discretization of the second order ODE:

$$\ddot{x}(t) + \alpha \dot{x}(t) + \nabla F(x(t)) = 0$$

where:

$$s = h^2, \quad a = 1 - \alpha h.$$

$$x_{k+1} = x_k + a(x_k - x_{k-1}) \rightarrow \nabla F(x_k) = 0$$

$$\underbrace{\frac{1}{h^2}(x_{k+1} - 2x_k + x_{k-1}))}_{\ddot{x}(t)} + \frac{(1-a)}{h} \underbrace{(x_k - x_{k-1}))}_{\dot{x}(t)} + \frac{s}{h^2} \underbrace{\nabla F(x_k)}_{\nabla F(x(t))} = 0$$

# The Heavy Ball method

## Convergence results - In the continuous case

$$\ddot{x}(t) + \alpha \dot{x}(t) + \nabla F(x(t)) = 0$$

### Theorem (Global convergence - Polyak 1964)

Let  $F : \mathbb{R}^n \rightarrow \mathbb{R}$  be a  $\mu$ -strongly convex function of class  $C^2$  and having a  $L$ -Lipschitz continuous gradient.

- If  $\alpha \leq 2\sqrt{\mu}$  then:

$$F(x(t)) - F^* = \mathcal{O}(e^{-\alpha t}).$$

- If  $\alpha > 2\sqrt{\mu}$  then:

$$F(x(t)) - F^* = \mathcal{O}\left(e^{-(\alpha - \sqrt{\alpha^2 - 4\mu})t}\right).$$

$$\bullet F(x) = \frac{\mu}{2} \|x\|^2 \quad \leadsto \quad \ddot{x}(t) + \alpha \dot{x}(t) + \mu x(t) = 0$$

# The Heavy Ball method

## Convergence results - In the discrete case

$$\begin{aligned} y_k &= x_k + a(x_k - x_{k-1}) \\ x_{k+1} &= y_k - s \nabla F(x_k) \end{aligned} \quad \Bigg/$$

with (Polyak's choice):

$$a = \left( \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \right)^2, \quad s = \left( \frac{2}{\sqrt{L} + \sqrt{\mu}} \right)^2.$$

### Theorem (Global convergence - [Polyak 1964])

Let  $F : \mathbb{R}^n \rightarrow \mathbb{R}$  be a  $\mu$ -strongly convex function of class  $C^2$  and having a  $L$ -Lipschitz continuous gradient. If  $s < \frac{2}{L}$  then:

$$F(x_k) - F^* \leq \left( \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \right)^k (F(x_0) - F^*).$$

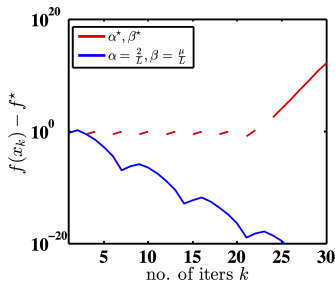
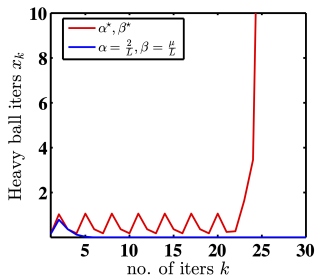
# The Heavy Ball method

## Convergence results - Without the $C^2$ assumption

**Counter example** [Ghadimi et al. 2015] Let  $F$  be a  $C^1$   $\mu$ -strongly convex and  $L$ -smooth function (with  $\mu = 5$  and  $L = 50$ ) such that:

$$\nabla F(x) = \begin{cases} 50x + 45 & \text{if } x < -1 \\ 5x & \text{if } -1 \leq x < 0 \\ 50x & \text{if } x \geq 0 \end{cases}$$

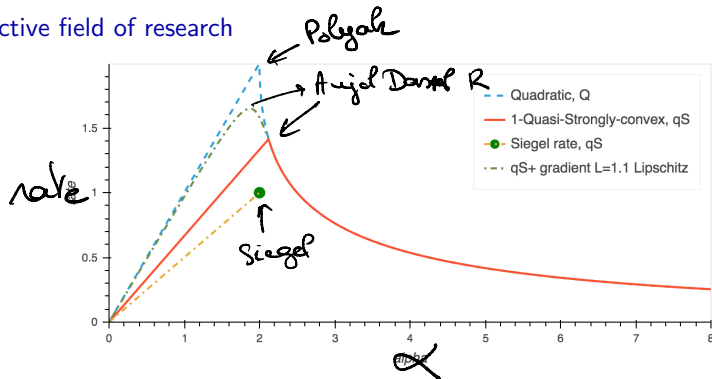
anf  $F$  is not of class  $C^2$ .



# The Heavy Ball method

## Convergence results - Without the $C^2$ assumption

An active field of research



- Nesterov's variant (2013):  $F(x_n) - F^* = \mathcal{O}((1 - \sqrt{\kappa})^n)$ .
- Changing the step and the inertia, [Ghadimi et al. 2015] prove the linear cv for  $C^1$  strongly convex functions having a Lipschitz continuous gradient.
- For strongly convex functions of class  $C^1$  having a  $L$ -Lipschitz gradient [Siegel 2019]: when  $\alpha = 2\sqrt{\mu}$ ,  $F(x(t)) - F^* = \mathcal{O}(e^{\sqrt{\mu}t})$ .

# The Nesterov's accelerated gradient method

## Outline

- 1 An inertial method - Historical choice of Nesterov
- 2 The dynamical system intuition
- 3 Proof for the convex case

# The Nesterov's accelerated gradient method

## The historical scheme

Historically Nesterov proposes in 1983 the following inertial method:

$$\begin{aligned}y_k &= x_k + \frac{t_k - 1}{t_{k+1}}(x_k - x_{k-1}) \\x_{k+1} &= y_k - s \nabla F(y_k)\end{aligned}$$

where the sequence  $(t_k)_{k \in \mathbb{N}}$  is defined by:  $t_1 = 1$  and:

$$t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}.$$

With such an algorithm Nesterov proves that:

$$\forall k \in \mathbb{N}, F(x_k) - F^* \leq \frac{2\|x_0 - x^*\|^2}{sk^2}$$

but he did not prove the convergence of the iterates  $(x_k)_{k \in \mathbb{N}}$ .



# The Nesterov's accelerated gradient method

A simplification of the first scheme

## Nesterov inertial scheme

$$\begin{cases} y_n &= x_n + \frac{n}{n + \alpha} (x_n - x_{n-1}) \\ x_{n+1} &= y_n - h \nabla F(y_n). \end{cases}$$

- Initially, Nesterov (1984) proposes  $\alpha = 3$ .
- Adapted by Beck and Teboulle to composite nonmooth functions (FISTA)

## Link between the ODE and the optimization scheme (1/2)

### Discretization of an ODE, Su Boyd and Candès (15)

The scheme defined by

$$x_{n+1} = y_n - \frac{1}{n+\alpha} \nabla F(y_n) \text{ with } y_n = x_n + \frac{n}{n+\alpha} (x_n - x_{n-1})$$

can be seen as a semi-implicit discretization of a solution of

$$\ddot{x}(t) + \frac{\alpha}{t} \dot{x}(t) + \nabla F(x(t)) = 0 \quad (\text{ODE})$$

With  $\dot{x}(t_0) = 0$ . Move of a solid in a potential field with a vanishing viscosity  $\frac{\alpha}{t}$ .

(Discretization step:  $h = \sqrt{s}$  and  $x_n \simeq x(n\sqrt{s})$ )

# The dynamical system intuition

Link with the ODEs - A guideline to study optimization algorithms

## General methodology to analyze optimization algorithms

- Interpreting the optimization algorithm as a discretization of a given ODE:

Nesterov iteration:  $x_{n+1} = y_n - h\nabla F(y_n)$  with  $y_n = x_n + \frac{n}{n+\alpha}(x_n - x_{n-1})$

Associated ODE:  $\ddot{x}(t) + \frac{\alpha}{t}\dot{x}(t) + \nabla F(x(t)) = 0$ .

- Analysis of ODEs using a Lyapunov approach:

$$\mathcal{E}(t) = F(x(t)) - F^* + \frac{1}{2}\|\dot{x}(t)\|^2.$$

$$\mathcal{E}(t) = t^2(F(x(t)) - F(x^*)) + \frac{1}{2}\|(\alpha - 1)(x(t) - x^*) + t\dot{x}(t)\|^2 ..$$

- Building a sequence of discrete Lyapunov energies adapted to the optimization scheme to get the same decay rates

# Convergence analysis of the Nesterov gradient method

## Convergence rate in the continuous setting

Let  $F : \mathbb{R}^N \rightarrow \mathbb{R}$  be a differentiable convex function and  $x^* \in \arg \min(F) \neq \emptyset$ .

- If  $\alpha \geq 3$ ,

$$F(x(t)) - F(x^*) = \mathcal{O}\left(\frac{1}{t^2}\right)$$

[Attouch, Chbani,  
Peypouquet, Redont 2016]

- If  $\alpha > 3$ , then  $x(t)$  cv to a minimizer of  $F$  and:

$$F(x(t)) - F(x^*) = o\left(\frac{1}{t^2}\right)$$

[Su, Boyd, Candes 2016]  
[Chambolle, Dossal 2017]  
[May 2017]

- If  $\alpha < 3$  then no proof of cv of  $x(t)$  but:

$$F(x(t)) - F(x^*) = \mathcal{O}\left(\frac{1}{t^{\frac{2\alpha}{3}}}\right)$$

[Attouch, Chbani, Riahi 2019]  
[Aujol, Dossal 2017]

# The Nesterov's accelerated gradient method

## State of the art results

Let  $F : \mathbb{R}^N \rightarrow \mathbb{R}$  be a differentiable convex function with  $X^* := \arg \min(F) \neq \emptyset$ .

$$\begin{cases} y_n &= x_n + \frac{n}{n+\alpha}(x_n - x_{n-1}) \\ x_{n+1} &= y_n - h \nabla F(y_n) \end{cases}, \quad \alpha > 0, \quad h < \frac{1}{L}$$

- If  $\alpha \geq 3$

$$F(x_n) - F(x^*) = \mathcal{O}\left(\frac{1}{n^2}\right)$$

[Nesterov 1984, Su, Boyd, Candes 2016,  
Chambolle Dossal 2015, Attouch et al. 2018]

- If  $\alpha > 3$ , then  $(x_n)_{n \geq 1}$  cv and:

$$F(x_n) - F(x^*) = o\left(\frac{1}{n^2}\right)$$

[Chambolle, Dossal 2014]  
[Attouch, Peypouquet 2015]

- If  $\alpha \leq 3$

$$F(x_n) - F(x^*) = \mathcal{O}\left(\frac{1}{n^{\frac{2\alpha}{3}}}\right).$$

[Attouch, Chbani, Riahi 2018]  
[Apidopoulos, Aujol, Dossal 2018]

# Convergence analysis of the Nesterov gradient method

Proof of the convergence rate  $\mathcal{O}\left(\frac{1}{t^2}\right)$  when  $\alpha \geq 3$

A first Lyapunov energy

$$\ddot{x}(t) + \frac{\alpha}{t} \dot{x}(t) + \nabla f(x(t)) = 0$$

$$\mathcal{E}(t) = F(x(t)) - F(x^*) + \frac{1}{2} \|\dot{x}(t)\|^2$$

be the mechanical energy associated to the ODE.

$$\begin{aligned} \mathcal{E}'(t) &= \langle \nabla F(x(t)), \dot{x}(t) \rangle + \langle \ddot{x}(t), \dot{x}(t) \rangle \\ &= -\frac{\alpha}{t} \|\dot{x}(t)\|^2 \leq 0 \end{aligned}$$

$$\hookrightarrow F(x(t)) - F^* \leq \frac{M}{t}$$

# Convergence analysis of the Nesterov gradient method

Proof of the convergence rate  $\mathcal{O}\left(\frac{1}{t^2}\right)$  when  $\alpha \geq 3$

A second Lyapunov energy to get the rate  $\mathcal{O}\left(\frac{1}{t^2}\right)$  Can we prove that the energy:

$$E(t) = t^2 (F(x(t)) - F(x^*)) + \frac{t^2}{2} \|\dot{x}(t)\|^2$$

is bounded ? The answer is : NO

# The Nesterov's accelerated gradient method

Proof of the convergence rate  $\mathcal{O}\left(\frac{1}{t^2}\right)$  under convexity (Su Boyd Candes 2014)

We define :

$$\mathcal{E}(t) = t^2(F(x(t)) - F(x^*)) + \overbrace{\frac{1}{2} \|(\alpha - 1)(x(t) - x^*) + t\dot{x}(t)\|^2}^{\geq 0}.$$

Using (ODE), a straightforward computation shows that:

$$\begin{aligned}\mathcal{E}'(t) &= -(\alpha - 1)t \underbrace{\langle \nabla F(x(t)), x(t) - x^* \rangle}_{\geq F(x(t)) - F(x^*) \text{ by convexity}} + 2t(F(x(t)) - F(x^*)) \\ &\leq \underbrace{(3 - \alpha)t}_{\leq 0} (F(x(t)) - F(x^*)).\end{aligned}$$

If  $\alpha \geq 3$  then  $\mathcal{E}'(t) \leq 0 \Rightarrow \forall t \geq t_0, \mathcal{E}(t) \leq \mathcal{E}(t_0)$

$$\Rightarrow \forall t \geq t_0, t^2(F(x(t)) - F^*) \leq \mathcal{E}(t) \leq \mathcal{E}(t_0)$$

$$\Rightarrow \forall t \geq t_0, F(x(t)) - F^* \leq \frac{\mathcal{E}(t_0)}{t^2}$$

$$\underline{\text{If } \alpha > 3} \quad \mathcal{E}'(t) \leq \frac{3 - \alpha}{t} \mathcal{E}(t) \Rightarrow \mathcal{E}(t) \leq \frac{\mathcal{E}(t_0)}{t^{\alpha-3}}, t \geq t_0$$



$$\forall r \geq t_0, \quad \varepsilon(r) \leq \frac{\varepsilon(t_0)}{r^{\alpha-3}} \Rightarrow \forall r \geq t_0, \quad t^2(F(x(t)) - F^*) \leq \frac{\varepsilon(t_0)}{r^{\alpha-3}}$$

$$\Rightarrow t(F(x(t)) - F^*) \leq \frac{\varepsilon(t_0)}{r^{\alpha-2}} \quad \begin{array}{l} \alpha > 3 \\ \alpha - 2 > 1 \end{array}$$

$$\Rightarrow \int_{t_0}^{+\infty} t(F(x(t)) - F^*) < +\infty$$

# The continuous, a guideline to analyse the Nesterov scheme

## $\mathcal{C}_\alpha$ for the class of differentiable convex functions

- Continuous setting:

$$\mathcal{E}(t) = t^2(F(x(t)) - F(x^*)) + \frac{1}{2} \|(\alpha - 1)(x(t) - x^*) + t\dot{x}(t)\|^2.$$

- Discrete setting:

$$\mathcal{E}_n = n^2(F(x_n) - F(x^*)) + \frac{1}{2h} \|(\alpha - 1)(x_n - x^*) + n(x_n - x_{n-1})\|^2$$

Using the definition of  $(x_n)_{n \geq 1}$  and the following convex inequality

$$F(x_n) - F(x^*) \leq \langle x_n - x^*, \nabla F(x_n) \rangle$$

we get

$$\mathcal{E}_{n+1} - \mathcal{E}_n \leq (3 - \alpha)n(F(x_n) - F(x^*)) \quad (1)$$

- 1 If  $\alpha \geq 3$ ,  $\forall n \geq 1$ ,  $n^2(F(x_n) - F(x^*)) \leq \mathcal{E}_1$
- 2 If  $\alpha > 3$ ,  $\sum_{n \geq 1} (\alpha - 3)n(F(x_n) - F(x^*)) \leq \mathcal{E}_1$

## Natural extension to composite optimization

$$\min_x F(x) = \underbrace{f(x)}_{\substack{\text{convex} \\ \text{differentiable} \\ L\text{-Lipschitz gradient}}} + \underbrace{g(x)}_{\substack{\text{convex p.s.c. "simple"} \\ \text{possibly non smooth.}}}$$

- 1 Notion of convex subdifferential
- 2 Proximal operator
- 3 The Forward Backward algorithm and FISTA

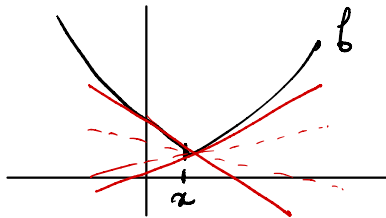
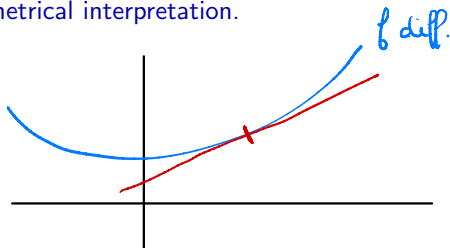
# Subdifferential of a convex function

## Definition

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  be a convex function and  $x_0 \in \text{dom}(f)$ . The subdifferential of  $f$  at  $x_0$ , denoted by  $\partial f(x_0)$ , is defined by:

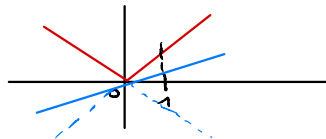
$$\partial f(x_0) = \{g \in \mathbb{R}^n \mid \forall x \in \text{dom}(f), f(x) \geq f(x_0) + \langle g, x - x_0 \rangle\}.$$

Geometrical interpretation.



Ex  $f(x) = |x|, x \in \mathbb{R}$

$$\partial f(1) = \{+1\} \quad \partial f(0) = [-1, 1]$$



# Subdifferential of a convex function

## Properties

$$\partial f(x_0) = \{g \in \mathbb{R}^n \mid \forall x \in \text{dom}(f), f(x) \geq f(x_0) + \langle g, x - x_0 \rangle\}$$

- 1 If  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  is convex then  $\partial f(x)$  is also convex and closed for any  $x \in \text{dom}(f)$ .
- 2 If  $f$  is convex l.s.c then  $\partial f(x)$  is additionally non empty and bounded for any  $x \in \text{dom}(f)$ .

# Subdifferential of a convex function

## Subdifferential calculus rules

Let  $f$  be a convex l.s.c. function. Then:

- 1 If  $f$  is differentiable on  $\text{dom} f$  then:

$$\forall x \in \text{int}(\text{dom}(f)), \partial f(x) = \{\nabla f(x)\}.$$

# Subdifferential of a convex function

## Subdifferential calculus rules

Let  $f$  be a convex l.s.c. function. Then:

- ① If  $f$  is differentiable on  $\text{dom} f$  then:

$$\forall x \in \text{int}(\text{dom}(f)), \partial f(x) = \{\nabla f(x)\}.$$

- ② Let  $\text{dom}(f) \subseteq \mathbb{R}^n$ ,  $A : \mathbb{R}^m \rightarrow \mathbb{R}^n$  be a linear operator and  $b \in \mathbb{R}^n$  then  $x \mapsto \phi(x) = f(Ax + b)$  is convex l.s.c. and:

$$\forall x \in \text{int}(\text{dom}(\phi)), \partial \phi(x) = A^\top \partial f(Ax + b).$$

# Subdifferential of a convex function

## Subdifferential calculus rules

Let  $f$  be a convex l.s.c. function. Then:

- ① If  $f$  is differentiable on  $\text{dom} f$  then:

$$\forall x \in \text{int}(\text{dom}(f)), \partial f(x) = \{\nabla f(x)\}.$$

- ② Let  $\text{dom}(f) \subseteq \mathbb{R}^n$ ,  $A : \mathbb{R}^m \rightarrow \mathbb{R}^n$  be a linear operator and  $b \in \mathbb{R}^n$  then  $x \mapsto \phi(x) = f(Ax + b)$  is convex l.s.c. and:

$$\forall x \in \text{int}(\text{dom}(\phi)), \partial \phi(x) = A^\top \partial f(Ax + b).$$

- ③ Si  $f(x) = \alpha_1 f_1(x) + \alpha_2 f_2(x)$  avec  $f_1$  et  $f_2$  convex l.s.c. on  $\mathbb{R}^n$  and  $(\alpha_1, \alpha_2) \in \mathbb{R}_+ \times \mathbb{R}_+$ . Then:

$$\begin{aligned} \partial f(x) &= \alpha_1 \partial f_1(x) + \alpha_2 \partial f_2(x) \\ &= \{g \in \mathbb{R}^n \mid \exists (x_1, x_2) \in \partial f_1(x_1) \times \partial f_2(x_2), g = \alpha_1 x_1 + \alpha_2 x_2\}. \end{aligned}$$



# Subdifferential of a convex function

## Subdifferential calculus rules

Let  $f$  be a convex l.s.c. function. Then:

- ① If  $f$  is differentiable on  $\text{dom} f$  then:

$$\forall x \in \text{int}(\text{dom}(f)), \partial f(x) = \{\nabla f(x)\}.$$

- ② Let  $\text{dom}(f) \subseteq \mathbb{R}^n$ ,  $A : \mathbb{R}^m \rightarrow \mathbb{R}^n$  be a linear operator and  $b \in \mathbb{R}^n$  then  $x \mapsto \phi(x) = f(Ax + b)$  is convex l.s.c. and:

$$\forall x \in \text{int}(\text{dom}(\phi)), \partial \phi(x) = A^\top \partial f(Ax + b).$$

- ③ Si  $f(x) = \alpha_1 f_1(x) + \alpha_2 f_2(x)$  avec  $f_1$  et  $f_2$  convex l.s.c. on  $\mathbb{R}^n$  and  $(\alpha_1, \alpha_2) \in \mathbb{R}_+ \times \mathbb{R}_+$ . Then:

$$\begin{aligned} \partial f(x) &= \alpha_1 \partial f_1(x) + \alpha_2 \partial f_2(x) \\ &= \{g \in \mathbb{R}^n \mid \exists (x_1, x_2) \in \partial f_1(x_1) \times \partial f_2(x_2), g = \alpha_1 x_1 + \alpha_2 x_2\}. \end{aligned}$$

- ④ Let  $g : \mathbb{R} \cup \{+\infty\} \rightarrow \mathbb{R} \cup \{+\infty\}$  a non-decreasing convex function. Let:  $h = g \circ f$ . Then:

$$\forall x \in \text{int}(\text{dom}(f)), \partial h(x) = \{\eta_1 \eta_2 \mid \eta_1 \in \partial g(f(x)), \eta_2 \in \partial f(x)\}. \quad (2)$$

# Subdifferential of a convex function

## Exercices

$$① \quad f(x) = |x| \qquad \partial f(x) = [-1, 1]$$

$$\begin{aligned} ② \quad f(x, y) &= x + 2|y| \qquad \partial f(x, 0) = \left\{ \begin{pmatrix} 1 \\ 2g \end{pmatrix} ; g \in \partial | \cdot | (0) \right\} \\ &= \left\{ \begin{pmatrix} 1 \\ 2g \end{pmatrix} / g \in [-1, 1] \right\} \end{aligned}$$

$$③ \quad f(x, y) = y^2 + |x^2 + 3y|$$

}

$$④ \quad f(x) = \|x\|_2^2, \quad x \in \mathbb{R}^N. \qquad \partial f(x) = \{2x\}$$

## CNS of optimality

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  be a convex function. The point  $x^* \in \mathbb{R}^n$  is a global minimum point of  $f$  si et seulement si:

$$0 \in \partial f(x^*).$$

Elementary proof: If  $f$  differentiable,  $\partial f(x^*) = \{\nabla f(x^*)\}$   
 $0 \in \partial f(x^*) \Leftrightarrow \nabla f(x^*) = 0$

$x^*$  global minimum pt of  $f$

$$\Leftrightarrow \forall x \in \mathbb{R}^n, \quad f(x) \geq f(x^*)$$

$$\Leftrightarrow \forall x \in \mathbb{R}^n, \quad f(x) \geq \underline{f(x^*) + \langle 0, x - x^* \rangle}$$

$$\Leftrightarrow 0 \in \partial f(x^*)$$

# Proximal operator

## Definition

$$\text{prox}_f(x) = \arg \min_{u \in \mathbb{R}^n} f(u) + \frac{1}{2} \|u - x\|_2^2.$$

## Examples

- $f(x) = \chi_X(x) = \begin{cases} 0 & \text{if } x \in X \\ +\infty & \text{otherwise,} \end{cases}$   $X$  convex

$$\min_{u \in \mathbb{R}^n} f(u) + \frac{1}{2} \|u - x\|_2^2 = \min_{u \in X} \frac{1}{2} \|u - x\|_2^2 \rightsquigarrow p_X(x)$$

- $f(x) = 0$   $\min_u 0 + \frac{1}{2} \|u - x\|_2^2 \rightsquigarrow u^* = x$

$$\text{prox}_0(x) = \text{id}(x)$$

- $f(x) = \alpha \|x\|_1 \rightsquigarrow \text{FISTA}.$

$$\text{prox}_{\gamma f}(x) = \max(0, |x| - \gamma) \times \text{sign}(x)$$

# Proximal operator

## Properties

- ① Let  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  be a convex l.s.c. function. Then  $\text{prox}_f(x)$  exists and is unique for any  $x \in \mathbb{R}^n$ . Moreover

$$p = \text{prox}_f(x) \Leftrightarrow x - p \in \partial f(p)$$

- ② Let  $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  be a convex function and  $t > 0$ . Then:

$$\bar{x} \in \arg \min_{x \in \mathbb{R}^n} f(x) \iff \bar{x} = \text{prox}_{tf}(\bar{x}).$$

## Proximal point algorithm

$$\begin{aligned} x_0 &\in \mathbb{R}^n \\ x_{k+1} &= \text{prox}_{tf}(x_k), \quad t > 0. \end{aligned}$$

# Forward-Backward algorithm

Let:

$$\min_{x \in \mathbb{R}^n} f(x) = g(x) + h(x)$$

where:

- $g : \mathbb{R}^n \rightarrow \mathbb{R}$  convex differentiable function having a  $L$ -Lipschitz gradient:

$$\|\nabla g(x_1) - \nabla g(x_2)\| \leq L\|x_1 - x_2\|, \quad \forall (x_1, x_2) \in \mathbb{R}^n \times \mathbb{R}^n. \quad (3)$$

- $h : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  a convex l.s.c. function: *simple = we can compute its prox.*

## Optimality condition

$$0 \in \partial f(x) \Leftrightarrow x = \text{prox}_{sh}(x - s\nabla g(x)).$$

$$T = \text{prox}_{sh} \circ (\mathcal{I} - s\nabla g)$$

# Forward-Backward and FISTA algorithms

## Forward-Backward algorithm

$$\begin{aligned}x_0 &\in \mathbb{R}^n \\x_{k+1} &= \text{prox}_{sh}(x_k - s\nabla g(x_k)), \quad s > 0.\end{aligned}$$

If  $s < \frac{2}{L}$  then the sequence  $(x_k)_{k \in \mathbb{N}}$  converge to a global minimum point of  $f$  and

$$\forall k \in \mathbb{N}, \quad F(x_k) - F(x^*) \leq \frac{2\|x_0 - x^*\|^2}{sk}$$

## FISTA algorithm

$$\begin{aligned}y_k &= x_k + \alpha_k(x_k - x_{k-1}) \\x_{k+1} &= \text{prox}_{sg}(y_k - s\nabla f(y_k))\end{aligned}$$

If  $s < \frac{1}{L}$  then the sequence  $(x_k)_{k \in \mathbb{N}}$  converge to a global minimum point of  $f$  and

$$F(x_k) - F(x^*) = o\left(\frac{1}{k^2}\right)$$

## Some numerical experiments



- 1 Preliminaries: local geometry of convex functions
- 2 Gradient descent methods
- 3 Accelerated gradient methods
  - The Heavy ball method
  - The Nesterov's accelerated gradient method
  - Natural extension to composite optimization
  - Some numerical experiments
- 4 Improving the state of the art results with the help of the geometry
- 5 Conclusion

# Improving the state of the art results with the help of the geometry

Joint work with J.-F. Aujol, Ch.Dossal

Let  $F : \mathbb{R}^N \rightarrow \mathbb{R}$  be a differentiable convex function and  $x^* \in \arg \min(F) \neq \emptyset$ .

- If  $\alpha \geq 3$ ,

$$F(x(t)) - F(x^*) = \mathcal{O}\left(\frac{1}{t^2}\right)$$

[Attouch, Chbani,  
Peypouquet, Redont 2016]

- If  $\alpha > 3$ , then  $x(t)$  cv to a minimizer of  $F$  and:

$$F(x(t)) - F(x^*) = o\left(\frac{1}{t^2}\right)$$

[Su, Boyd, Candes 2016]  
[Chambolle, Dossal 2017]  
[May 2017]

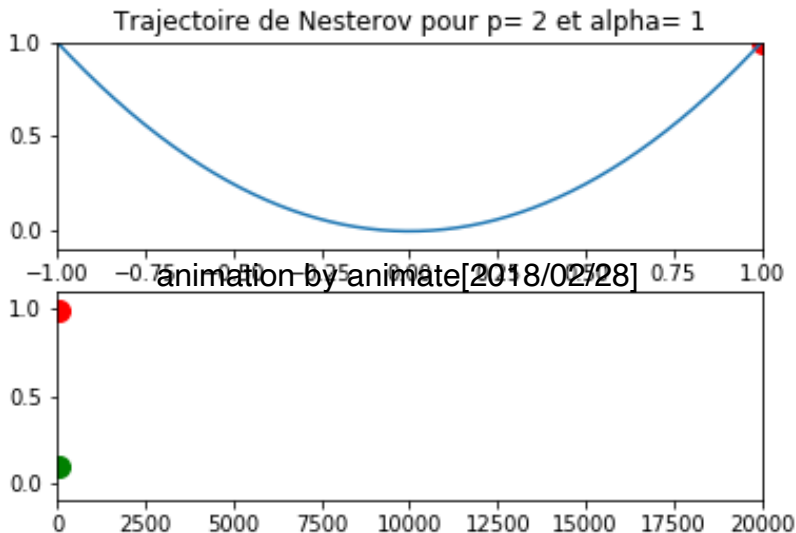
- If  $\alpha < 3$  then no proof of cv of  $x(t)$  but:

$$F(x(t)) - F(x^*) = \mathcal{O}\left(\frac{1}{t^{\frac{2\alpha}{3}}}\right)$$

[Attouch, Chbani, Riahi 2019]  
[Aujol, Dossal 2017]

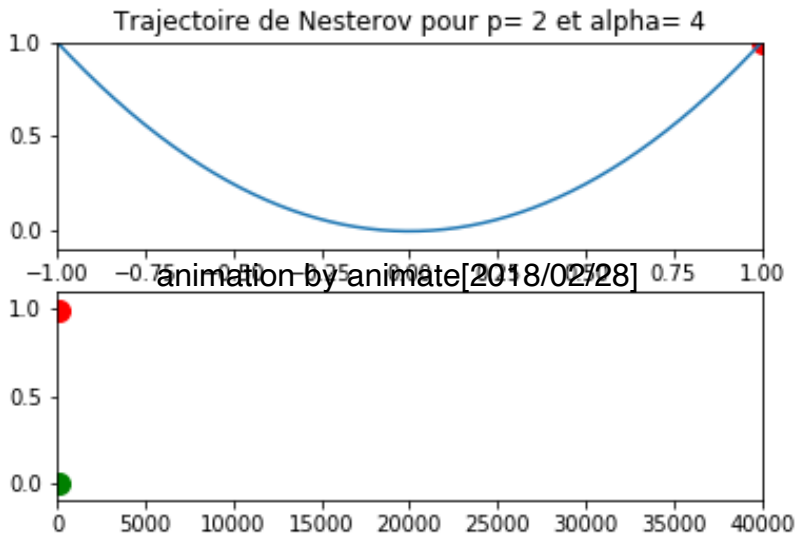
**First Example :**  $F(x) = x^2$  and  $\alpha = 1$  - **State of the art rate:**  $\mathcal{O}(\frac{1}{n^{2/3}})$

In blue  $F(x_n)$ , in orange  $n \times (F(x_n) - F^*)$



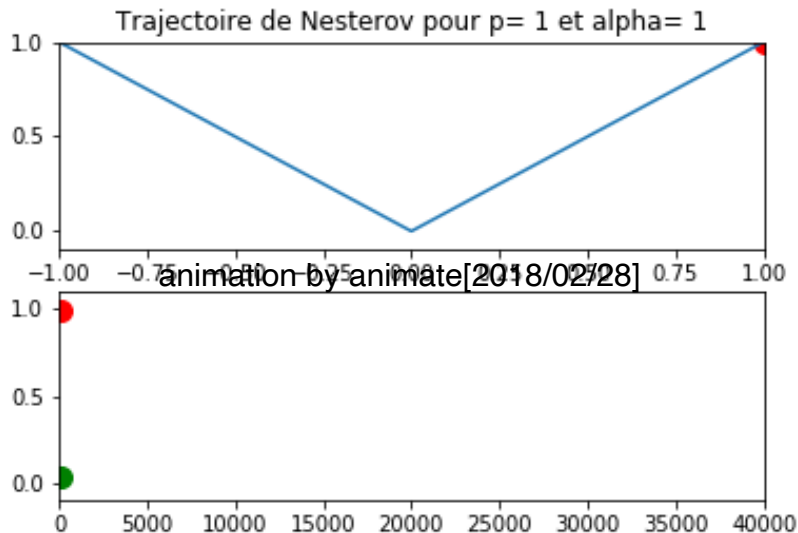
**Second Example :**  $F(x) = x^2$  and  $\alpha = 4$  - **State of the art rate:**  $\mathcal{O}(\frac{1}{n^2})$

In blue  $F(x_n)$ , in orange  $n^4 \times (F(x_n) - F^*)$



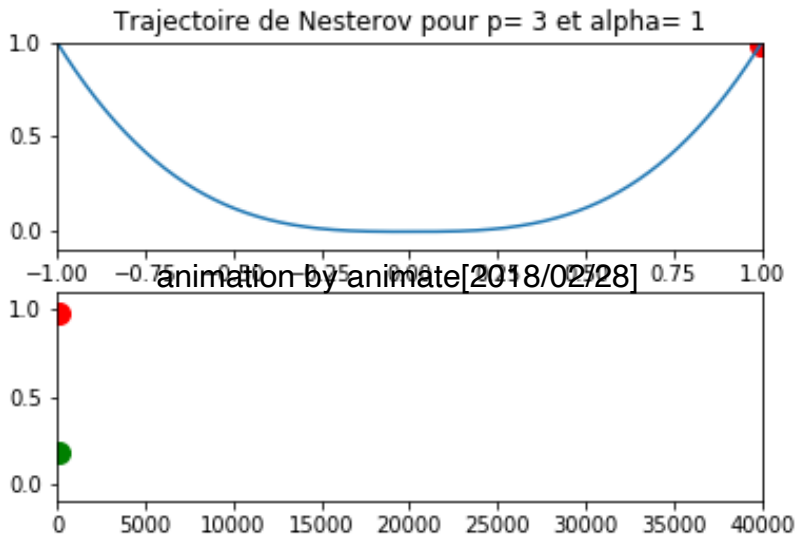
### A Third Example FISTA : $F(x) = |x|$ and $\alpha = 1$

In blue  $F(x_n)$ , in orange  $n^{\frac{2}{3}} \times F(x_n)$



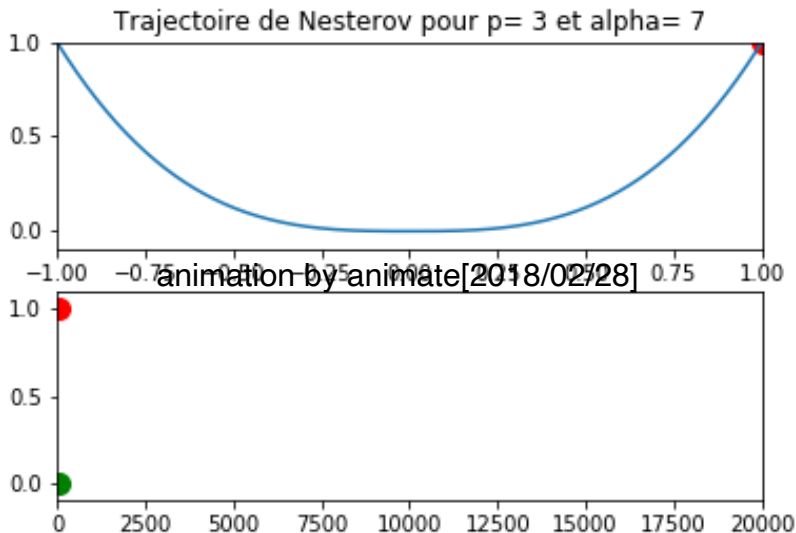
**Fourth Example :**  $F(x) = |x|^3$  and  $\alpha = 1$  - **State of the art rate:**  
 $\mathcal{O}(\frac{1}{n^{2/3}})$

In blue  $F(x_n)$ , in orange  $n^{\frac{6}{5}} \times (F(x_n) - F^*)$



**Fifth Example :**  $F(x) = |x|^3$  and  $\alpha = 7$  - State of the art rate:  $\mathcal{O}(\frac{1}{n^2})$

In blue  $F(x_n)$ , in orange  $n^6 \times (F(x_n) - F^*)$



# Local geometry of convex function

## Flatness

### Definition

Let  $F : \mathbb{R}^N \rightarrow \mathbb{R}$  be a differentiable function.  $F$  has the  $\mathcal{H}(\gamma)$  property for some  $\gamma \geq 1$  if

$$\forall x \in \mathbb{R}^n, F(x) - F(x^*) \leq \frac{1}{\gamma} \langle \nabla F(x), x - x^* \rangle.$$

### Flatness properties

- If  $(F - F^*)^{\frac{1}{\gamma}}$  is convex, then  $F$  satisfies  $\mathcal{H}(\gamma)$ .
- If  $F$  satisfies  $\mathcal{H}(\gamma)$  then for any  $x^* \in X^*$ , there exist  $C > 0$  and  $\eta > 0$  such that

$$\forall x \in B(x^*, \eta), F(x) - F(x^*) \leq C \|x - x^*\|^\gamma.$$

- Ensures that  $F$  is sufficiently flat (*at least as flat as  $\|x\|^\gamma$* ) around its set of minimizers.
- May prevent from bad oscillations of the solution.



## Theorems for sharp functions [Aujol, Dossal, R. 2018]

$$\ddot{x}(t) + \frac{\alpha}{t} \dot{x}(t) + \nabla F(x(t)) = 0.$$

Assume now that  $F$  is  $\mu$ -strongly convex, satisfies the flatness condition  $\mathcal{H}(\gamma)$  and admits a unique minimizer  $x^*$ . Then:

$$F(x(t)) - F(x^*) = \mathcal{O}\left(\frac{1}{t^{\frac{2\alpha\gamma}{\gamma+2}}}\right) \quad (4)$$

$$x_{n+1} = y_n - h \nabla F(y_n) \quad \text{where:} \quad y_n = x_n + \frac{n}{n+\alpha}(x_n - x_{n-1}), \quad \alpha > 0, \quad h < \frac{1}{L}$$

### Theorem for sharp functions (Apidopoulos, Aujol, Dossal, R. (2018))

Assume that  $F$  is strongly convex and satisfies  $\mathcal{H}(\gamma)$  for some  $\gamma \in [1, 2]$ .

$$\forall \alpha > 0, \quad F(x_n) - F(x^*) = \mathcal{O}\left(\frac{1}{n^{\frac{2\gamma\alpha}{\gamma+2}}}\right). \quad (5)$$

## Theorems for sharp functions [Aujol, Dossal, R. 2018]

### Comments

#### Theorem for sharp functions (Apidopoulos, Aujol, Dossal, R. (2018))

Assume that  $F$  is strongly convex and satisfies  $\mathcal{H}(\gamma)$  for some  $\gamma \in [1, 2]$ .

$$\forall \alpha > 0, F(x_n) - F(x^*) = \mathcal{O}\left(\frac{1}{n^{\frac{2\gamma\alpha}{\gamma+2}}}\right). \quad (6)$$

### Comments

- For  $\gamma = 1$  we recover the decay  $\mathcal{O}\left(\frac{1}{n^{\frac{2\alpha}{3}}}\right)$  from [Attouch, Cabot 2018].
- Since  $\nabla F$  is  $L$ -Lipschitz and satisfies  $\mathcal{L}(2)$ ,  $F$  automatically satisfies  $\mathcal{H}(\gamma)$  for some  $\gamma > 1$  and thus

$$\frac{2\gamma\alpha}{\gamma+2} > \frac{2\alpha}{3}.$$

- For quadratic functions (i.e. for  $\gamma = 2$ ), we get  $\mathcal{O}\left(\frac{1}{n^\alpha}\right)$ .

# Theorems for sharp functions [Aujol, Dossal, R. 2018]

## Sketch of proof in the continuous case

- ① We define for  $(p, \xi, \lambda) \in \mathbb{R}^3$

$$\mathcal{H}(t) = t^p \left( t^2(F(x(t)) - F^*) + \frac{1}{2} \|(\lambda(x(t) - x^*) + t\dot{x}(t))\|^2 + \frac{\xi}{2} \|x(t) - x^*\|^2 \right)$$

- ② We choose  $(p, \xi, \lambda) \in \mathbb{R}^3$  depending on the hypotheses to ensure that  $\mathcal{H}$  is bounded.  $\mathcal{H}$  may not be non increasing.

- ▶ For the class of convex functions, take:  $p = 0, \lambda = \alpha - 1, \xi = 0$ .
- ▶ For the class of sharp convex functions, take:

$$p = \frac{2\alpha\gamma}{\gamma+2} - 2, \lambda = \frac{2\alpha}{\gamma+2}, \xi = \lambda(\lambda + 1 - \alpha).$$

- ③ We deduce that there exists  $A \in \mathbb{R}$  such that

$$t^{2+p}(F(x(t)) - F(x^*)) \leq A - t^p \frac{\xi}{2} \|x(t) - x^*\|^2$$

- ④ If  $\xi \geq 0$  then  $F(x(t)) - F(x^*) = \mathcal{O}\left(\frac{1}{t^{p+2}}\right)$ .
- ⑤ If  $\xi \leq 0$  we must use the strong convexity to conclude.

## Convergence rates for flat functions

### Theorem for flat functions (Apidopoulos, Aujol, Dossal, R. (2018))

Let  $\gamma > 2$ . If  $F$  has a unique minimizer  $x^*$ , if  $F$  satisfies the flatness condition  $\mathcal{H}(\gamma)$  and the growth condition:

$$\forall x \in \mathbb{R}^n, \quad \frac{\mu}{2} \|x - x^*\|^\gamma \leq F(x) - F^*$$

Then if  $\alpha > \frac{\gamma+2}{\gamma-2}$

$$F(x_n) - F(x^*) = O\left(\frac{1}{n^{\frac{2\gamma}{\gamma-2}}}\right).$$

### Comments

- Better rate than  $o(\frac{1}{n^2})$ .
- Better rate than for the Gradient descent: if  $F$  satisfies  $\mathcal{L}(\gamma)$  with  $\gamma > 2$ , then

$$F(x_n) - F(x^*) = O\left(\frac{1}{n^{\frac{\gamma}{\gamma-2}}}\right)$$

[Garrigos et al. 2017].

## Application to the linear Least Square problem

Let  $A : \mathbb{R}^N \rightarrow \mathbb{R}^N$  a positive definite bounded linear operator and  $y \in \mathbb{R}^N$ . Consider

$$\min_{x \in \mathbb{R}^N} F(x) := \frac{1}{2} \|Ax - y\|^2.$$

- $F$  is convex and has a  $L$ -Lipschitz continuous gradient ( $L = \|A^*A\|$ ).
- As a convex quadratic function, we have:

$$F(x) - F(x^*) = \frac{1}{2} \langle \nabla F(x), x - x^* \rangle = \frac{1}{2} \|A(x - x^*)\|^2.$$

►  $F$  satisfies  $\mathcal{H}(\gamma)$  for any  $\gamma \in [1, 2]$ , and  $\mathcal{L}(2)$ .

- $\forall n, x_n \in \{x_0\} + \text{Im}(A^*)$ .

Since this problem has a unique solution on the space  $\{x_0\} + \text{Im}(A^*)$ , our theorem is still applicable and:

$$F(x_n) - F^* = \mathcal{O}\left(\frac{1}{n^\alpha}\right).$$

## To sum up

### Two ingredients to get better convergence rates on $F(x_n) - F^*$

- A **sharpness** condition
  - ▶ Ensuring that the magnitude of the gradient is not too low in the neighborhood of the minimizers.
- A **flatness** condition.
  - ▶ Ensuring that  $F$  is not too sharp in the neighborhood of its minimizers to prevent from bad oscillations of the solution.

Optimal convergence rates for Nesterov acceleration. J.-F. Aujol, Ch. Dossal, A. Rondepierre. May 2018.

Convergence rates of an inertial gradient descent algorithm under growth and flatness conditions. V. Apidopoulos, J.-F. Aujol, Ch. Dossal, A. Rondepierre. December 2018.

- 1 Preliminaries: local geometry of convex functions
- 2 Gradient descent methods
- 3 Accelerated gradient methods
  - The Heavy ball method
  - The Nesterov's accelerated gradient method
  - Natural extension to composite optimization
  - Some numerical experiments
- 4 Improving the state of the art results with the help of the geometry
- 5 Conclusion

# Is the Nesterov's method really an acceleration of the GD ?

## A first conclusion

- If  $F$  is sharp, Gradient Descent is faster than Nesterov.
- If  $F$  is flat, Nesterov is faster than Gradient Descent.
- Choose  $\alpha$  as large as possible

## A second conclusion : it's more complicated

- Constants in big  $\mathcal{O}$  or in geometric decays may be important.

For example in the convex case ( $\gamma = 1$ ), the constant in  $\mathcal{O}\left(t^{-\frac{2\alpha}{3}}\right)$  is of the form:

$$\forall t \geq \frac{\alpha}{\sqrt{\mu}}, \quad F(x(t)) - F(x^*) \leq CE_m(t_0) \left( \frac{\alpha}{t\sqrt{\mu}} \right)^{\frac{2\alpha}{3}}$$

- Nesterov with restart and backtracking may outperform Conjugate Gradient on the least square problem.