

# Complexity in continuous optimization (5/6)

February 20, 2025

Today: High-order methods

Mathematical Programming (2020) 184:71–120  
<https://doi.org/10.1007/s10107-019-01406-y>

FULL LENGTH PAPER

Series A



## Lower bounds for finding stationary points I

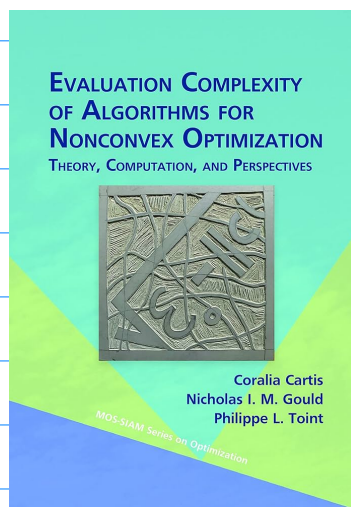
Yair Carmon<sup>1</sup> · John C. Duchi<sup>2</sup> · Oliver Hinder<sup>3</sup> · Aaron Sidford<sup>3</sup>

SIAM J. OPTIM.  
Vol. 23, No. 2, pp. 1092–1125

© 2013 Society for Industrial and Applied Mathematics

## AN ACCELERATED HYBRID PROXIMAL EXTRAGRADIENT METHOD FOR CONVEX OPTIMIZATION AND ITS IMPLICATIONS TO SECOND-ORDER METHODS\*

RENATO D. C. MONTEIRO<sup>†</sup> AND B. F. SVAITER<sup>‡</sup>



Math. Program., Ser. B (2008) 112:159–181  
DOI 10.1007/s10107-006-0089-x

FULL LENGTH PAPER

## Accelerating the cubic regularization of Newton's method on convex problems

Yu. Nesterov

# HIGH-ORDER OPTIMIZATION

Motivation: → we have seen algorithms/complexity for problems with a  $C^{1,1}$  objective and  $C^{2,2}$  objective

Most of the complexity literature falls into this setting

→ We have seen complexity guarantees for first-order optimality ( $\|\nabla f(x)\| \leq \epsilon_1$ )

second-order optimality ( $\|\nabla f(x)\| \leq \epsilon_1$  and  $\lambda_{\min}(\nabla^2 f(x)) \geq -\epsilon_2$ )

⋮

Q) What can we say when the objective function is  $C^{p,p}$  for some  $p \in \mathbb{N}^*$ ?

Yes!

→ Does the use of  $p$ th-order derivatives improve complexity with respect to 1st and 2nd order guarantees? (when  $p=2$ , we know this is the case for 1st-order guarantees)

Yes in theory but...

→ Can we get stronger optimality guarantees using those derivatives?

## ① Nonconvex setting

minimize  $f(x)$   
 $x \in \mathbb{R}^m$

$f \in C_L^{p,p}$  for some  $L > 0$

$\forall x \in \mathbb{R}^m, \nabla^i f(x) \in \mathbb{R}^{\overbrace{m \times m \times \dots \times m}^{i \text{ times}}}$   
 $\therefore i$ th order derivative of  $f$  at  $x$

$$f \in C_L^{PP} \Rightarrow \forall (x, s) \in (\mathbb{R}^m)^2$$

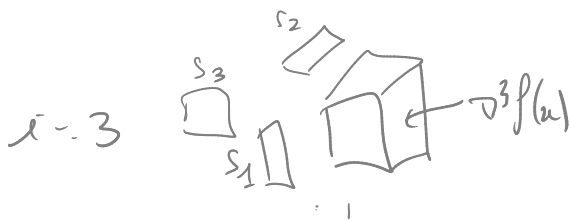
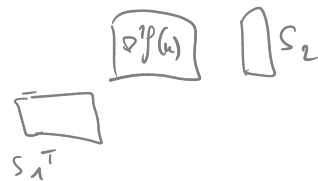
$$f(x+s) \leq f(x) + \underbrace{\sum_{i=1}^P \frac{1}{i!} \nabla^i f(x) [s, \dots, s]}_{i \text{ times}} + \frac{L}{(p+1)!} \|s\|^{p+1}$$

p<sup>th</sup> order Taylor expansion of f at x

$\nabla^i f(x) [s_1, \dots, s_i]$  :  $\nabla^i f(x)$  applied to  $i$  vectors  $s_1, \dots, s_i$

$$i=1 \quad \nabla^1 f(x) [s_1] \equiv \nabla f(x)^T s_1 \quad \boxed{\quad} \quad \boxed{s_1}$$

$$i=2 \quad \nabla^2 f(x) [s_1, s_2] \equiv s_1^T \nabla^2 f(x) s_2 \quad (= s_2^T \nabla^2 f(x) s_1)$$



$$i=1 \quad \nabla f(x) [s_1] = \sum_{j=1}^m [\nabla f(x)]_j [s_1]_j$$

$$i=2 \quad \nabla^2 f(x) [s_1, s_2] = \sum_{j_1=1}^m \sum_{j_2=1}^m [\nabla^2 f(x)]_{j_1 j_2} [s_1]_{j_1} [s_2]_{j_2}$$

$$i \quad \nabla^i f(x) [s_1, \dots, s_i] = \sum_{j_1=1}^m \dots \sum_{j_i=1}^m [\nabla^i f(x)]_{j_1 j_2 \dots j_i} [s_1]_{j_1} \dots [s_i]_{j_i}$$

NB:  $m=1$

$$f(x+s) \leq f(x) + \sum_{i=1}^P f^{(i)}(x) s^i + \frac{L}{(p+1)!} |s|^{p+1}$$

↳ By analogy with the cases  $p=1$  and  $p=2$ , we can build algorithms based on the  $p$ th order Taylor expansion.

## Trust-region method

At every iteration  $k$ , given  $(x_k, \delta_k)$ , compute

$$f(x_k), \nabla f(x_k), \dots, \nabla^p f(x_k)$$

and  $s_k$  as a solution of the subproblem

minimize  $T_p(x_k, s)$  Model of  $f$  around  $x_k$

Trust region  
( $\delta_k$ : radius)

→  $\|s\| \leq \delta_k$

where  $T_p(x_k, s) = f(x_k) + \sum_{i=1}^p \frac{1}{i!} \nabla^i f(x_k)[s, \dots, s]$

Compute  $x_{k+1}$  and  $\delta_{k+1}$  according to  $\frac{f(x_k) - f(x_k + s_k)}{T_p(x_k, 0) - T_p(x_k, s_k)}$

↑ change in the true function

↑ change in the model from  $x_k$  to  $x_{k+1}$

## Adaptive regularization (Cubic regularization when $p=2$ )

At iteration  $k$ , given  $(x_k, \sigma_k)$ , compute

$$f(x_k), \dots, \nabla^p f(x_k)$$

and compute  $s_k$  as a solution of the subproblem

$$\text{minimize}_{s \in \mathbb{R}^m} \left\{ T_p(x_k, s) + \frac{\sigma_k}{p+1} \|s\|^{p+1} \right\}$$

→ unconstrained problem  
 → Regularizer  $\|s\|^{p+1}$  guarantees that there exists a solution

Compute  $x_{k+1}$  and  $\sigma_{k+1}$  according to

$$\frac{f(x_k) - f(x_k + s_k)}{T_p(x_k, 0) - T_p(x_k, s_k)}$$

→ Both TR (Trust Region) and AR (Adaptive Regularization) only differ from the case  $p=2$  in the definition of the subproblem (and the number of derivatives they use)

## Complexity guarantees

(Basic) TR: Reach  $\| \nabla f(x) \| \leq \epsilon_1$  in at most  $O(\epsilon_1^{-(p+1)})$  iterations  
 Before:  $p=2$ ,  $O(\epsilon_1^{-2})$

Will guarantee than a second-order method for  $p > 2$   
 Related to the decrease guarantee in TR methods:

$p=1$  If  $\| \nabla f(x) \| \geq \epsilon_1$ , then

$$T_2(x_k, 0) - T_2(x_k, s_k) \geq O(\| \nabla f(x) \|^2) \geq O(\epsilon_1^2)$$

→ decrease for a step  $s = -t \nabla f(x)$

$$p > 2 \quad \nexists \|\nabla f(x_k)\| \geq \varepsilon_1$$

$$T_p(x_k, 0) - T_2(x_k, s_k) \geq \min_{0 \leq t \leq \frac{\delta_k}{\|\nabla f(x_k)\|}} (T_p(x_k, 0) - T_p(x_k, t \nabla f(x_k)))$$

$$\geq O(\|\nabla f(x_k)\|^{p+1})$$

$$\geq O(\varepsilon_1^{p+1})$$

$$(p=2, \quad \exists \|\nabla f(x_k)\| \geq \varepsilon_1$$

$$T_p(x_k, 0) - T_2(x_k, s_k) \geq O(\|\nabla f(x_k)\|^2 \varepsilon_1, \varepsilon_1^3) \\ \geq O(\varepsilon_1^3)$$

AR:

Reach  $x_k$  such that  $\|\nabla f(x_k)\| \leq \varepsilon_1$  at most

$$O\left(\varepsilon_1^{-\frac{p+1}{p}}\right) \text{ iterations}$$

$$p=1 \quad \varepsilon_1^{-2}$$

$$p=2 \quad \varepsilon_1^{-3/2}$$

$$p=3 \quad \varepsilon_1^{-4/3}$$

$$\varepsilon_1^{-\frac{p+1}{p}} \xrightarrow{p \rightarrow \infty} (\varepsilon_1^{-1})$$

Analysis: At every iteration,

$$\rightarrow T_p(x_k, 0) - T_p(x_k, s_k) \geq O(\|s_k\|^{p+1})$$

$$\rightarrow \|s_k\| \geq O(\|\nabla f(x_{k+1})\|^{1/p})$$

# Lower bounds (Carnot et al 2020)

## Generic analysis

• Problem class  $\mathcal{F} = \left\{ f: \mathbb{R}^d \rightarrow \mathbb{R} \text{ for some } d \geq 1, f \in C^{P,P}, f(0) = \inf_{x \in \mathbb{R}^d} f(x) \leq \Delta \right\}$   
*f orthogonally invariant*

• Algorithmic class  $\mathcal{A}$

$x_0 = 0_{\mathbb{R}^m}$  (without loss of generality)

$\forall k \geq 1, x_k = A_k \left( f(x_0), \nabla f(x_0), \dots, \nabla^P f(x_0) \right)$

$\uparrow$   
Mapping to  $\mathbb{R}^m$   $f(x_{k-1}), \nabla f(x_{k-1}), \dots, \nabla^P f(x_{k-1})$

$\forall a \in \mathcal{A}, \forall f \in \mathcal{F}, \forall \varepsilon_1 > 0$

$$T_{\varepsilon_1}(a, f) = \inf \left\{ k \mid \|\nabla f(x_k)\| \leq \varepsilon_1 \right\}$$

Th 1  $\forall \varepsilon_1 > 0$

$$\inf_{a \in \mathcal{A}} \sup_{f \in \mathcal{F}} T_{\varepsilon_1}(a, f) \leq C \Delta L^{1/P} \varepsilon_1^{-P/P}$$

$C > 0$

and  $\inf_{a \in \mathcal{A}} \sup_{f \in \mathcal{F}} T_{\varepsilon_1}(a, f) \geq \hat{C} \Delta (L/l)^{1/P} \varepsilon_1^{-P/P}$

$$l = O(p \log p)$$

## Orthogonal invariance of $f$

$$\forall f \in F, f: \mathbb{R}^d \rightarrow \mathbb{R} \text{ and } \forall U \in \mathbb{R}^{d' \times d}, U^T U = \text{Id},$$

$$f': \mathbb{R}^{d'} \rightarrow \mathbb{R}$$

$$x \mapsto f(U^T x) \in F$$

Because of orthogonal invariance,

$$\inf_{a \in A} \sup_{f \in F} T_\varepsilon(a, f) \geq \inf_{a \in A^{32}} \sup_{f \in F} T_\varepsilon(a, f)$$

$A^{32} \subseteq A$ : "zero-respecting algorithms"

Inspiration: Nesterov's example for accelerated gradient

$\Rightarrow$  Any algorithm would only update 1 coordinate of the iterate per iteration

$\Rightarrow$  By design:

$\rightarrow$  Start with  $x_0 = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix}$

$x^0 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \rightarrow$  It 1:  $x_1 = \begin{bmatrix} [x^*]_1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

$\rightarrow$  It  $k$ :  $x_k = \begin{bmatrix} [x^*]_1 \\ \vdots \\ [x^*]_k \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

A zero-respecting algorithm satisfies ( $x = \begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix}$ )

$$\forall k \geq 1, \text{supp}(x_k) \subseteq \bigcup_{1 \leq q \leq p} \bigcup_{j < k} \text{supp}(\nabla^q f(x_j))$$



where  $\text{supp}(\cdot)$  denotes the support

$$\forall x \in \mathbb{R}^m, \quad \text{supp}(x) = \{i \mid 1 \leq i \leq m, x_i \neq 0\}$$

Recursive definition of the support

$$\forall T \in \mathbb{R}^{\overbrace{m \times \dots \times m}^{j \text{ times}}},$$

$$\forall i = 1 \dots m, \text{ define } T^{(i)} \in \mathbb{R}^{\overbrace{m \times \dots \times m}^{j-1}}$$

$$[T^{(i)}]_{m_1, \dots, m_{j-1}} = T_{i, m_1, \dots, m_{j-1}}$$

$$\text{supp}(T) = \max_{1 \leq i \leq m} \text{supp}(T^{(i)})$$

→ A zero-respecting algorithm does not update coordinates of the iterate that correspond to zero coordinates

### Zero-respecting function

$f: \mathbb{R}^p \rightarrow \mathbb{R}$  is zero-respecting if

$\exists x$  such that  $\text{supp}(x) \subseteq \{1, \dots, i-1\}$ ,

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_{i-1} \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\bigcup_{0 \leq q \leq p} \text{supp}(\nabla^q f(x)) \subseteq \{1, \dots, i\}$$

Proof technique: Build a zero-respecting function such that any zero-respecting algorithm will take at least  $O\left(\frac{p+1}{\epsilon}\right)$  iterations to satisfy  $\|\nabla f(x)\| \leq \epsilon$

# Function construction

$$\forall T \geq 1, \bar{f}_T: \mathbb{R}^T \rightarrow \mathbb{R}$$

$$x \mapsto -\psi(x) \phi(x) + \sum_{i=2}^T \left[ \psi(-x_{i-1}) \Phi(-x_i) - \psi(x_{i-1}) \Phi(x_i) \right]$$

$$\psi(x) = \begin{cases} \exp\left(1 - \frac{1}{(2x-1)^2}\right) & \text{if } x > \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

$$\Phi(x) = \sqrt{e} \int_{-\infty}^x e^{-t^2/2} dt$$

can show that  $\bar{f}_T \in C_L^{p,p}$  with  $L \leq \exp\left(\frac{T}{2} p \log p + p\right) = o(p \exp p)$

then, define

$$f_{\varepsilon_1}: \mathbb{R}^{T_{\varepsilon_1}} \rightarrow \mathbb{R}$$

$$x \mapsto \frac{L}{\varepsilon_1} \left( \frac{\Delta \varepsilon_1}{L} \right)^{\frac{p+1}{p}} \bar{f}_{T_{\varepsilon_1}} \left( \left( \frac{L}{\Delta \varepsilon_1} \right)^{\frac{1}{p}} x \right)$$

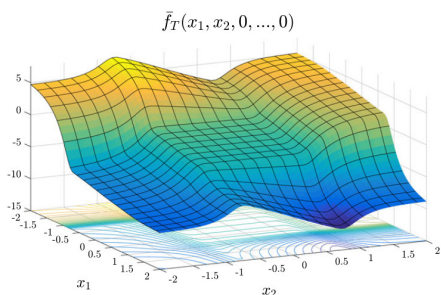
$$T_{\varepsilon_1} = \left\lfloor \frac{\Delta L^{1/p}}{12 \varepsilon_1^{1/p}} \varepsilon_1^{-\frac{p+1}{p}} \right\rfloor$$

$$f_{\varepsilon_1} \text{ is } C_{\frac{L}{\varepsilon_1}}^{p,p}, f \in \mathcal{F}$$

and

$\forall x \in \mathbb{R}^{T_{\varepsilon_1}}$  such that

$$x_{T_{\varepsilon_1}} = 0, \quad \|\nabla f_{\varepsilon_1}(x)\| > \varepsilon_1$$



→ If the algorithm only updates coordinate  $i$  at iteration  $i$ , it takes at least  $T_\varepsilon$  iterations to satisfy  $\|D/\varepsilon_i(x)\| \leq \varepsilon_i$

Current challenge: Implement the AR method for  $p \geq 3$

↳ Lack of good method for solving the subproblem  
minimize  $T_p(x, s) + \frac{\sigma}{p+1} \|s\|^{p+1}$   
 $x \in \mathbb{R}^m$

↳ Believed to be impractical for  $p \geq 4$ , the case  $p=3$  is subject of current research

## ② The convex setting

Starting point: minimize  $f(x)$ ,  $f \in C^{2,2}$   
 $x \in \mathbb{R}^m$

Adaptive cubic regularization ( $p=2$ ):  $O(\varepsilon^{-3/2})$  complexity

If  $f$  is convex, for the same algorithm:  $O(\varepsilon^{-1/2})$  complexity

same complexity than accelerated gradient, but only known to be optimal for  $C^1$  convex functions

Nesterov (2008):

For  $f \in C^{2,2}$  convex, an accelerated variant of cubic regularization has complexity  $O(\varepsilon^{-1/3})$

Monteiro and Svaiter (2013)

Their algorithm uses regularized Newton steps and for  $C^{2,2}$  convex functions, has complexity  $O(\epsilon^{-2/7})$

$\Rightarrow$  This is the optimal complexity for  $C^{2,2}$  convex functions  $\frac{2}{7} < \frac{2}{6} = \frac{1}{3}$

$\Rightarrow$  Two families of algorithms

- Optimal methods  $O(\epsilon^{-2/7})$
- Near-optimal methods  $O(\epsilon^{-1/3})$

$\Rightarrow$  Around 2018, the results were extended to high-order  $f$  ( $P$ -IP convex)

Near-optimal :  $O(\epsilon^{-\frac{1}{P+1}})$

Optimal :  $O(\epsilon^{-\frac{2}{3P+1}})$

$$\begin{array}{l} P=1 \quad \epsilon^{-\frac{1}{2}} = \epsilon^{-\frac{2}{4}} \\ P=2 \quad \epsilon^{-1/3} \quad \text{vs} \quad \epsilon^{-2/7} \end{array}$$

Lower bound :  $O(\epsilon^{-\frac{2}{3P+1}})$  (2017)

The main results in that area appeared in COLT (Conference on Learning Theory)

④ High-order conditions

1<sup>st</sup> order condition: If  $x$  minimum of  $f \in C^1$ , then  $\|\nabla f(x)\| = 0$

2<sup>nd</sup> order — : —————  $C^2$  then  $\|\nabla f(x)\| = 0$   
 $\nabla^2 f(x) \succeq 0$

3<sup>rd</sup> order — : —————  $C^3$  then  $\|\nabla f(x)\| = 0$   
 $\nabla^2 f(x) \succeq 0$

⋮  
p<sup>th</sup> order  $\nabla^3 f(x)[s, s, s] = 0$   
 $\forall s, \nabla^2 f(x)s = 0$

4<sup>th</sup> order condition  $\equiv$  Checking copositivity (Hard)