

Generalized Nesterov's accelerated proximal gradient algorithms with convergence rate of order $o(1/k^2)$

Huynh Van Ngai

University of Quy Nhon

Variational Analysis and Optimisation Webinar
March 24, 2021



Outline

Outline

- 1 An overview- gradient descent and proximal point methods
- 2 Nesterov's accelerated gradient descent and some extensions
- 3 Generalized accelerated proximal gradient algorithm
- 4 Generalized accelerated forward-backward scheme

Convergence

- If f is L -smooth (i.e., ∇f is L -Lipschitz for $L > 0$), and if step size is small enough ($\alpha_k \leq 2/L$), then the sequence $\{x_k\}$ converges to a stationary point (if it exists) of f . As f is convex, it converges to the global minimizer x^* of f .
- Convergence rate: $O(1/k)$, i.e., for some $c > 0$,

$$f(x_k) - f(x^*) \leq c/k.$$

- If f is not differentiable, $\nabla f(x_k)$ is replaced by a sub-gradient $x_k^* \in \partial f(x_k)$.

Gradient projection (GP) method

Consider the constrained optimization problem:

$$\begin{aligned} \min f(x) \\ x \in C \subseteq \mathbb{R}^n \end{aligned} \quad (2)$$

where $C \subseteq \mathbb{R}^n$ is a closed convex subset, and f is continuously differentiable.

In the constrained optimality theorem, if $\bar{x} \in C$ is a local minimum of (2), then

$$\langle \nabla f(\bar{x}), x - \bar{x} \rangle \geq 0, \quad \forall x \in C. \quad (3)$$

The (GP) method consists of the iteration:

$$x_{k+1} = P_C[x_k - \alpha_k \nabla f(x_k)]. \quad (4)$$

$P_C(z)$: the (unique) projection of $z \in \mathbb{R}^n$ on C .

The (GP) algorithm has been proposed by Goldstein in 1964. The same method was independently proposed by Levitin and Polyak one year later. This method is nowadays referred as [Goldstein-Levitin-Polyak gradient projection method](#).

Refs.

- A. Goldstein. Convex programming in Hilbert space. *Bull. Amer. Math. Soc.*, 70(5), 709-710, 1964.
- E. S. Levitin and B. T. Polyak. Constrained minimization problems. *USSR Comput. Math. Math. Phys.* 6, 1-50, 1966 (English transl. of paper in *Zh. Vychisl. Mat. i Mat. Fiz.*, 6, 787-823, 1965).

Proximal point algorithm (PPA)

Consider the optimization problem (1), where $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is a proper convex function.

The *proximal point algorithm* is the following iteration:

$$x_{k+1} = \operatorname{Argmin} \left\{ f(x) + \frac{1}{2\alpha_k} \|x - x_k\|^2 : x \in \mathbb{R}^n \right\} := \operatorname{prox}_{\alpha_k f}(x_k).$$

(PPA) is related closely to the celebrated [Tikhonov Regularization](#).

Historical References

- B. Martinet. Régularisation d'inéquations variationnelles par approximations successives. *Revue Française d'Informatique et Recherche Opérationnelle*, 1970.
- T. R. Rockafellar. Monotone operators and the proximal point algorithm, *SIAM Journal on Control and Optimization*, 14(5), 877-898, 1976.

Fast iterative shrinkage-thresholding algorithm (FISTA)-Beck & Teboulle(2009)

Consider again the composite convex optimization (5):

$$\min_{x \in \mathbb{R}^n} F(x) := f(x) + \Phi(x),$$

where f is L -smooth on \mathbb{R}^n , and Φ is l.s.c convex, possibly non-differentiable.

(FISTA) Algorithm

- $x_{k+1} = \text{prox}_{L^{-1}\Phi}(y_k - L^{-1}\nabla f(y_k)).$
- $y_{k+1} = (1 - \gamma_{k+1})x_{k+1} + \gamma_k x_k,$
- $\gamma_k = \frac{1 - \lambda_k}{\lambda_{k+1}},$
- $\lambda_k = \frac{1 + \sqrt{1 + 4\lambda_{k-1}^2}}{2}.$

Convergence result:

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

References

- Nesterov Y., A method for unconstrained convex minimization problem with the rate of convergence $O(1/k^2)$, *Doklady AN SSSR* (translated as Soviet Math.Docl.) **269**, 543-547, (1983).
- Nesterov Y., Smooth minimization of non-smooth functions, *Math. Program. Ser A.*, **103**, 127-152 (2005).
- Nesterov Y., Gradient methods for minimization composite objective functions, *Math. Prog. Ser. A*, **140**, 125-161, (2013).
- Nesterov Y., Universal Gradient methods for convex optimization problems, *Math. Prog. Ser. A*, **152**, 381-404, (2015).
- Beck A., Teboulle M., A fast iterative shrinkage-thresholding algorithm for linear inverse problems, *SIAM J. Imag. Sci.*, **2**(1), 183-202, (2009).
- Attouch H., Peypouquet J., The rate of convergence of Nesterov's accelerated forward-backward method is actually faster $O(1/k^2)$, *SIAM J. Optim.*, **26**(3), 1824-1834, (2016).

uniformly convex functions

A function $\varphi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is called p -uniformly convex with parameter μ , for some $\mu \geq 0$, $p \geq 2$, or called (μ, p) -uniformly convex if for all $x, y \in \mathbb{R}^n$, $\lambda \in [0, 1]$ one has

$$\varphi(\lambda x + (1 - \lambda)y) \leq \lambda\varphi(x) + (1 - \lambda)\varphi(y) - \frac{\mu}{p}\lambda(1 - \lambda)\|x - y\|^p.$$

When $p = 2$, the function φ is called strongly convex (with parameter μ .) Note that if φ is (μ, p) -uniformly convex, then for all $x, y \in \mathbb{R}^n$, all $x^* \in \partial\varphi(x)$, one has

$$\langle x^*, y - x \rangle \leq \varphi(y) - \varphi(x) - \frac{\mu}{p}\|y - x\|^p. \quad (8)$$

Generalized accelerated proximal gradient algorithm (GAPGA)

Parameters.

- Given a ρ -strongly convex function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ ($\rho > 0$) which attains minimum at $y_0 \in \mathbb{R}^n$:

$$h(y) \geq h(y_0) + \frac{\rho}{2} \|y - y_0\|^2, \quad \forall y \in \mathbb{R}^n, \quad (9)$$

- parameters $C, \mu > 0, 0 < \kappa \leq 1/L$,
- a sequence of positive reals $\{\alpha_k\}$; sequences of nonnegative reals $\{\beta_k\}$, and $\{\gamma_k\}$ as in Section 2. Set

$$A_k = \sum_{i=0}^k \alpha_i, \quad B_k = \sum_{i=0}^k \beta_i,$$

and also assume that $A_k \geq B_k$ for all $k \in \mathbb{N}$, and denote $A_{-1} = B_{-1} = 0$.

Algorithm 1(GAPGA1)

Initialization: Initial data: y^0 as in (9). Set $k = 0$.

Repeat: For $k = 0, 1, \dots$,

1. Find

$$\begin{aligned} x_k &= \operatorname{argmin} \left\{ \Phi(y) + \langle \nabla f(y_k), y - y_k \rangle + \frac{1}{2\kappa} \|y - y_k\|^2 : y \in \mathbb{R}^n \right\} \\ &= \operatorname{prox}_{\kappa\Phi} (y_k - \kappa \nabla f(y_k)). \end{aligned} \quad (10)$$

2. Find

$$\begin{aligned} z_k &= \operatorname{argmin}_{x \in \mathbb{R}^n} \left\{ Ch(x) + \sum_{i=0}^{k-1} \alpha_i [f(y_i) + \langle \nabla f(y_i), x - y_i \rangle + \Psi_{z_i}(x) \right. \\ &\quad \left. + \frac{1}{2} \mu \gamma_i \|x - y_i\|^2] + \alpha_k [f(y_k) + \langle \nabla f(y_k), x - y_k \rangle \right. \\ &\quad \left. + \Phi(x) + \frac{1}{2} \mu \gamma_k \|x - y_k\|^2] \right\} \end{aligned} \quad (11)$$

Algorithm 1-continued

3. Set Ψ_{z_k} is a support function to Φ at z_k such that

$$\begin{aligned}
 & \min_{x \in \mathbb{R}^n} \left\{ Ch(x) + \sum_{i=0}^{k-1} \alpha_i [f(y_i) + \langle \nabla f(y_i), x - y_i \rangle \right. \\
 & \quad \left. + \Psi_{z_i}(x) + \frac{1}{2} \mu \gamma_i \|x - y_i\|^2] \right. \\
 & \quad \left. + \alpha_k [f(y_k) + \langle \nabla f(y_k), x - y_k \rangle + \Phi(x) + \frac{1}{2} \mu \gamma_k \|x - y_k\|^2] \right\} \\
 & = \min_{x \in \mathbb{R}^n} \left\{ Ch(x) + \sum_{i=0}^k \alpha_i [f(y_i) + \langle \nabla f(y_i), x - y_i \rangle \right. \\
 & \quad \left. + \Psi_{z_i}(x) + \frac{1}{2} \mu \gamma_i \|x - y_i\|^2] \right\}.
 \end{aligned} \tag{12}$$

4. Set

$$\tau_k := \frac{\alpha_{k+1}}{A_{k+1} - B_k}, \quad y_{k+1} = \tau_k z_k + (1 - \tau_k) x_k.$$

Remarks

- In Nesterov's original accelerated schemes, $\tau_k := \frac{\alpha_{k+1}}{A_{k+1}}$, which is a particular case of Algorithm 1 with $\beta_k := 0$, $k \in \mathbb{N}$.
- In Step 3 of Algorithm 1, we can take $\Psi_{z_k} = \Phi$. If we set $\Psi_{z_k} = \Phi$, for all $k \in \mathbb{N}$, Algorithm 1 gives a generalized variant of Nesterov's accelerated dual averaging algorithm.
- An another way to choose Ψ_{z_k} is as follows. As in Step 2, z_k is a minimizer of the convex function in the right hand of (11), then there is $z_k^* \in \partial\Phi(z_k)$ such that

$$0 \in C\partial h(z_k) + \sum_{i=0}^{k-1} \alpha_i [\nabla f(y_i) + \partial\Psi_{z_i}(z_k)] + \alpha_k [\nabla f(y_k) + z_k^*] + \mu \sum_{i=0}^k \alpha_i \gamma_i (z_i - y_i). \quad (13)$$

Then the support function

$$\Psi_{z_k}(x) := \langle z_k^*, x - z_k \rangle + \Phi(z_k), \quad x \in \mathbb{R}^n, \quad (14)$$

Remarks-continued

- When $h(x) := \frac{1}{2}\|x - y_0\|^2$, and for all $k \in \mathbb{N}$, the support function Ψ_{z_k} is defined by (14) for all $k \in \mathbb{N}$, then

$$z_{k+1} = \text{prox}_{\frac{\alpha_{k+1}}{C + \mu\alpha_{k+1}\gamma_{k+1}}\Phi} \left[\frac{1}{C + \mu\alpha_{k+1}\gamma_{k+1}} W_{k+1} \right];$$

$$W_{k+1} := (C + \mu\alpha_k\gamma_k)z_k - \mu\alpha_k\gamma_k y_k + \alpha_{k+1}\gamma_{k+1}y_{k+1} - \alpha_{k+1}\nabla f(y_{k+1}). \quad (15)$$

- In particular, when $\mu = 0$, the sequence $\{z_k\}$ is defined recurrently by

$$z_{k+1} = \text{prox}_{\frac{\alpha_{k+1}}{C}\Phi} \left[z_k - \frac{\alpha_{k+1}}{C} \nabla f(y_{k+1}) \right]. \quad (16)$$

This is exactly the accelerated scheme of the proximal gradient methods.

Convergence analysis

Define the *estimate* function:

$$\begin{aligned}
 F_k(x) = & Ch(x) + \sum_{i=0}^{k-1} \alpha_i [f(y_i) + \langle \nabla f(y_i), x - y_i \rangle + \Psi_{z_i}(x) + \frac{1}{2} \mu \gamma_i \|x - y_i\|^2] \\
 & + \alpha_k [f(y_k) + \langle \nabla f(y_k), x - y_k \rangle + \Phi(x) + \frac{1}{2} \mu \gamma_k \|x - y_k\|^2], \quad x \in \mathbb{R}^n.
 \end{aligned}
 \tag{17}$$

The following theorem gives an estimate for function values $f(x_k) + \Phi(x_k)$, and it is crucial to derive the subsequent convergence rates.

Theorem

Let $\{x_k\}$ and $\{y_k\}$ be sequences generated by Algorithm 1. Suppose that $\kappa \leq 1/L$ and the sequences $\{\alpha_k\}$, $\{\beta_k\}$ and $\{\gamma_k\}$ satisfy the condition

$$\left(C\rho + \mu \sum_{i=0}^{k-1} \alpha_i \gamma_i \right) (A_k - B_{k-1}) \geq \alpha_k^2 / \kappa, \quad \text{for all } k \in \mathbb{N}. \quad (18)$$

Then one has for all $k \in \mathbb{N}$,

$$\begin{aligned} & \sum_{i=0}^k \beta_i [f(x_i) + \Phi(x_i)] + (A_k - B_k) [f(x_k) + \Phi(x_k)] \\ & + \frac{1}{2} (1/\kappa - L) \sum_{i=0}^k (A_i - B_{i-1}) \|x_i - y_i\|^2 \leq \min_{x \in \mathbb{R}^n} F_k(x). \end{aligned} \quad (19)$$

Theorem-continued

Moreover, if f is μ -strong convex, then (19) holds if $\gamma_k = 1$, $k \in \mathbb{N}$, and the sequences $\{\alpha_k\}$, $\{\beta_k\}$ verifying the condition

$$\left(C\rho + \mu \sum_{i=0}^{k-1} \alpha_i \right) (A_k - B_{k-1}) \geq \alpha_k^2 (\kappa^{-1} - \mu), \quad \text{for all } k \in \mathbb{N}. \quad (20)$$

Convergence 1

Theorem

In Algorithm 1, pick $\alpha_k = k$; $\beta_k = k/2$; $\mu = 0$, and $C, \kappa > 0$ such that $C\rho \geq \kappa^{-1} \geq L$. Then condition (18) is satisfied, and therefore for a minimizer x^ of problem (7), one has*

$$\lim_{k \rightarrow \infty} \min_{i=[k/2], \dots, k} k^2 [f(x_i) + \Phi(x_i) - f(x^*) - \Phi(x^*)] = 0, \quad (21)$$

where $[k/2]$ stands for the integer part of $k/2$. Therefore if $\{f(x_k) + \Phi(x_k)\}$ is a decreasing sequence, then

$$\lim_{k \rightarrow \infty} k^2 [f(x_k) + \Phi(x_k) - f(x^*) - \Phi(x^*)] = 0. \quad (22)$$

Convergence 2- Uniformly convex case

Theorem

Let f is (μ, p) -uniformly convex with $p > 2$, $\mu > 0$. Let $0 < \kappa \leq L^{-1}$, and $C, \rho, m > 0$ such that

$$m_{\mu\kappa} \geq \begin{cases} 2^{\frac{4}{p-2}} \frac{8\rho}{(p-2)^2} & \text{if } 2 < p < 6, \\ \frac{8\rho}{(p-2)^2} & \text{if } p \geq 6; \end{cases} \quad (23)$$

$$C\rho \geq \begin{cases} \kappa^{-1} & \text{if } 2 < p < 6, \\ \frac{p-2}{4} m_{\mu} & \text{if } p \geq 6. \end{cases} \quad (24)$$

Convergence 3: strongly convex case

Theorem

Let f is μ -strongly convex for some $\mu > 0$, and let q, C such as (??). Then for the sequence $\{x_k\}$ generated by Algorithm 1 with sequences $\alpha_k := q^k$, $\beta_k = 0$, and $\gamma_k = 1$, $k \in \mathbb{N}$, and a minimizer x^* of problem (7), one has

$$f(x_k) + \Phi(x_k) - f(x^*) - \Phi(x^*) \leq \frac{(q-1)Ch(x^*)}{q^{k+1} - 1}, \text{ for all } k \in \mathbb{N}. \quad (26)$$



(GAPGA) has convergence rate $o(1/k^2)$

Let $\{x_k\}$ be the sequences defined by Algorithm (GAPGA). Let x^* be a minimizer of problem (7).

Convergence results.

(i). For $\mu = 0$, and any two sequences of positive reals $\{\alpha_k\}$ and $\{\beta_k\}$ with $\alpha_k \geq \beta_k$ for $k \in \mathbb{N}$ and

$$0 < \liminf_{k \rightarrow \infty} \frac{\beta_k}{k} \leq \limsup_{k \rightarrow \infty} \frac{\alpha_k}{k} < +\infty, \quad \limsup_{k \rightarrow \infty} \frac{\beta_k}{\alpha_k} < 1,$$

then we can find $C_0 > 0$ such that for all $C \geq C_0$, one has

$$\lim_{k \rightarrow \infty} k^2 \min_{i=[k/2], \dots, k} [f(x_i) + \Phi(x_i) - f(x^*) - \Phi(x^*)] = 0.$$

This talk is based on the paper:

- Huynh Van Ngai & Ta Anh Son, Generalized Nesterov's accelerated proximal gradient algorithms with convergence rate of order $o(1/k^2)$, *submitted* (2020).

THANKS!