

Catalyst Meta Acceleration Framework: The history and the gist of it

Hongda Li

UBC Okanagan

November 12, 2024

1 Introduction

- The History and a Series of Papers
- Nesterov's Estimating Sequence
- Example: Accelerated proximal gradient

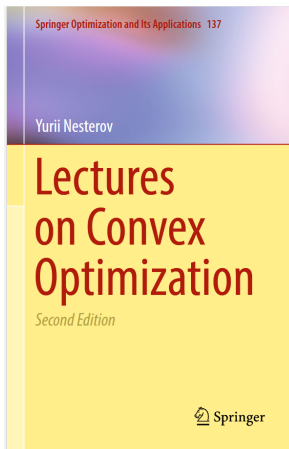
2 Guler 1993

- Exact Accelerated PPM
- Inexact accelerated PPM

3 Lin 2015

4 Paquette 2018

5 References



- Yurri Nesterov's book: "Lectures on Convex Optimization" 2018, Springer [1].



SIAM J. OPTIMIZATION
Vol. 2, No. 4, pp. 649–664, November 1992

© 1992 Society for Industrial and Applied Mathematics
007

NEW PROXIMAL POINT ALGORITHMS FOR CONVEX MINIMIZATION*

OSMAN GÜLER†

Abstract. This paper introduces two new proximal point algorithms for minimizing a proper, lower-semicontinuous convex function $f: \mathbf{R}^n \rightarrow \mathbf{R} \cup \{\infty\}$. Under this minimal assumption on f , the first algorithm possesses the global convergence rate estimate $f(x_k) - \min_{x \in \mathbf{R}^n} f(x) = O(1/(\sum_{j=0}^{k-1} \sqrt{\lambda_j})^2)$, where $\{\lambda_k\}_{k=0}^\infty$ are the proximal parameters. It is shown that this algorithm converges, and global convergence rate estimates for it are provided, even if minimizations are performed inexactly at each iteration. Both algorithms converge even if f has no minimizers or is unbounded from below. These algorithms and results are valid in infinite-dimensional Hilbert spaces.

Key words. proximal point algorithms, global convergence rates, augmented Lagrangian algorithms, convex programming

AMS(MOS) subject classifications. primary 90C25; secondary 49D45, 49D37

- Osman Guler's, "New proximal point algorithm for convex optimization", SIAM J. Optimization 1992. [2]

A Universal Catalyst for First-Order Optimization

Hongzhou Lin¹, Julien Mairal¹ and Zaid Harchaoui^{1,2}
¹Inria ²NYU
{hongzhou.lin, julien.mairal}@inria.fr
zaid.harchaoui@nyu.edu

Abstract

We introduce a generic scheme for accelerating first-order optimization methods in the sense of Nesterov, which builds upon a new analysis of the accelerated proximal point algorithm. Our approach consists of minimizing a convex objective by approximately solving a sequence of well-chosen auxiliary problems, leading to faster convergence. This strategy applies to a large class of algorithms, including gradient descent, block coordinate descent, SAG, SAGA, SDCA, SVRG, Finito/MISO, and their proximal variants. For all of these methods, we provide acceleration and explicit support for non-strongly convex objectives. In addition to theoretical speed-up, we also show that acceleration is useful in practice, especially for ill-conditioned problems where we measure significant improvements.

(a) Lin 2015

Catalyst Acceleration for Gradient-Based Non-Convex Optimization

Courtney Paquette Hongzhou Lin Dmitry Drusvyatskiy
University of Waterloo MIT University of Washington
c2paquette@uwaterloo.ca hongzhou@mit.edu ddrusv@uw.edu

Julien Mairal Zaid Harchaoui
Inria* University of Washington
julien.mairal@inria.fr zaid@uw.edu

January 3, 2019

Abstract

We introduce a generic scheme to solve nonconvex optimization problems using gradient-based algorithms originally designed for minimizing convex functions. Even though these methods may originally require convexity to operate, the proposed approach allows one to use them on weakly convex objectives, which covers a large class of non-convex functions typically appearing in machine learning and signal processing. In general, the scheme is guaranteed to produce a stationary point with a worst-case efficiency typical of first-order methods, and when the objective turns out to be convex, it automatically accelerates in the sense of Nesterov and achieves near-optimal convergence rate in function values. These properties are achieved without assuming any knowledge about the convexity of the objective, by automatically adapting to the unknown weak convexity constant. We conclude the paper by showing promising experimental results obtained by applying our approach to incremental algorithms such as SVRG and SAGA for sparse matrix factorization and for learning neural networks.

(b) Paquette 2018

- Honzhou Lin et al. “Universal Catalyst for first order optimization” 2015 JLMR [3].
- Paquette et al. “Catalyst for gradient-based nonconvex optimization” 2018 JLMR [4].

Objectives of the Talk

List of objectives

- 1 Introduce the technique of Nesterov's estimating sequence for convergence proof of algorithms.
- 2 Understand the historical context for the inspirations of the Catalyst algorithm.
- 3 Understand the theories behind the Catalyst meta acceleration.
- 4 Understand key innovations for controlling the errors in Catalyst accelerations.
- 5 Introduce the Non-convex extension of the method.

A note on the scope

Specific applications and algorithms are outside of the scope because variance reduced stochastic method is itself a big topic.

Objectives of the Talk

List of objectives

- 1 Introduce the technique of Nesterov's estimating sequence for convergence proof of algorithms.
- 2 Understand the historical context for the inspirations of the Catalyst algorithm.
- 3 Understand the theories behind the Catalyst meta acceleration.
- 4 Understand key innovations for controlling the errors in Catalyst accelerations.
- 5 Introduce the Non-convex extension of the method.

A note on the scope

Specific applications and algorithms are outside of the scope because variance reduced stochastic method is itself a big topic.

Nesterov's Estimating Sequence

Definition (Nesterov's estimating sequence)

Let $(\phi_k : \mathbb{R}^n \mapsto \mathbb{R})_{k \geq 0}$ be a sequence of functions. We call this sequence of function a Nesterov's estimating sequence when it satisfies the conditions that:

- 1 There exists another sequence $(x_k)_{k \geq 0}$ such that for all $k \geq 0$ it has $F(x_k) \leq \phi_k^*$.
- 2 There exists a sequence of $(\alpha_k)_{k \geq 0}$ such that for all $x \in \mathbb{R}^n$, $\phi_{k+1}(x) - \phi_k(x) \leq -\alpha_k(\phi_k(x) - F(x))$.

Nesterov's Estimating Sequence and Convergence

Observations

If we define $\phi_k, \Delta_k(x) := \phi_k(x) - F(x)$ for all $x \in \mathbb{R}^n$ and assume that F has minimizer x^* . Then observe that $\forall k \geq 0$:

$$\begin{aligned}\phi_{k+1}(x) - \phi_k(x) &\leq -\alpha_k(\phi_k(x) - F(x)) \\ \iff \phi_{k+1}(x) - F(x) - (\phi_k(x) - F(x)) &\leq -\alpha_k(\phi_k(x) - F(x)) \\ \iff \Delta_{k+1}(x) - \Delta_k(x) &\leq -\alpha_k \Delta_k(x) \\ \iff \Delta_{k+1}(x) &\leq (1 - \alpha_k) \Delta_k(x).\end{aligned}$$

Unroll the recurrence, by setting $x = x^*$, $\Delta_k(x^*)$ is non-negative and using the property of Nesterov's estimating sequence it gives:

$$\begin{aligned}F(x_k) - F(x^*) &\leq \phi_k^* - F(x^*) \leq \Delta_k(x^*) = \phi_k(x^*) - F(x^*) \\ &\leq \left(\prod_{i=0}^k (1 - \alpha_i) \right) \Delta_0(x^*).\end{aligned}$$

Example: accelerated proximal gradient

Quick Notations

Assume that: $F = f + g$ where f is L -Lipschitz smooth and $\mu \geq 0$ strongly convex and g is convex. Define

$$\mathcal{M}^{L^{-1}}(x; y) := g(x) + f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2} \|x - y\|^2,$$

$$\tilde{\mathcal{J}}_{L^{-1}} y := \operatorname{argmin}_x \mathcal{M}^{L^{-1}}(x; y),$$

$$\mathcal{G}_{L^{-1}}(y) := L \left(I - \tilde{\mathcal{J}}_{L^{-1}} \right) y.$$

Example: accelerated proximal gradient

Definition (Accelerated proximal gradient estimating sequence)

Define $(\phi_k)_{k \geq 0}$ be the Nesterov's estimating sequence recursively given by:

$$l_F(x; y_k) := F\left(\tilde{\mathcal{J}}_{L^{-1}} y_k\right) + \langle \mathcal{G}_{L^{-1}} y_k, x - y_k \rangle + \frac{1}{2L} \|\mathcal{G}_{L^{-1}} y_k\|^2,$$
$$\phi_{k+1}(x) := (1 - \alpha_k) \phi_k(x) + \alpha_k \left(l_F(x; y_k) + \frac{\mu}{2} \|x - y_k\|^2 \right).$$

The Algorithm generates a sequence of vectors y_k, x_k , and scalars α_k satisfies the following:

$$x_{k+1} = \tilde{\mathcal{J}}_{L^{-1}} y_k,$$
$$\text{find } \alpha_{k+1} \in (0, 1) : \alpha_{k+1} = (1 - \alpha_{k+1}) \alpha_k^2 + (\mu/L) \alpha_{k+1}$$
$$y_{k+1} = x_{k+1} + \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}} (x_{k+1} - x_k).$$

One of the possible base case can be $x_0 = y_0$ and any $\alpha_0 \in (0, 1)$.

Accelerated proximal point method

Guler in 1993 discovered the following:

- 1 The method of proximal point can be accelerated via Nesterov's estimating sequence.
- 2 The accelerated convergence rate retains for certain magnitude of errors on inexact evaluation of proximal point method.

Quick notations

We use the following list of notations:

$$\begin{aligned}\mathcal{M}^\lambda(x; y) &:= F(x) + \frac{1}{2\lambda} \|x - y\|^2 \\ \mathcal{J}_\lambda y &:= \operatorname{argmin}_x \mathcal{M}^\lambda(x; y) \\ \mathcal{G}_\lambda &:= \lambda^{-1}(I - \mathcal{J}_\lambda).\end{aligned}$$

We use $\mathcal{G}_k, \mathcal{J}_k, \mathcal{M}_k$ as a short for $\mathcal{G}_{\lambda_k}, \mathcal{J}_{\lambda_k}, \mathcal{M}_{\lambda_k}$. $(\lambda_k)_{k \geq 0}$ is a sequence that controls proximal operator.

Estimating sequence of accelerated PPM

Definition (Accelerated PPM estimating sequence)

$(\phi_k)_{k \geq 0}$ has for all $k \geq 0$, any $A \geq 0$:

$$\phi_0 := f(x_0) + \frac{A}{2} \|x - x_0\|^2,$$

$$\phi_{k+1}(x) := (1 - \alpha_k)\phi_k(x) + \alpha_k(\textcolor{red}{F}(\mathcal{J}_k y_k) + \langle \textcolor{red}{G}_k y_k, x - \mathcal{J}_k y_k \rangle).$$

$(\lambda_k)_{k \geq 0}$, $x_k = \mathcal{J}_{\lambda} y_k$. auxiliary vectors (y_k, v_k) , and $(\alpha_k, A_k)_{k \geq 0}$ satisfies $k \geq 0$:

$$\alpha_k = \frac{1}{2} \left(\sqrt{(A_k \lambda_k)^2 + 4A_k \lambda_k} - A_k \lambda_k \right)$$

$$y_k = (1 - \alpha_k)x_k + \alpha_k v_k$$

$$v_{k+1} = v_k - \frac{\alpha_k}{A_{k+1} \lambda_k} (y_k - \mathcal{J}_k y_k)$$

$$A_{k+1} = (1 - \alpha_k)A_k.$$

Convergence of accelerated PPM

An accelerated rate

The accelerated PPM generate $(x_k)_{k \geq 0}$ such that $F(x_k) - F^*$ converges at a rate of:

$$\mathcal{O} \left(\frac{1}{\left(\sum_{i=1}^k \sqrt{\lambda_i} \right)^2} \right).$$

Note, PPM without accelerate converges at a rate of $\mathcal{O}((\sum_{i=1}^k \lambda_i)^{-1})$.

Accelerated Inexact PPM

Guler cited Rockafellar 1976 [5] for condition (A'):

$$\begin{aligned}x_{k+1} \approx \mathcal{J}_k y_k \text{ be such that: } \text{dist} \left(\mathbf{0}, \partial \mathcal{M}^k(x_{k+1}; y_k) \right) &\leq \frac{\epsilon_k}{k} \\ \implies \|x_{k+1} - \mathcal{J}_k y_k\| &\leq \epsilon_k.\end{aligned}$$

Putting things into the context of accelerated PPM, the thereoem follows is pivotal:

Theorem (Guler's inexact proximal point error bound (Lemma 3.1))

Define the minimum of the Moreau Enevelope: $\mathcal{M}_k^ := \min_z \mathcal{M}^{\lambda_k}(z; y_k)$. If x_{k+1} is an inexact evaluation under condition (A'), then the estimating sequence admits the conditions that:*

$$\frac{1}{2\lambda_k} \|x_{k+1} - \mathcal{J}_k y_k\|^2 = \mathcal{M}_k(x_{k+1}, y_k) - \mathcal{M}_k^* \leq \frac{\epsilon_k^2}{2\lambda_k}.$$

Guler's Major Results

Theorem (Guler's accelerated inexact PPM convergence (Theorem 3.3))

If the error sequence $(\epsilon_k)_{k \geq 0}$ for condition A' is bounded by $\mathcal{O}(1/k^\sigma)$ for some $\sigma > 1/2$, then the accelerated proximal point method has for any feasible $x \in \mathbb{R}^n$:

$$f(x_k) - f(x) \leq \mathcal{O}(1/k^2) + (1/k^{2\sigma-1}) \rightarrow 0.$$

If $\sigma \geq 3/2$, the method converges at a rate of $\mathcal{O}(1/k^2)$. It looks exciting but it's not because:

- 1 Determining $(\epsilon_k)_{k \geq 0}$ requires knowledge on ϕ_k^* .
- 2 ϕ_k^* is expressed with untracable quantity: $F(\mathcal{J}_k y_k)$.

So the algorithm contains intractable quantities: $F(\mathcal{J}_k y_k)$.

Guler's Major Results

Theorem (Guler's accelerated inexact PPM convergence (Theorem 3.3))

If the error sequence $(\epsilon_k)_{k \geq 0}$ for condition A' is bounded by $\mathcal{O}(1/k^\sigma)$ for some $\sigma > 1/2$, then the accelerated proximal point method has for any feasible $x \in \mathbb{R}^n$:

$$f(x_k) - f(x) \leq \mathcal{O}(1/k^2) + (1/k^{2\sigma-1}) \rightarrow 0.$$

If $\sigma \geq 3/2$, the method converges at a rate of $\mathcal{O}(1/k^2)$. It looks exciting but it's not because:

- 1 Determining $(\epsilon_k)_{k \geq 0}$ requires knowledge on ϕ_k^* .
- 2 ϕ_k^* is expressed with untracable quantity: $F(\mathcal{J}_k y_k)$.

So the algorithm contains intractable quantities: $F(\mathcal{J}_k y_k)$.

Hongzhou Lin 2015 [3] did the following:

- ① Improved the proof from Guler 1993 to include strong convexity objective.
- ② Showed that $(\epsilon_k)_{k \geq 0}$ can be determined algorithmically and an accelerated rate can be achieved.
- ③ Invented his own accelerated varianced reduced incremental method called: “Prox MISO” to demonstrate the Catalyst Framework.

Quick notations

References



Y. Nesterov, *Lectures on Convex Optimization*, ser. Springer Optimization and Its Applications. Cham: Springer International Publishing, 2018, vol. 137. [Online]. Available: <http://link.springer.com/10.1007/978-3-319-91578-4>



O. Güler, "New Proximal Point Algorithms for Convex Minimization," *SIAM Journal on Optimization*, vol. 2, no. 4, pp. 649–664, Nov. 1992, publisher: Society for Industrial and Applied Mathematics. [Online]. Available: <https://epubs.siam.org/doi/10.1137/0802032>



H. Lin, J. Mairal, and Z. Harchaoui, "A Universal Catalyst for First-Order Optimization." MIT Press, Dec. 2015, p. 3384. [Online]. Available: <https://inria.hal.science/hal-01160728>



C. Paquette, H. Lin, D. Drusvyatskiy, J. Mairal, and Z. Harchaoui, "Catalyst for Gradient-based Nonconvex Optimization," in *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*. PMLR, Mar. 2018, pp. 613–622, iSSN: 2640-3498. [Online]. Available: <https://proceedings.mlr.press/v84/paquette18a.html>



R. T. Rockafellar, "Monotone operators and the proximal point algorithm," vol. 14, no. 5, pp. 877–898. [Online]. Available: <http://epubs.siam.org/doi/10.1137/0314056>