

Catalyst Meta Acceleration Framework: The Gist of its Theories

Hongda Li

November 8, 2024

Abstract

Nesterov’s accelerated gradient first appeared back in the 1983 has sparked numerous theoretical and practical advancements in Mathematics programming literatures. The idea behind Nesterov’s acceleration is universal in the convex case it has concrete extension in the non-convex case. In this paper we survey specifically the Catalyst Acceleration that incorporated ideas from the Accelerated Proximal Point Method proposed by Guler back in 1993. The paper reviews Nesterov’s classical analysis of accelerated gradient in the convex case. The paper will describe key aspects of the theoretical innovations involved to achieve the design of the algorithm in convex, and non-convex case.

2010 Mathematics Subject Classification: Primary 65K10, 90c25, 90C30; Secondary 65Y20. **Keywords:** Nesterov acceleration, Proximal point method.

1 Introduction

Nesterov first proposed the idea of an optimal algorithm named accelerated gradient descent method in his seminal work back in 1983 [7]. It was seminal at the time because the algorithm’s upper bound on the iteration complexity sealed the gap between the lower bound for all first order Lipschitz smooth convex function and the upper bound for this class of functions. For a specific definition of the class of algorithms that are considered “First Order”, we refer reader to Chapter 2 of Nesterov’s new book [8] for more information. In brief the method of gradient descent has an upper bound of $\mathcal{O}(1/k)$ in iteration complexity. It doesn’t achieve the $\mathcal{O}(1/k^2)$ lower iteration complexity bound for first order optimization

algorithms. The method of accelerated gradient descent has an upper bound of $\mathcal{O}(1/k^2)$, making it optimal.

On first judgement, it's tempting to think that the existence of this optimal algorithm sealed the ceiling for the theoretical development for the entire class of convex first-order smooth optimization. The judgement is correct but lacks the nuance in understanding. The missing piece here is the fact that Nesterov's accelerated gradient is a system of analysis technique instead of any specific design patterns in algorithms.

To demonstrate, the introduction of Guler's works in 1993 [4] proposed an accelerated scheme using the technique of Nesterov's estimating sequence for Proximal Point Method (PPM) in the convex case. Let $(\lambda_k)_{k \geq 0}$ be the sequence of scalars used for regularizing the proximal point method which generates sequence $(x_k)_{k \geq 0}$ given any initial guess x_0 . Guler's prior work [3] showed that convergence of PPM method in the convex case has $\mathcal{O}(1/\sum_{i=1}^n \lambda_i)$. His new algorithm using the technique introduced in Nesterov's accelerated gradient achieves a convergence rate of $\mathcal{O}(1/(\sum_{i=1}^n \sqrt{\lambda_i})^2)$. In addition, he also proposed together an inexact Accelerated PPM method using conditions described in Rockafellar's works in 1976 [10].

One would be tempting to conclude that this has sealed the ceiling for research on the topic of extending Nesterov's acceleration. That is indeed correct, but not from a practical point of view. Let $F : \mathbb{R}^n \mapsto \mathbb{R}$ be our objective function, $\mathcal{J}_\lambda := (I + \lambda \partial F)^{-1}$ and $\mathcal{M}^\lambda(x; y) := F(x) + \frac{1}{2\lambda} \|x - y\|^2$ then the inexact proximal point considers with error ϵ_k has the following characterizations of inexactness as put forward by Guler [4]:

$$\begin{aligned} \tilde{x} &\approx \mathcal{J}_\lambda y \\ \text{dist}(\partial \mathcal{M}^\lambda(\tilde{x}; y)) &\leq \frac{\epsilon}{\lambda} \end{aligned}$$

However, this is troublesome because if we need to approximate the resolvent operator \mathcal{J}_λ , then it's probably difficult to compute the subgradient $\partial \mathcal{M}(\cdot; y)$, which make it difficult to know when we achieved the required exactness for a PPM evaluation. Otherwise, if we already know the subgradient well, then why approximate it in the first place?

Introduced in Lin et al. [5][6] is a series of papers on a concrete meta algorithm called Catalyst (It's called 4WD Catalyst for the non-convex extension in works by Paquette, Lin et al. [9]). It's called a meta algorithm because it uses other first order algorithm to evaluate inexact proximal point method and then performs the accelerated PPM using Nesterov's acceleration. Their innovations are tracking and controlling the errors made in the inexact PPM throughout the algorithm and some original example usages of the Catalyst framework.

One would be tempting to assert that this has sealed the ceiling for both theories and practice of Nesterov's acceleration hence it must be the center of discussion in this report. The conclusion is indeed correct which it will happen in the sections that follow while the assertion remains open.

1.1 Contributions

The writing is expository and won't contain major results. We reviewed the literatures and faithfully reproduced some claims, in addition we give insights into understanding the claim in relations to other papers and foundational ideas in optimization.

2 Preliminaries

Throughout the entire writing, let our ambient space is \mathbb{R}^n . We assume the optimization problem of:

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

In this section we introduce the idea of Nesterov's estimating sequence. Nesterov's estimating sequence is fundamental to works in Guler's accelerated PPM method, and Catalyst meta acceleration as a whole.

2.1 Method of Nesterov's Estimating Sequence

Definition 2.1 (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:*

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

Observation 2.2 *If we define ϕ_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 $\forall k \geq 0$:*

$$\begin{aligned} \Delta_k(x) &= \phi_k(x) - f(x) \geq \phi_k^* - f(x) \\ x = x_k &\implies \Delta_k(x_k) \geq \phi_k^* - f(x_k) \geq 0 \\ x = x_* &\implies \Delta_k(x_*) \geq \phi_k^* - f_* \geq f(x_k) - f_* \geq 0 \end{aligned}$$

The function $\Delta_k(x)$ is non-negative specifically at the points: x_, x_k . Additionally, we can*

derive the convergence rate of $\Delta_k(x^*)$ because $\forall x \in \mathbb{R}^n$:

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

Unrolling the above recursion it yields:

$$\Delta_{k+1}(x) \leq (1 - \alpha_k) \Delta_k(x) \leq \dots \leq \left(\prod_{i=0}^k (1 - \alpha_i) \right) \Delta_0(x).$$

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

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

Therefore, it yields a convergence of the sequence $f(x_k) \rightarrow f(x^*)$ with a rate relates to sequence $(\alpha_k)_{k \in \mathbb{N}}$.

Much of the analysis of convergence Nesterov's type accelerated gradient method inherit the idea of Nesterov's estimating sequence. Such a proof won't result in simple proof because the construction of ϕ_k is non-trivial, but it comes with the advantage too because we can put creativity into the construction of the estimating sequence $(\phi_k)_{k \geq 0}$.

3 Nesterov's Accelerated Proximal Gradient

This section swiftly exposes the constructions of the Nesterov's estimating sequence for the FISTA algorithm by Beck[2], which is specific case of Algorithm (2.2.63), in Nesterov's book [8]. Discussion on these algorithms are relevant because they share the same format as the Catalyst Acceleration framework and accelerated PPM.

Throughout this section we assume that: $F = f + g$ where f is L -Lipschitz smooth and $\mu \geq 0$ strongly convex and g is convex. Define

$$\begin{aligned}
\mathcal{M}^{L^{-1}}(x; y) &:= g(x) + f(y) + \langle \nabla f(x), x - y \rangle + \frac{L}{2} \|x - y\|^2, \\
\tilde{\mathcal{J}}_{L^{-1}} y &:= \underset{x}{\operatorname{argmin}} \mathcal{M}^{L^{-1}}(x; y), \\
\mathcal{G}_{L^{-1}}(y) &:= L \left(I - \tilde{\mathcal{J}}_{L^{-1}} \right) y.
\end{aligned}$$

In the literature, $\mathcal{G}_{L^{-1}}$ is commonly known as the gradient mapping. The definition follows, we define the Nesterov's estimating sequence used to derive the accelerated proximal gradient method.

Definition 3.1 (Accelerated proximal gradient estimating sequence) Define $(\phi_k)_{k \geq 0}$ be the Nesterov's estimating sequence recursively given by:

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

And the sequence of vector y_k, x_k , and scalars α_k satisfies the following:

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

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

Observation 3.2 One key component of the Nesterov's estimating sequence is the use of the proximal gradient inequality: $l_F(x; y_k) + \mu/2 \|x - y_k\|^2$. In the convex case the function has the property $l_F(\cdot, y) \leq F(\cdot)$ for all y . More precisely, if $f \equiv 0$ then $\tilde{\mathcal{J}}_{L^{-1}} y_k$ becomes resolvent $(I + L^{-1} \partial F)^{-1}$, which makes x_k being an exact evaluation of PPM. And we have

$$\begin{aligned} l_F(x; y_k) &= F(\mathcal{J}_{L^{-1}} y_k) + \langle L(y - \mathcal{J}_{L^{-1}} y), x - y_k \rangle + \frac{L}{2} \|y_k - \mathcal{J}_{L^{-1}} y_k\|^2 \\ &= F(\mathcal{J}_{L^{-1}} y_k) + \langle L(y - \mathcal{J}_{L^{-1}} y), x - \mathcal{J}_{L^{-1}} y_k \rangle. \end{aligned}$$

This is the proximal inequality. Observe that the inequality with proximal gradient term can be interpreted as an example of inexact evaluation of the PPM and the inequality.

To demonstrate the usage of Nesterov's estimating sequence here, consider sequence $(x_k)_{k \geq 0}$ such that $F(x_k) \leq \phi_k^*$. Assume the existence of minimizer x^* for F , by definition of ϕ_k let $x = x^*$ then $\forall k \geq 0$:

$$\begin{aligned} \phi_{k+1}(x^*) &= (1 - \alpha_k) \phi_k(x^*) + \alpha_k \left(l_F(x^*; y_k) + \frac{\mu}{2} \|x^* - y_k\|^2 \right) \\ \phi_{k+1}(x^*) - \phi_k(x^*) &= -\alpha_k \phi_k(x^*) + \alpha_k \left(l_F(x^*; y_k) + \frac{\mu}{2} \|x^* - y_k\|^2 \right) \\ \implies \phi_{k+1}(x^*) - F(x^*) + F(x^*) - \phi_k(x^*) &\leq -\alpha_k (\phi_k(x^*) - F(x^*)) \\ \implies F(x_{k+1}) - F(x^*) &\leq \phi_{k+1}^* - F(x^*) \leq \phi_{k+1}(x^*) - F(x^*) \leq (1 - \alpha_k) (\phi_k(x^*) - F(x^*)). \end{aligned}$$

On the first inequality we used the fact that $l_F(x; y_k) + \mu/2 \|x - y_k\|^2 \leq F(x)$. Unrolling the recurrence, we can get the convergence rate of $F(x_k) - F(x^*)$ to be on Big O of $\prod_{i=1}^k (1 - \alpha_i)$.

Remark 3.3 The definition is a generalization of Nesterov’s estimating sequence comes from (2.2.63) from Nesterov’s book [8]. Compare to Nesterov’s work, we used proximal gradient operator instead of projected gradient. The same inequality is called “Fundamental Proximal Gradient Inequality” in Amir Beck’s book [1], Theorem 10.16.

4 Guler 1993

This section introduces the setup of the Nesterov’s estimating sequence used in Guler’s accelerated Proximal Point method. In addition, this section will highlight some observations and theoretical results accordingly.

Definition 4.1 (Accelerated PPM estimating sequence)

5 Lin 2015

6 Non-convex Extension of Catalyst Acceleration

References

- [1] A. BECK, *First-order Methods in Optimization*, MOS-SIAM Series in Optimization, SIAM, israel, 2017.
- [2] A. BECK AND M. TEBoulLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*, SIAM Journal on Imaging Sciences, 2 (2009), pp. 183–202.
- [3] O. GULER, *On the convergence of the proximal point algorithm for convex minimization*, SIAM Journal on Control and Optimization, 29 (1991), p. 17. Num Pages: 17 Place: Philadelphia, United States Publisher: Society for Industrial and Applied Mathematics.
- [4] O. GÜLER, *New Proximal Point Algorithms for Convex Minimization*, SIAM Journal on Optimization, 2 (1992), pp. 649–664. Publisher: Society for Industrial and Applied Mathematics.
- [5] H. LIN, J. MAIRAL, AND Z. HARCHAOUI, *A Universal Catalyst for First-Order Optimization*, MIT Press, Dec. 2015, p. 3384.
- [6] —, *Catalyst Acceleration for First-order Convex Optimization: from Theory to Practice*, in Journal of Machine Learning Research, vol. 18, 2018, pp. 1–54.

- [7] Y. NESTEROV, *A method for solving the convex programming problem with convergence rate $O(1/k^2)$* , Proceedings of the USSR Academy of Sciences, (1983).
- [8] Y. NESTEROV, *Lectures on Convex Optimization*, vol. 137 of Springer Optimization and Its Applications, Springer International Publishing, Cham, 2018.
- [9] 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.
- [10] R. T. ROCKAFELLAR, *Monotone operators and the proximal point algorithm*, 14, pp. 877–898.

A Postponed proofs

A.1 Theorems and claims for accelerated proximal gradient

Theorem A.1 (Fundamental theorem of proximal gradient)

Theorem A.2 (Canonical form of proximal gradient estimating sequence)