

Adaptive Accelerated Gradient Converging Method under Hölderian Error Bound Condition

Mingrui Liu, Tianbao Yang

Department of Computer Science, The University of Iowa

Problem of Interest

Smooth Composite Optimization Problem:

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) \triangleq f(\mathbf{x}) + g(\mathbf{x}) \quad (1)$$

- $f(\mathbf{x})$: continuously differentiable convex with L -Lipschitz continuous gradient
- $g(\mathbf{x})$: proper lower semi-continuous convex
- proximal mapping: $P_g(\mathbf{u}) = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{x} - \mathbf{u}\|_2^2 + g(\mathbf{x})/L$.
- proximal gradient $G(\mathbf{x}) = L(\mathbf{x} - P_g(\mathbf{x} - \nabla f(\mathbf{x})/L))$.

Accelerated Proximal Gradient Methods:

- Nesterov's APG:

$$\mathbf{x}_{\tau+1} = P_g(\mathbf{y}_\tau - \nabla f(\mathbf{y}_\tau)/L), \mathbf{y}_{\tau+1} = \mathbf{x}_{\tau+1} + \beta_\tau(\mathbf{x}_{\tau+1} - \mathbf{x}_\tau).$$

$$\beta_\tau = \frac{\tau}{\tau+3}; \text{ iteration complexity (IC): } O(1/\sqrt{\epsilon})$$
- if $f(\mathbf{x})$ is α -strongly convex: $\beta_\tau = \frac{\sqrt{L}-\sqrt{\tau}}{\sqrt{L}-\sqrt{\tau}}$, IC: $O(\sqrt{L/\alpha} \log(1/\epsilon))$
- if $g(\mathbf{x})$ is α -strongly convex: Nesterov's ADG, same IC as above

Algorithm 1 ADG

```

1:  $\mathbf{x}_0 \in \Omega, A_0 = 0, \mathbf{v}_0 = \mathbf{x}_0$ 
2: for  $t = 0, \dots, T$  do
3:   Find  $a_{t+1}$  from quadratic equation  $\frac{a^2}{A_t+a} = 2^{1+\alpha A_t}/L$ 
4:   Set  $A_{t+1} = A_t + a_{t+1}$ 
5:   Set  $\mathbf{y}_t = \frac{A_t}{A_{t+1}}\mathbf{x}_t + \frac{a_{t+1}}{A_{t+1}}\mathbf{v}_t$ 
6:   Compute  $\mathbf{x}_{t+1} = P_g(\mathbf{y}_t - \nabla f(\mathbf{y}_t)/L)$ 
7:   Compute  $\mathbf{v}_{t+1} = \arg \min_{\mathbf{x}} \sum_{\tau=1}^{t+1} a_\tau \nabla f(\mathbf{x}_\tau)^\top \mathbf{x} + A_{t+1} g(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2$ 
8: end for

```

Recent Advances

Linear Convergence under weaker conditions, e.g., quadratic error bound condition (QEB, or quadratic growth condition).

$$\text{dist}(\mathbf{x}, \Omega_*) \leq c(F(\mathbf{x}) - F_*)^{1/2},$$

Examples: LASSO

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{n} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

- PG [4]: IC = $O(Lc^2 \log(1/\epsilon))$
- a restarting version of APG [3]: IC = $O(\sqrt{Lc^2} \log(1/\epsilon))$ provided that the value of c is known.
- Issue: the value of c is usually unknown

Questions: how to develop algorithms with improved IC when c is unknown? What if considering a weaker condition than QEB?

Hölderian Error Bound (HEB)

Definition 1: A function $F(\mathbf{x})$ is said to satisfy a Hölderian error bound condition on ξ -sublevel set if there exist $\theta \in (0, 1]$ and $0 < c < \infty$ such that for any $\mathbf{x} \in \mathcal{S}_\xi$,

$$\text{dist}(\mathbf{x}, \Omega_*) \leq c(F(\mathbf{x}) - F_*)^\theta,$$

where Ω_* denotes the set of optimal solution.

- closely related to the Kurdyka - Łojasiewicz (KL) inequality in real algebraic geometry.
- when functions are semi-algebraic and continuous, the above inequality is known to hold on any compact set
- $\theta = 0$ can be considered as a special case

PG and restarting APG under HEB

Algorithm 2 PG under HEB

```

1: Input:  $\mathbf{x}_0 \in \Omega$  such that  $F(\mathbf{x}_0) - F_*) \leq \epsilon_0$ 
2: for  $\tau = 1, \dots, t$  do
3:    $\mathbf{x}_{\tau+1} = P_g(\mathbf{x}_\tau - \nabla f(\mathbf{x}_\tau)/L)$ 
4: end for
5: Option I: return  $\mathbf{x}_{t+1}$ 
6: Option II: return  $\mathbf{x}_k$  s.t.  $G(\mathbf{x}_k) = \min_{\tau} \|G(\mathbf{x}_\tau)\|_2$ 

```

- Option I for achieving $F(\mathbf{x}_t) - F_* \leq \epsilon$,
- IC is $O(\max\{\frac{Lc^2}{\epsilon^{1-2\theta}}, Lc^2 \log(\frac{\epsilon_0}{\epsilon})\})$ if $\theta \leq 1/2$, and $O(Lc^2 \epsilon_0^{2\theta-1})$ otherwise.
- Option II is for achieving $G(\mathbf{x}_k) \leq \epsilon$
- IC is $O\left(Lc^{\frac{1}{1-\theta}} \max\left\{\frac{1}{\epsilon^{1-2\theta}}, \log\left(\frac{\epsilon_0}{\epsilon}\right)\right\}\right)$ if $\theta \leq 1/2$, and is $O(c^2 L \epsilon_0^{2\theta-1})$ otherwise.

Adaptive Accelerated Gradient Converging Methods (adaAGC)

Algorithm 4 adaAGC

```

1: Input:  $\mathbf{x}_0 \in \Omega$  and  $c_0$  and  $\gamma > 1$ 
2: Let  $c_e = c_0$  and  $\epsilon_0 = \|G(\mathbf{x}_0)\|_2$ ,
3: for  $k = 1, \dots, K$  do
4:   Let  $\delta_k$  be given in (2) and  $g_{\delta_k}(\mathbf{x}) = g(\mathbf{x}) + \frac{\delta_k}{2} \|\mathbf{x} - \mathbf{x}_{k-1}\|_2^2$ 
5:    $\mathbf{x}_k^k = \mathbf{x}_{k-1}$  and  $\mathbf{y}_k = \mathbf{x}_{k-1}$ 
6:   for  $s = 1, \dots$  do
7:     for  $\tau = 1, \dots$  do
8:       Apply one step of ADG to  $f(\mathbf{x}) + g_{\delta_k}(\mathbf{x})$  to generate  $\mathbf{x}_{\tau+1}^k$ 
9:       if  $\|G(\mathbf{x}_{\tau+1}^k)\|_2 \leq \epsilon_{k-1}/2$  then
10:        let  $\mathbf{x}_k = \mathbf{x}_{\tau+1}^k$  and  $\epsilon_k = \epsilon_{k-1}/2$ .
11:        break the two enclosing for loops
12:      else if  $\tau = \left\lceil 2\sqrt{\frac{L+\delta_k}{\delta_k}} \log \frac{\sqrt{L(L+\delta_k)}}{\delta_k} \right\rceil$  then
13:        let  $c_e = \gamma c_e$  and break the enclosing for loop
14:      end if
15:    end for
16:  end for
17: end for
18: Output:  $\mathbf{x}_K$ 

```

Main Theorem: Suppose $F(\mathbf{x}_0) - F_* \leq \epsilon_0$, $F(\mathbf{x})$ satisfies HEB on \mathcal{S}_{ϵ_0} and $c_0 \leq c$. Let $\epsilon_0 = \|G(\mathbf{x}_0)\|_2$, $K = \lceil \log_2(\frac{\epsilon_0}{\epsilon}) \rceil$, $p = (1-2\theta)/(1-\theta)$ for $\theta \in (0, 1/2]$. The IC of adaAGC for having $\|G(\mathbf{x}_K)\|_2 \leq \epsilon$ is (where $\tilde{O}(\cdot)$ suppresses a log term depending on c, c_0, L, γ)

$$\text{IC} = \begin{cases} \tilde{O}\left(\sqrt{L}c^{\frac{1}{2(1-\theta)}} \max\left(\frac{1}{\epsilon^{1-\theta}}, \log(\epsilon_0/\epsilon)\right)\right) & \text{if } \theta \in (0, 1/2] \\ \tilde{O}\left(\sqrt{L}c \log(\epsilon_0/\epsilon)\right) & \text{if } \theta = 1/2 \\ \tilde{O}\left(\sqrt{L}c \epsilon_0^{\theta-1/2}\right) & \theta \in (1/2, 1] \end{cases}$$

Table 1: Summary of iteration complexities in this work under the HEB condition with $\theta \in (0, 1/2]$, where $G(\mathbf{x})$ denotes the proximal gradient, $C(1/\epsilon^\alpha) = \max(1/\epsilon^\alpha, \log(1/\epsilon))$ and $\tilde{O}(\cdot)$ suppresses a logarithmic term. If $\theta > 1/2$, all algorithms can converge with finite steps of proximal mapping. rAPG stands for restarting APG. * mark results available for certain subclasses of problems.

algo.	PG	rAPG	adaAGC
$F(\mathbf{x}) - F_* \leq \epsilon$	$O\left(c^2 LC\left(\frac{1}{\epsilon^{1-2\theta}}\right)\right)$	$O\left(c\sqrt{LC}\left(\frac{1}{\epsilon^{1/2-\theta}}\right)\right)$	$O\left(c\sqrt{LC}\left(\frac{1}{\epsilon^{1/2-\theta}}\right)\right)^*$
$\ G(\mathbf{x})\ _2 \leq \epsilon$	$O\left(c^{\frac{1}{1-\theta}} LC\left(\frac{1}{\epsilon^{1-2\theta}}\right)\right)$	-	$\tilde{O}\left(c^{\frac{1}{2(1-\theta)}} \sqrt{LC}\left(\frac{1}{\epsilon^{1-2\theta}}\right)\right)$
requires θ	No	Yes	Yes
requires c	No	Yes	No

Applications and Experimental Results

Applications: Consider the regularized problems with a smooth loss:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}^\top \mathbf{a}_i, b_i) + \lambda R(\mathbf{x}), \quad (3)$$

where (\mathbf{a}_i, b_i) , $i = 1, \dots, n$ denote a set of training examples, $R(\mathbf{x})$ is the regularizer.

Examples of $\theta = 1/2$: piecewise quadratic convex function [1]

- square loss, squared hinge loss, huber loss
- ℓ_1, ℓ_∞ norm, Huber norm and $\ell_{1,\infty}$ norm.

Examples of $\theta = 1/2$: structured smooth composite functions

- $f(\mathbf{x}) = h(A\mathbf{x})$: h is smooth and strongly convex on any compact set
- $g(\mathbf{x})$ is a polyhedral function (e.g., ℓ_1 norm)

Examples of $\theta < 1/2$: ℓ_1 constrained ℓ_p regression: $\theta = 1/p$

$$\min_{\|\mathbf{x}\|_1 \leq s} F(\mathbf{x}) \triangleq \frac{1}{n} \sum_{i=1}^n (\mathbf{a}_i^\top \mathbf{x} - b_i)^p, \quad p \in 2\mathbb{N}. \quad (4)$$

Experimental Results:

Table 2: squared hinge loss with ℓ_1 norm (left) and ℓ_∞ norm (right) regularization on splice data

Algorithm	$\epsilon = 10^{-4}$	$\epsilon = 10^{-5}$	$\epsilon = 10^{-6}$	$\epsilon = 10^{-7}$	$\epsilon = 10^{-4}$	$\epsilon = 10^{-5}$	$\epsilon = 10^{-6}$	$\epsilon = 10^{-7}$
PG	2040	2040	2040	2040	3514	3724	3724	3724
FISTA	1289	1289	1289	1289	5526	5526	5526	5526
urFISTA	1666	2371	2601	3480	1674	2379	2605	3488
adaAGC	1410	1410	1410	2382	2382	2382	2382	2382
					FISTA > adaAGC > PG > urFISTA	adaAGC > urFISTA > PG > FISTA		

Table 3: square loss with ℓ_1 norm (left) and ℓ_∞ norm (right) regularization on cpusmall data

Algorithm	$\epsilon = 10^{-4}$	$\epsilon = 10^{-5}$	$\epsilon = 10^{-6}$	$\epsilon = 10^{-7}$	$\epsilon = 10^{-4}$	$\epsilon = 10^{-5}$	$\epsilon = 10^{-6}$	$\epsilon = 10^{-7}$
PG	109298	159908	170915	170915	139505	204120	210874	210874
FISTA	6781	16387	23779	23779	6610	16418	20082	20082
urFISTA	18278	26706	35173	43603	18276	26704	35169	43601
adaAGC	9571	12623	13575	9881	13033	13632	13632	13632
					adaAGC > FISTA > urFISTA > PG	adaAGC > urFISTA > PG > FISTA		

Table 4: ℓ_1 regularized huber loss (left) and $\ell_$