The Classic Optimization Problem: $\operatorname{argmin}_{x \in X} f(x)$

Simple first order algorithm: $x^{k} \leftarrow y^{k-1} + \eta \nabla f(y^{k-1}), y^{k} = x^{k} + \alpha (x^{k} - x^{k-1})$



Strongly convex $f: \alpha = \frac{\sqrt{1/\mu\eta} - 1}{\sqrt{1/\mu\eta} + 1}$ General convex $f: \alpha = \frac{k-1}{k+1}$

Iteration complexity:

Mthods	Error After k Iters	Assu
VGD	$\exp(-k/\kappa)$	Strong
VGD	1/k	Сс
NAG	$\exp(-k/\sqrt{\kappa})$	Strong
NAG	$1/k^2$	Сс

*conditioning number $\kappa = L/\mu$, where μ is the strong convexity constant of f and L is the Lipchitz constant

How to understand these algorithms via an intuitive way?

Our answer: Study their continuous time limit and the physical systems corresponding to these algorithms!



The iterates of the algorithm The continuous-time limit (step_size $\rightarrow 0$) of iterates: a curve described by an ordinary differential equation (ODE)

The Physical Systems Behind Optimization Algorithms

Lin F. Yang¹, Raman Arora², Vladimir Braverman², and Tuo Zhao³ ¹ Princeton University, ² Johns Hopkins University, ³ GaTech

How to connect discrete iterates to continuous iterates?

□ Set (1)
$$k(t)$$
: = $\left|\frac{t}{\sqrt{\eta}}\right|$ or (2) $k(t)$: =
□ [1] chooses (1) for AGD and choose

Our Solution: a unified framework for time-scaling Set $k(t) := \left| \frac{t}{h} \right|$ for all algorithms. Let the algorithm tell us what is h!

Taylor expansion the algorithm iterates:

 $(x^{(k+1)} - x^{(k)}) =$ $(x^{(k)} - x^{(k-1)}) =$ and $\eta \nabla f \left[x^{(k)} + \alpha \left(x^{(k)} \right) \right]$

 $m\ddot{X}(t) +$

A unified ODE that describe a damped oscillator system:

m



NAG: massive system

 \Rightarrow

uptions

ly convex

onvex

ly Convex

onvex



Damping coefficient: c

 $\left|\frac{t}{n}\right|$ and define $X(t) = x^{k(t)}$ ses (2) for NAG

$$\begin{aligned} &\dot{X}(t)h + \frac{1}{2}\ddot{X}(t)h^2 + o(h), \\ &\dot{X}(t)h - \frac{1}{2}\ddot{X}(t)h^2 + o(h), \\ &\dot{X}^{(k)} - x^{(k-1)}\Big] = \eta \nabla f(X(t)) + O(\eta h). \end{aligned}$$

$$c\dot{X}(t) + \nabla f(X(t)) = 0.$$

the *particle* mass, as the damping coefficient, the *potential field*. as

A unified view of the choice of *h*:

VGD: massless system $\dot{X} + \nabla f(X) = 0$ $\Rightarrow m = 0, c = 1$ $\Rightarrow h = \Theta(\eta)$

$$m\ddot{X} + c\dot{X} + \nabla f(X) = 0$$

$$h = \Theta(\sqrt{\eta}), 1 - \alpha = \Theta(h)$$

Insights:

- the ODE
- damping

The energy decreasing of the system:

 $\mathcal{E}(t) \propto e$

Our framework naturally extends to the PL-condition *f* :

PL-condition: $\Box K = L/\mu$

Our framework naturally extends to other optimization methods:

□ Randomized accelerated coordinate gradient descent (ARCG)

[1] / [2] /Ours	VC
General Convex	/
Strongly Convex	/
Proximal Variants	/
PL Condition	/
Physical Systems	/

References:

• [1] Su, W., Boyd, S. and Candes, E (2014);

The energy of the physical system is a Lyapunov function of

• Energy decreases fastest when the system is under critical

Consider quadratic function $f(x) = K ||x - x^*||^2/2$

$$\exp\Big(-\frac{1}{2}\Big[\frac{c}{m}-\sqrt{\frac{c^2}{m^2}-\frac{4\mathcal{K}}{m}}\Big]t\Big).$$

Critical damping: $c^2 = 4mK$, which corresponds to NAG



Randomized coordinate gradient descent (RCGD) $x_{i}^{(k)} = x_{i}^{(k-1)} - \eta \nabla_{j} f(x^{(k-1)})$ and $x_{j}^{(k)} = x_{j}^{(k-1)}$

$$x_{j}^{(k)} = y_{j}^{(k-1)} - \eta \nabla_{j} f(y^{(k-1)}), \quad x_{\backslash j}^{(k)} = y_{\backslash j}^{(k-1)}, \text{ and } y^{(k)} = x^{(k)} + \alpha \left(x^{(k)} - x^{(k-1)}\right)$$

$$\frac{1}{\left[2\right] / \text{Ours}} \quad V\text{GD} \quad \text{NAG} \quad \text{RCGD} \quad \text{ARCG} \quad \text{Newton}$$

$$\frac{1}{\text{neral Convex}} \quad \frac{1}{-\sqrt{--1}} \quad \forall \sqrt{-1} \quad \forall \sqrt{-1} \quad \sqrt{-1} \quad \frac{1}{-\sqrt{--1}} \quad \sqrt{-1} \quad \sqrt{-1}$$