

ODE-based Learning to Optimize

Zhonglin Xie

Beijing International Center for Mathematical Research
Peking University

Joint work with Zaiwen Wen, Wotao Yin

June 5, 2024

Two important questions

- ▶ How to translate the fast convergence properties of ODEs to algorithms?

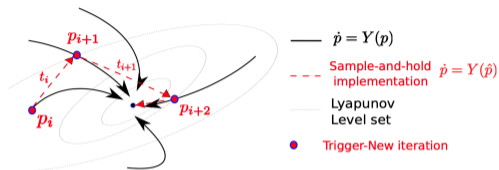
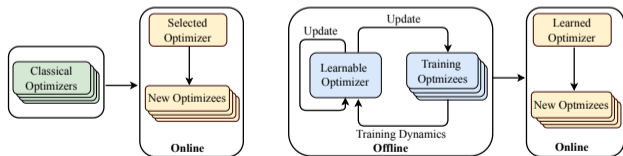


Figure: Rate-matching discretization

Combine **error analysis** in ODE and **complexity analysis** in optimization

- ▶ How to select the best coefficients for (ISHD)?



A learning to optimize framework with **theoretical guarantee**

Figure: Learning to optimize

A continuous-time viewpoint of acceleration methods: $\min f(x)$

- ▶ Gradient descent method corresponds to gradient flow

$$x_{k+1} = x_k - \sqrt{s} \nabla f(x_k) \quad \Leftrightarrow \quad \dot{x}(t) = -\nabla f(x(t))$$

- ▶ Nesterov accelerated gradient method corresponds to

$$\begin{cases} x_k = y_{k-1} - s \nabla f(y_{k-1}) \\ y_k = x_k + \frac{k-1}{k+2} (x_k - x_{k-1}) \end{cases} \Leftrightarrow \begin{cases} \ddot{x}(t) + \frac{3}{t} \dot{x}(t) + \sqrt{s} \nabla^2 f(x(t)) \dot{x}(t) \\ + \left(1 + \frac{3\sqrt{s}}{2t}\right) \nabla f(x(t)) = 0 \end{cases}$$

- ▶ Inertial system with Hessian-driven damping

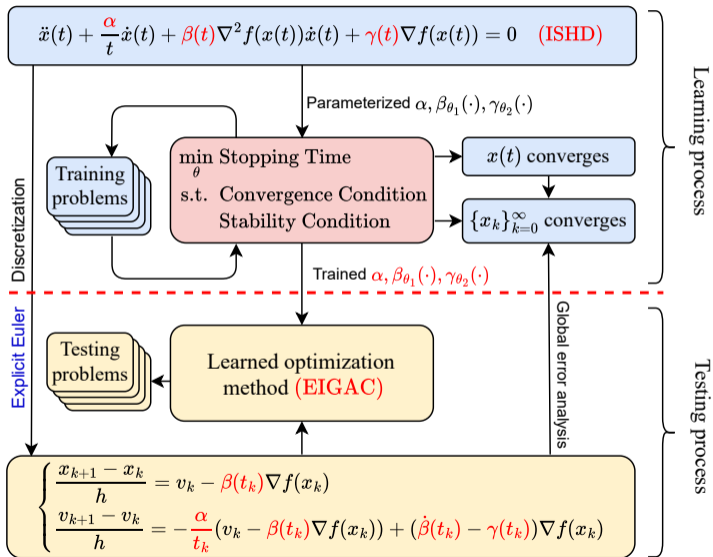
$$\ddot{x}(t) + \frac{\alpha}{t} \dot{x}(t) + \beta(t) \nabla^2 f(x(t)) \dot{x}(t) + \gamma(t) \nabla f(x(t)) = 0 \quad (\text{ISHD})$$

- ▶ Let $w(t) = \gamma(t) - \dot{\beta}(t) - \beta(t)/t$. Convergence condition for (ISHD)

$$\gamma(t) > \dot{\beta}(t) + \frac{\beta(t)}{t}, \quad t\dot{w}(t) \leq (\alpha - 3)w(t), \quad \text{for all } t \geq t_0$$

Convergence rate: $f(x(t)) - f_\star = \mathcal{O}(1/(t^2 w(t)))$

Our training and testing framework



Outline

- 1 Conditions for stability-preserving discretization
- 2 Select the best coefficients using learning to optimize
- 3 Computation of the conservative gradients
- 4 Convergence analysis
- 5 Numerical results

A fundamental result: an enhanced convergent condition for ISHD

Theorem 1

Given $\kappa \in (0, 1]$, $\lambda \in (0, \alpha - 1]$ and f is twice differentiable convex

$$\begin{aligned} \delta(t) &= t^2(\gamma(t) - \kappa\dot{\beta}(t) - \kappa\beta(t)/t) + (\kappa(\alpha - 1 - \lambda) - \lambda(1 - \kappa))t\beta(t), \\ w(t) &= \gamma(t) - \dot{\beta}(t) - \beta(t)/t, \quad \delta(t) > 0, \quad \text{and} \quad \dot{\delta}(t) \leq \lambda tw(t), \end{aligned} \quad (\text{CVG-CDT})$$

$\alpha \geq 3$, $t_0 > 0$, $\varepsilon > 0$ are real numbers, β and γ are nonnegative continuously differentiable functions defined on $[t_0, +\infty)$. Then $x(t)$ is bounded and

$$\begin{aligned} f(x(t)) - f_* &\leq \mathcal{O}\left(\frac{1}{\delta(t)}\right), \quad \|\nabla f(x(t))\| \leq \mathcal{O}\left(\frac{1}{t\beta(t)}\right), \quad \|\dot{x}(t)\| \leq \mathcal{O}\left(\frac{1}{t}\right), \\ \int_{t_0}^{\infty} (\lambda tw(t) - \dot{\delta}(t))(f(x(t)) - f_*) dt &\leq \infty, \quad \int_{t_0}^{\infty} t(\alpha - 1 - \lambda)\|\dot{x}(t)\|^2 dt \leq \infty, \\ \int_{t_0}^{\infty} t^2\beta(t)w(t)\|\nabla f(x)\|^2 dt &\leq \infty, \quad \int_{t_0}^{\infty} t^2\beta(t)\langle \nabla^2 f(x(t))\dot{x}(t), \dot{x}(t) \rangle dt \leq \infty \end{aligned}$$

Proof: Lyapunov function and term cancelling

- ▶ Construct the Lyapunov function

$$\begin{aligned} E(t) = & \delta(t) (f(x(t)) - f_*) + \frac{1}{2} \|\lambda(x(t) - x_*) + t(\dot{x}(t) + \kappa\beta(t)\nabla f(x(t)))\|^2 \\ & + \lambda(1 - \kappa)t\beta(t)\langle \nabla f(x(t)), x(t) - x_* \rangle + \frac{\kappa(1 - \kappa)}{2} \|t\beta(t)\nabla f(x)\|^2 \quad (\text{Lya}) \\ & + \frac{\lambda(\alpha - 1 - \lambda)}{2} \|x(t) - x_*\|^2 \end{aligned}$$

- ▶ Differentiating through t , we set the term with brown color to 0:

$$\begin{aligned} \frac{d}{dt} E(t) = & \dot{\delta}(t)(f(x(t)) - f_*) - \lambda t w(t)\langle \nabla f(x(t)), x(t) - x_* \rangle - (\alpha - 1 - \lambda)t\|\dot{x}(t)\|^2 \\ & + \left(\delta(t) - (t^2 u(t) + (\kappa(\alpha - 1 - \lambda) - \lambda(1 - \kappa))t\beta(t)) \right) \langle \nabla f(x(t)), \dot{x}(t) \rangle \\ & - \kappa t^2 \beta(t) w(t) \|\nabla f(x(t))\|^2 - (1 - \kappa)t^2 \beta(t) \langle \nabla^2 f(x(t)) \dot{x}(t), \dot{x}(t) \rangle \leq 0 \end{aligned}$$

- ▶ Integrating the inequality above from t_0 to t completes the proof of Theorem 1

Applying forward Euler scheme to (ISHD)

- ▶ Let $v(t_0) = x(t_0) + \beta(t_0)\nabla f(x(t_0))$ and

$$\psi_{\Xi}(x(t), v(t), t) = \begin{pmatrix} v(t) - \beta(t)\nabla f(x(t)) \\ -\frac{\alpha}{t}(v(t) - \beta(t)\nabla f(x(t))) + (\dot{\beta}(t) - \gamma(t))\nabla f(x(t)) \end{pmatrix} \quad (1)$$

- ▶ The equation (ISHD) can be reformulated as the first-order system

$$\begin{pmatrix} \dot{x}(t) \\ \dot{v}(t) \end{pmatrix} = \psi_{\Xi}(x(t), v(t), t), \text{ notice that } \nabla^2 f(x(t))\dot{x}(t) = \frac{d}{dt}\nabla f(x(t))$$

- ▶ Let h be the step size, $t_k = t_0 + kh, k \geq 0$. The forward Euler scheme of the (ISHD) is

$$\begin{cases} \frac{x_{k+1} - x_k}{h} = v_k - \beta(t_k)\nabla f(x_k), \\ \frac{v_{k+1} - v_k}{h} = -\frac{\alpha}{t}(v_k - \beta(t_k)\nabla f(x_k)) + (\dot{\beta}(t_k) - \gamma(t_k))\nabla f(x_k) \end{cases} \quad (\text{F-Euler})$$

Explicit Discretization with Fixed Stepsize is Unstable

Set $p = 5, \alpha = 2p + 1, \beta(t) \equiv 0$ and $\gamma(t) = p^2 t^{p-2}$ in (ISHD). Consider

$$\min_{x \in \mathbb{R}^n} f(x) = \frac{1}{N} \sum_{i=1}^N \log(1 + \exp(-b_i \langle a_i, w \rangle))$$

where the data pairs $\{a, b_i\} \in \mathbb{R}^n \times \{0, 1\}, i \in [N]$.

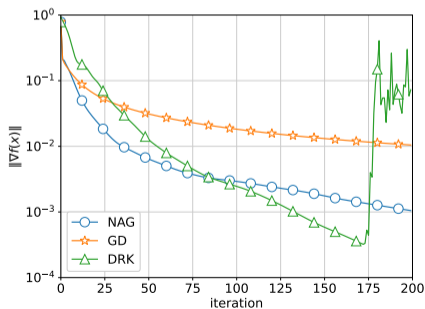


Figure: Directly applying 4-th Runge-Kutta diverges

Conditions for stable discretization

Theorem 2

Suppose the assumptions in Theorem 1 and (CVG-CDT) hold. Given t_0 , s_0 , and h , the sequence $\{x_k\}_{k=0}^{\infty}$ is generated by the equation (F-Euler) and $\bar{x}(t)$ is defined as

$$\bar{x}(t) = x_k + \frac{x_{k+1} - x_k}{h}(t - t_k), \quad t \in [t_k, t_{k+1}). \quad (2)$$

Assume three constants $0 \leq C_1$, $0 < C_2 \leq 1/h - 1/t_0$, and $0 < C_3$ fulfill

$$|\dot{\beta}(t)| \leq C_1\beta(t), \quad |\dot{\gamma}(t) - \ddot{\beta}(t)| \leq C_2(\gamma(t) - \dot{\beta}(t)), \quad \beta(t) \leq C_3w(t). \quad (3)$$

Then, it holds $f(x_k) - f_* \leq \mathcal{O}(1/k)$ under the following stability condition:

$$\Lambda(x, f) \geq \|\nabla^2 f(x)\|, \quad \alpha\beta(t)/t \leq \gamma(t) - \dot{\beta}(t) \leq \beta(t)/h, \quad (\text{STB-CDT})$$

$$\sqrt{\int_0^1 \Lambda((1-\tau)X(t, \Xi, f) + \tau\bar{x}(t), f) d\tau} \leq \frac{\sqrt{\gamma(t) - \dot{\beta}(t)} + \sqrt{\gamma(t) - \dot{\beta}(t) - \frac{\alpha}{t}\beta(t)}}{\beta(t)}.$$

Key technique: error decomposition

- ▶ Local Truncated Error:

$$\varphi(t) = \begin{pmatrix} x(t+h) - x(t) \\ v(t+h) - v(t) \end{pmatrix} - h \begin{pmatrix} v(t) - \beta(t)\nabla f(x(t)) \\ -\frac{\alpha}{t}v(t) + \left(\frac{\alpha}{t}\beta(t) + \dot{\beta}(t) - \gamma(t)\right)\nabla f(x(t)) \end{pmatrix}. \quad (4)$$

- ▶ Global error: $r_k = x(t_k) - x_k$, $s_k = v(t_k) - v_k$, and $e_k = (r_k, s_k)$
- ▶ We only need to control e_{k+1} , which has two resources

$$\begin{aligned} \begin{pmatrix} r_{k+1} \\ s_{k+1} \end{pmatrix} &= \begin{pmatrix} x(t_k) \\ v(t_k) \end{pmatrix} + \int_{t_k}^{t_k+h} \psi(s) ds - \begin{pmatrix} x_k \\ v_k \end{pmatrix} - h\psi(t_k) \\ &= \underbrace{\begin{pmatrix} I - h\beta(t_k)G(t_k) & hl \\ (\alpha\beta(t_k)/t_k + \dot{\beta}(t_k) - \gamma(t_k))G(t_k) & (1 - \alpha h/t_k)I \end{pmatrix}}_{W(t_k, G(t_k))} \begin{pmatrix} r_k \\ s_k \end{pmatrix} + h\varphi(t_k) \end{aligned}$$

where $G(t_k) = \int_0^1 \nabla^2 f(x(t_k) + \tau r_k) d\tau$. Abbreviate $W_k = W(t_k, G(t_k))$

Control $\rho(t)$ and e_k simultaneously

- Define the contraction factor $\rho(t) = \|W(t, G(t))\|$. We need only to estimate the **nonlinear** recurrence relation (**convolution series**)

$$\|e_{n+1}\| \leq \|W_n\| \|e_n\| + h \|\varphi(t_n)\| \leq \prod_{k=0}^n \|W_k\| \|e_0\| + \|\varphi(t_n)\| + \sum_{k=0}^{n-1} \prod_{l=k+1}^n \|W_l\| \|\varphi(t_l)\|$$

- We enhance Theorem 2 and prove the following propositions for each k :

$$\underbrace{\rho(t_k) \leq 1 - \frac{M}{2t_k}}_{P_0(k)}, \quad \underbrace{\|e_k\| \leq \frac{M_1}{\sqrt{t_k}}}_{P_1(k)}, \quad \text{and} \quad \underbrace{f(x_k) - f_* \leq \mathcal{O}\left(\frac{1}{k}\right)}_{P_2(k)} \quad (5)$$

- $\{P_0(k)\}_{k \leq n}$ implies $P_1(n)$, while $P_0(n)$ relies on stable condition and $P_1(n)$
- $P_1(n)$ implies $P_3(n)$: The function value minimization rate is

$$\begin{aligned} f(x_k) - f_* &\leq |f(x_k) - f(x(t_k))| + |f(x(t_k)) - f_*| \\ &\leq \underbrace{\|\nabla f(x(t_k))\|}_{\mathcal{O}(1/(t\beta(t)))} \|e_k\| + \frac{1}{2} \left\| \int_0^1 \nabla^2 f(x(t_k) + \tau e_k) d\tau \right\| \underbrace{\|e_k\|^2}_{\mathcal{O}(1/t)} + \frac{E(t_0)}{t^2 w(t)} \end{aligned}$$

Proof of the convergence property of the forward Euler discretization

- ▶ Using Cauchy inequality and Theorem 1, for certain M_3 , we have

$$\|\varphi(t)\| \leq o(1/t) \quad \text{and} \quad \sum_{k=0}^n t_k^{\alpha/2} \|\varphi(t_k)\| \leq M_3 t_n^{\alpha/2-1/2}. \quad (6)$$

- ▶ Assume $P_0(k)$ and $P_1(k)$ are valid for $k \leq n$. For $k \leq n$, we have

$$\begin{aligned} \prod_{l=k}^n \|W_l\| &= \prod_{l=k}^n \rho(t_l) = \exp\left(\sum_{l=k}^n \ln(\rho_l - 1 + 1)\right) \leq \exp\left(\sum_{l=k}^n (\rho_l - 1)\right) \\ &\leq \exp\left(-\sum_{l=k}^n \frac{\alpha h}{2t_l}\right) \leq \exp\left(-\frac{\alpha}{2} \int_{t_k}^{t_{n+1}} \frac{1}{t} dt\right) = \left(\frac{t_k}{t_{n+1}}\right)^{\alpha/2} \end{aligned} \quad (7)$$

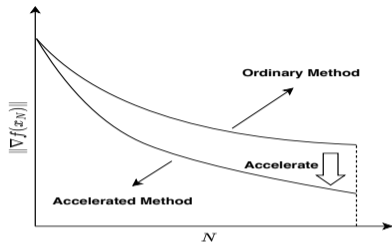
- ▶ $\{P_0(k)\}_{k \leq n}$ implies $P_1(n)$: Combine $\{P_0(k)\}_{k \leq n}$ and equation (6), we have

$$\|e_{n+1}\| \leq \|\varphi(t_n)\| + \sum_{k=0}^{n-1} \prod_{l=k+1}^n \|W_l\| \|\varphi(t_l)\| \leq \sum_{k=0}^n \left(\frac{t_{k+1}}{t_{n+1}}\right)^{\alpha/2} \|\varphi(t_k)\| \leq M_3 \frac{1}{\sqrt{t_{n+1}}}$$

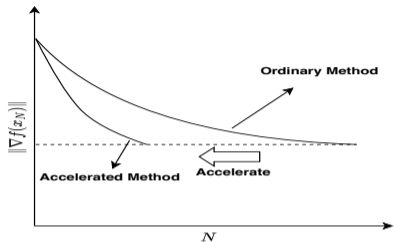
Outline

- 1 Conditions for stability-preserving discretization
- 2 Select the best coefficients using learning to optimize**
- 3 Computation of the conservative gradients
- 4 Convergence analysis
- 5 Numerical results

Stopping time: a differentiable continuous-time complexity



(a) Measure-based



(b) Complexity-based

Definition 3 (Stopping Time)

Given the initial time t_0 , the initial value x_0 , the initial velocity $\dot{x}(t_0)$, the trajectory $X(\Xi, t, f)$ of the system (ISHD), and a tolerance ε , the stopping time of the criterion $\|\nabla f(x)\| \leq \varepsilon$ is

$$T(\Xi, f) = \inf\{t \mid \|\nabla f(X(\Xi, t, f))\| \leq \varepsilon, t \geq t_0\}. \quad (8)$$

Tackle the point-wise constraints using integration

- ▶ With $w(t), \delta(t)$ defined in (CVG-CDT), we introduce

$$p(x, \bar{x}, \Xi, t, f) = \left[\beta(t) \sqrt{\int_0^1 \Lambda((1-\tau)x + \tau\bar{x}, f) d\tau} - \sqrt{\gamma(t) - \dot{\beta}(t)} \right. \\ \left. - \sqrt{\gamma(t) - \dot{\beta}(t) - \frac{\alpha}{t}\beta(t)} \right]_+, \\ q(\Xi, t) = \left[\gamma(t) - \dot{\beta}(t) - \beta(t)/h \right]_+ + \left[\dot{\beta}(t) + \alpha\beta(t)/t - \gamma(t) \right]_+ \\ + \left[\dot{\delta}(t) - \lambda tw(t) \right]_+ + [-\delta(t)]_+.$$

- ▶ Setting $P, Q \leq 0$ ensures (CVG-CDT) and (STB-CDT) hold for f

$$P(\Xi, f) = \int_{t_0}^{T(\Xi, f)} p(X(t, \Xi, f), \bar{x}(t), \Xi, t, f) dt, \quad Q(\Xi, f) = \int_{t_0}^{T(\Xi, f)} q(\Xi, t) dt$$

A L2O framework for selecting the best coefficients

- ▶ Induced distribution: Given a random variable $\xi \sim \mathbb{P}$. We say \mathbb{P} is the induced probability of the parameterized function $f(\cdot; \xi)$

$$\mathbb{E}_f[T(\Xi, f)] = \int_{\xi} T(\Xi, f(\cdot; \xi)) d\mathbb{P}(\xi) = \mathbb{E}_{\xi}[T(\Xi, f(\cdot; \xi))]$$

- ▶ Framework: minimize the expectation of stopping time under conditions of convergence and stable discretization

$$\min_{\Xi} \mathbb{E}_f[T(\Xi, f)], \tag{9}$$

$$\text{s.t. } \mathbb{E}_f[P(\Xi, f)] \leq 0, \quad \mathbb{E}_f[Q(\Xi, f)] \leq 0, \tag{10}$$

- ▶ **Parameterization:** $\beta \rightarrow \beta_{\theta_1}, \gamma \rightarrow \gamma_{\theta_2}$. Set $\theta = (\alpha, \theta_1, \theta_2)$.

Solve the L2O problem using exact penalty method

Given the penalty parameter ρ , the ℓ_1 exact penalty problem writes

$$\begin{aligned}\min_{\theta} \Upsilon(\theta) &= \mathbb{E}_f[T(\theta, f)] + \rho(\mathbb{E}_f[P(\theta, f)] + \mathbb{E}_f[Q(\theta, f)]) \\ &= \mathbb{E}_f[T(\theta, f) + \rho(P(\theta, f) + Q(\theta, f))]\end{aligned}\tag{11}$$

Algorithm Stochastic Penalty Method (StoPM) for L2O problem

- 1: **Input:** initial weight θ_0 , penalty coefficient ρ , training dataset \mathcal{F}
 - 2: **while** Not(Stopping Condition) **do**
 - 3: Sample a function: $f_k \in \mathcal{F}$.
 - 4: Computing the gradients J_T, J_P and J_Q correspond to T, P and Q
 - 5: Update variable: $\theta_{k+1} \leftarrow \theta_k - \eta(J_T + J_P + J_Q)$.
 - 6: Update index: $k \leftarrow k + 1$.
 - 7: **end while**
 - 8: **Output:** the trained weight θ_* .
-

Outline

- 1 Conditions for stability-preserving discretization
- 2 Select the best coefficients using learning to optimize
- 3 Computation of the conservative gradients**
- 4 Convergence analysis
- 5 Numerical results

Conservative gradient

- ▶ When parameterize α, β, γ using neural networks, they may be **nonsmooth**
- ▶ The output of *auto differentiation* in **nonsmooth** functions may not be **Clarke subdifferentials**, but they are certainly **conservative gradients**
- ▶ **Conservative gradient** generalizes subdifferentials while preserving **chain rule**
- ▶ Ψ is termed the **conservative Jacobian (gradient if $m = 1$)** of π if and only if

$$\frac{d}{dt}\pi(r(t)) = Ar(t), \quad \text{for all } A \in \Psi(r(t)), \text{ for almost all } t \in [0, 1],$$

for any absolutely continuous curve $r : [0, 1] \rightarrow \mathbb{R}^d$

- ▶ Consider the example:

$$f(s) = ([-s]_+ + s) - [s]_+ \equiv 0 \quad \xrightarrow{\text{autograd using TensorFlow}} \quad g(s) = \begin{cases} 0 & (s \neq 0) \\ 1 & (s = 0) \end{cases}$$

g is not the Clarke subdifferential of f but a conservative gradient

Differentiate through the ODE flow of (ISHD): $\partial X / \partial \theta$

- ▶ Reformulate (ISHD) as a first-order system (1) with a parameterized right-hand-side term ψ :

$$\psi: \mathbb{R}^{2n+1+p} \rightarrow \mathbb{R}^{2n}, \quad (x, v, t, \theta) \mapsto \psi_\theta(x, v, t). \quad (12)$$

Denote the flow of (1) with parameterized ψ as $X(x_0, v_0, \theta, t)$

- ▶ Denote $J^\psi: \mathbb{R}^{2n+1+p} \rightarrow \mathbb{R}^{2n \times (2n+1+p)}$ as a conservative Jacobian of ψ with respect to (x, v, t, θ) . The coordinate projection (partial derivative) writes

$$J_{x,v}^\psi = \Pi_{x,v} J^\psi, \quad J_t^\psi = \Pi_t J^\psi \quad \text{and} \quad J_\theta^\psi = \Pi_\theta J^\psi$$

- ▶ Applying the general result to the first-order system (1): $\theta \mapsto A(t_0)$ is a conservative Jacobian of $\theta \rightarrow X(x_0, v_0, \theta, t_1)$ (smooth version: $\partial X / \partial \theta$)

$$\dot{A}(t) = J_{x,v}^\psi(t)A(t) + J_\theta(t), \quad A(t_1) = 0_{2n \times p} \quad \text{for all } t \in [t_0, t_1] \quad (13)$$

- ▶ Smooth version:

$$\frac{\partial X}{\partial \theta} = \int_{t_0}^{t_1} \frac{\partial \psi_\theta}{\partial x} \frac{dX}{d\theta} + \frac{\partial \psi_\theta}{\partial \theta} dt$$

Evaluate the derivative of stopping time: $\nabla_{\theta} T(f, \theta)$

- ▶ Take limit by continuity: $\|\nabla f(X(T(f, \theta), f, \theta))\|^2 - \varepsilon^2 \equiv 0$
- ▶ Implicit function theorem (valid in nonsmooth case):

$$\nabla f(X)^{\top} \nabla^2 f(X) \left(\frac{\partial X}{\partial t} \Big|_{t=T} \nabla_{\theta} T(f, \theta) + \frac{\partial X}{\partial \theta} \right) = 0$$

where $T = T(f, \theta)$, $X = X(T(f, \theta), f, \theta)$

- ▶ Invoking the first-order form of (ISHD):

$$\frac{\partial X}{\partial t} \Big|_{t=T} = \dot{x}(T) = v(T) - x(T) - \beta(T) \nabla f(x(T))$$

where $x(t) = X(t, f, \theta)$

- ▶ The derivative:

$$\nabla_{\theta} T(f, \theta) = \left(\nabla f(X)^{\top} \nabla^2 f(X) (v(T) - X - \beta(T) \nabla f(X)) \right)^{-1} \nabla f(X)^{\top} \nabla^2 f(X) \frac{\partial X}{\partial \theta}$$

Conservative gradient of the constraints

- ▶ Recap:

$$P(\theta, f) = \int_{t_0}^{T(\theta, f)} p(X(t, \theta, f), \bar{x}(t), \theta, t, f) dt, \quad Q(\theta, f) = \int_{t_0}^{T(\theta, f)} q(\theta, t) dt$$

- ▶ Applying the chain rule gives

$$\frac{dP(\theta, f)}{d\theta} = \int_{t_0}^{T(\theta, f)} \frac{\partial \psi_{\theta}(s(t), t, f)}{\partial \theta} w(t) + \frac{\partial p(x(t), \bar{x}(t), \theta, t, f)}{\partial \theta} dt + p(x(T), \bar{x}(T), \theta, T, f) \frac{dT(\theta, f)}{d\theta},$$

$$\frac{dQ(\theta, f)}{d\theta} = \int_{t_0}^{T(\theta, f)} \frac{dq(\theta, t)}{d\theta} dt + q(\theta, T) \frac{dT(\theta, f)}{d\theta},$$

where $\bar{x}(\cdot)$ is the interpolation defined in (2), $w(\cdot)$ is the solution of

$$-\begin{pmatrix} \frac{\partial p(x(t), \bar{x}(t), \theta, t, f)}{\partial x} \\ 0_{n \times 1} \end{pmatrix} - \frac{\partial \psi_{\theta}(s(t), t, f)}{\partial s} w(t) = \dot{w}(t), \quad w(T(\theta, f)) = 0_{2n \times 1}.$$

Clarke subdifferential of the point-wise maximal function

- Let $\{f_\eta: \mathbb{R}^n \rightarrow (-\infty, +\infty]\}_{\eta \in A}$ be a family of *proper convex* functions and

$$f(x) = \sup_{\eta \in A} f_\eta(x)$$

- If $x_0 \in \bigcap_{\eta \in A} \text{int dom } f_\eta$, and $I(x_0) = \{\eta \in A \mid f_\eta(x_0) = f(x_0)\}$, then

$$\text{conv} \left(\bigcup_{\eta \in I(x_0)} \partial f_\eta(x_0) \right) = \partial f(x_0)$$

- $\lambda_{\max}(A) = \sup_{\|u\|=1} u^\top A u$. Set $v = \arg \max_{\|u\|=1} u^\top \nabla^2 f(x) u$. We have

$$\begin{aligned} \frac{\partial \Lambda(f, x)}{\partial x} &= \left\{ \left\langle \frac{\partial \lambda_{\max}(A)}{\partial A} \Big|_{A=\nabla^2 f(x)}, \frac{\partial \nabla^2 f(x)}{\partial x_k} \right\rangle \right\}_k = \left\{ \sum_{i,j} \partial_{ijk} f(x) v_i v_j \right\}_k \\ &= \frac{d}{d\eta_2} \left(\frac{d}{d\eta_1} \nabla f(x + \eta_1 v + \eta_2 v) \Big|_{\eta_1=0} \right) \Big|_{\eta_2=0} = D^3 f(x)[v, v], \end{aligned}$$

when $\Lambda(f, x) = \lambda_{\max}(\nabla^2 f(x))$. This enables the evaluation of $\partial p / \partial \theta, \partial p / \partial x$

Outline

- 1 Conditions for stability-preserving discretization
- 2 Select the best coefficients using learning to optimize
- 3 Computation of the conservative gradients
- 4 Convergence analysis**
- 5 Numerical results

Criteria for Clarke stationarity using directional derivative

- ▶ Let \varkappa be Lipschitz continuous near $\bar{\theta}$
- ▶ The Clarke directional derivative of \varkappa at $\bar{\theta}$ along a nonzero vector ϑ :

$$\varkappa^\circ(\bar{\theta}; \vartheta) \triangleq \limsup_{\substack{\theta \rightarrow \bar{\theta} \\ \tau \downarrow 0}} \frac{\varkappa(\theta + \tau\vartheta) - \varkappa(\theta)}{\tau} \quad (14)$$

- ▶ The Clarke subdifferential of \varkappa at θ is given by

$$\partial\varkappa(\theta) \triangleq \left\{ \mathbf{a} \in \mathbb{R}^{d_\theta} : \varkappa^\circ(\theta; \vartheta) \geq \mathbf{a}^\top \vartheta, \quad \forall \vartheta \in \mathbb{R}^{d_\theta} \right\}$$

- ▶ Clarke stationarity: $0 \in \partial\varkappa(\theta)$
- ▶ θ is a Clarke stationary point of \varkappa if and only if $\varkappa^\circ(\theta; \vartheta) \geq 0$ for all $\vartheta \in \mathbb{R}^{d_\theta}$

Precludes infeasible stationary point using sufficient decrease condition

- ▶ Given the training dataset \mathcal{F} , we denote the residual function as

$$R(\theta) = \mathbb{E}_f[P(\theta, f) + Q(\theta, f)]$$

- ▶ This function measures the constraints violation. The feasible set is defined by

$$S = \{\theta \mid P(\theta, f) \leq 0, Q(\theta, f) \leq 0, \forall f \in \mathcal{F}\} \quad (15)$$

Assumption 1 (Uniform sufficient decrease condition)

For each infeasible point θ , i.e. $\theta \notin S$, there exists a nonzero vector ϑ , such that $R^\circ(\theta; \vartheta) \leq -c\|\vartheta\|$. Here the constant c is uniform for each θ .

Theorem 4

Suppose $\mathbb{E}_f[T(\theta, f)]$ is globally Lipschitz continuous with Lipschitz constant L_T . Let Assumption 1 hold. Given the penalty parameter $\rho > L_T/c$, any infeasible point of the penalty function Υ can not be a D -stationary point.

Sufficient decrease condition precludes infeasible stationary point

- ▶ Consider the penalty function

$$\Upsilon(\theta) = \mathbb{E}_f [T(\theta, f) + \rho(P(\theta, f) + Q(\theta, f))]$$

- ▶ For any infeasible point θ , using Assumption 1, there must exist a direction ϑ , such that

$$\Upsilon^\circ(\theta; \vartheta) = \mathbb{E}_f [T(\cdot, f)]^\circ(\theta; \vartheta) + \rho R^\circ(\theta, \vartheta) \leq L_T \|\vartheta\| - c\rho \|\vartheta\| < 0.$$

- ▶ Invoking the criteria of Clarke stationary point, we know θ cannot be a Clarke stationary point of Υ
- ▶ Clarke subgradient is a minimal conservative gradient
- ▶ For any conservative gradient J^Υ of Υ , we have $\partial\Upsilon \subset J^\Upsilon$. Hence, θ can not be a D -stationary point of Υ

SGD converges with (nonsmooth) auto-differentiation: Assumptions

Assumption 2 (Assumptions of the SGD)

1. The step sizes $\{\eta_k\}_{k \geq 1}$ satisfy

$$\eta_k \geq 0, \quad \sum_{k=1}^{\infty} \eta_k = \infty, \quad \text{and} \quad \sum_{k=1}^{\infty} \eta_k^2 < \infty.$$

2. Almost surely, the iterates $\{\theta_k\}_{k \geq 1}$ are bounded, i.e., $\sup_{k \geq 1} \|\theta_k\| < \infty$.
3. $\{\xi_k\}_{k \geq 1}$ is a uniformly bounded difference martingale sequence with respect to the increasing σ -fields

$$\mathcal{F}_k = \sigma(\theta_j, \varrho_j, \xi_j : j \leq k).$$

In other words, there exists a constant $M_\xi > 0$ such that

$$\mathbb{E}[\xi_k \mid \mathcal{F}_k] = 0 \quad \text{and} \quad \mathbb{E}[\|\xi_k\|^2 \mid \mathcal{F}_k] \leq M_\xi \quad \text{for all } k \geq 1.$$

SGD converges with (nonsmooth) auto-differentiation

Assumption 3

The complementary of $\{\Upsilon(\theta) \mid 0 \in J^\Upsilon(\theta)\}$ is dense in \mathbb{R} .

Theorem 5 (SGD converges using conservative gradient)

Suppose that Assumptions 2 and 3 hold. Then every limit point of $\{\theta_k\}_{k \geq 1}$ is stationary and the function values $\{z(\theta_k)\}_{k \geq 1}$ converge.

Theorem 6 (Convergence guarantee for Algorithm 1)

Suppose Assumption 1, 2 and 3 hold, $\{\theta_k\}_{k \geq 1}$ is generated by Algorithm 1. Then almost surely, every limit point θ_ of $\{\theta_k\}_{k \geq 1}$ satisfies $\theta_* \in S_f$, $0 \in J^\Upsilon(\theta_*)$ and the sequence $\{\Upsilon(\theta_k)\}_{k \geq 1}$ converges.*

Outline

- 1 Conditions for stability-preserving discretization
- 2 Select the best coefficients using learning to optimize
- 3 Computation of the conservative gradients
- 4 Convergence analysis
- 5 Numerical results**

Setting and datasets

- ▶ Consider the logistic regression problem defined by

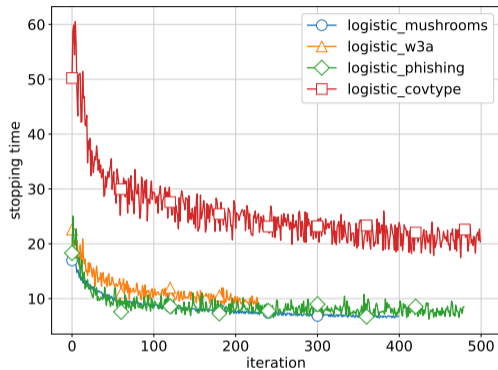
$$\min_{x \in \mathbb{R}^n} f_{\mathcal{D}}(x) = \frac{1}{|\mathcal{D}|} \sum_{(a_i, b_i) \in \mathcal{D}} \log(1 + \exp(-b_i \langle a_i, x \rangle)),$$

where \mathcal{D} is a subset of a given dataset Σ and $\{a_i, b_i\} \in \mathbb{R}^n \times \{0, 1\}$, $i \in [|\mathcal{D}|]$

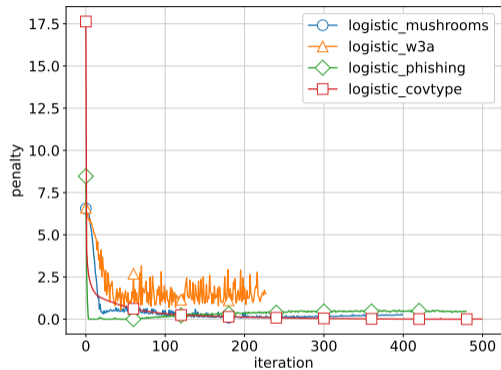
- ▶ The datasets are listed as below

Dataset	n	N_{train}	N_{test}	Separable
a5a	123	6,414	26,147	No
w3a	300	4,912	44,837	No
mushrooms	112	3,200	4,924	Yes
covtype	54	102,400	478,612	No
phishing	68	8,192	2,863	No
separable	101	20,480	20,480	Yes

Training results



(a) Stopping time on logistic regression



(b) Penalty on logistic regression

Figure: The training process in different tasks.

Testing: Compared methods

- ▶ **GD.** $x_{k+1} = x_k - h\nabla f(x_k)$. We set the stepsize as $h = 1/L$
- ▶ **NAG.** We choose $h = 1/L$ and employ the version for convex functions

$$y_{k+1} = x_k - h\nabla f(x_k), \quad x_{k+1} = y_{k+1} + \frac{k-1}{k+2}(y_{k+1} - y_k)$$

- ▶ **EIGAC.** Explicit inertial gradient algorithm with correction, i.e. Algorithm corresponds to (F-Euler). We provide two versions of EIGAC with default coefficients $\alpha = 6$, $\beta(t) = (4/h - 2\alpha/t)/L$, and $\beta(t) = h\gamma(t)$ and the coefficients learned by Algorithm 1. The numerical experiments effectively show that the EIGAC with default coefficients are sufficient to converge and the performance is comparable with NAG, while EIGAC with learned coefficients is superior over other methods.

Testing: Compared methods

- ▶ **IGAHD**. Inertial gradient algorithm with Hessian-driven damping. This method is obtained by applying a NAG inspired time discretization of

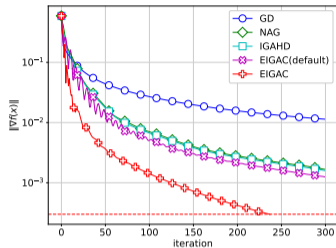
$$\ddot{x}(t) + \frac{\alpha}{t}\dot{x}(t) + \beta\nabla^2 f(x(t))\dot{x}(t) + \left(1 + \frac{\beta}{t}\right)\nabla f(x(t)) = 0. \quad (16)$$

Let $s = 1/L$. In each iteration, setting $\alpha_k = 1 - \alpha/k$, the method performs

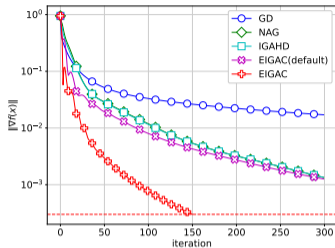
$$\begin{cases} y_k = x_k + \alpha_k (x_k - x_{k-1}) - \beta\sqrt{s}(\nabla f(x_k) - \nabla f(x_{k-1})) - \frac{\beta\sqrt{s}}{k}\nabla f(x_{k-1}), \\ x_{k+1} = y_k - s\nabla f(y_k). \end{cases} \quad (17)$$

It has been show that IGAHD owns $\mathcal{O}(1/k^2)$ convergence rate when $0 \leq \beta < 2/\sqrt{s}$ and $s \leq 1/L$. Its performance may not coincide with NAG due to the existence of the gradient correction term. In our experiments, IGAHD serves as a baseline of the optimization methods derived from the ODE viewpoint without learning.

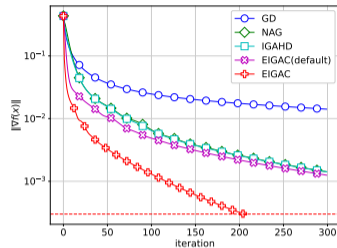
Testing results



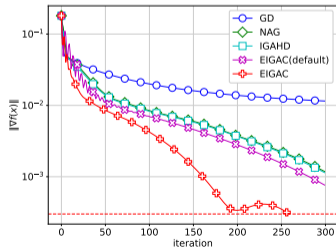
(a) a5a



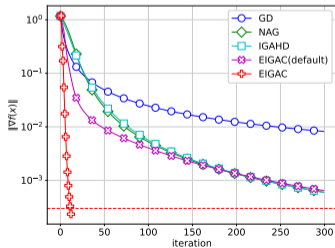
(b) mushrooms



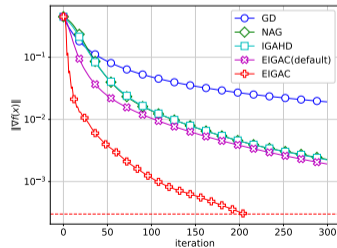
(c) w3a



(d) covtype



(e) separable



(f) phishing