

A geometric integration approach to non-smooth and non-convex optimisation

Martin Benning¹, Matthias Ehrhardt², GRW Quispel³,
Erlend Skaldehaug Riis⁴, Torbjørn Ringholm⁵, Carola-Bibiane Schönlieb⁴

¹Queen Mary University of London, UK

²University of Bath, UK

³La Trobe University, Australia

⁴University of Cambridge, UK

⁵Norwegian University of Science and Technology, Norway

Variational Methods and Optimization in Imaging
IHP, Paris



UNIVERSITY OF
CAMBRIDGE

- 1 Geometric numerical integration and the discrete gradient method
- 2 The DG method for nonsmooth, nonconvex optimisation
- 3 Beyond gradient flow

Geometric numerical integration and the discrete gradient method

- $\min_{x \in \mathbb{R}^n} V(x), \quad \dot{x}(t) = -\nabla V(x(t))$ (gradient flow)
- Forward Euler $\rightarrow x^{k+1} = x^k - \tau \nabla V(x^k),$
Backward Euler $\rightarrow x^{k+1} = x^k - \tau \nabla V(x^{k+1}).$
- Numerical integration and analysis of ODEs
 - Optimisation scheme \rightarrow ODE
Accelerated gradient descent $\rightarrow \ddot{x} + \frac{3}{t}\dot{x} = -\nabla V(x)$
[Su, Boyd, Candès (2016)]
 - Numerical integration tools to improve efficiency
Runge-Kutta methods for stiff ODEs \rightarrow Larger time steps
[Eftekhar, Vandereycken, Vilmart, Zygalakis (2018)]
 - Discretisation methods for structure preservation of ODE
Symmetry preservation \rightarrow acceleration phenomenon
[Betancourt, Jordan, Wilson (2018)]

Geometric integration and discrete gradients

- Geometric numerical integration
 - ODEs have structure (conservation laws, symplectic structure.)
 - Aim: **Preserve structure** when numerically solving ODEs
- Discrete gradient (DG) method
 - Preserves first integrals; energy conservation and **dissipation** laws, Lyapunov functions¹
- Optimisation methods
 - Want to solve $\min_{x \in \mathbb{R}^n} V(x)$.
 - Apply DG method to **gradient flow** to preserve dissipative structure

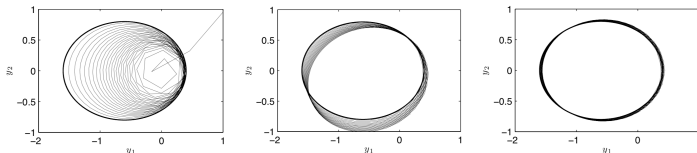


Figure: Dahlby, Owren, Yaguchi (2011).

'Why geometric numerical integration?' [Iserles, Quispel (2018)]

¹McLachlan, Quispel, and Robidoux (1999).

Discrete gradient method

Definition

For a smooth function $V : \mathbb{R}^n \rightarrow \mathbb{R}$, a **discrete gradient** $\bar{\nabla} V(x, y)$ satisfies

- (i) $\lim_{y \rightarrow x} \bar{\nabla} V(x, y) = \nabla V(x)$ (consistency),
- (ii) $\langle \bar{\nabla} V(x, y), y - x \rangle = V(y) - V(x)$ (mean value).

$$\min V : \mathbb{R}^n \rightarrow \mathbb{R}$$

Discrete gradient method.

$$x^{k+1} = x^k - \tau \bar{\nabla} V(x^k, x^{k+1})$$

Dissipative:

$$\begin{aligned} \frac{V(x^{k+1}) - V(x^k)}{\tau} &= \langle x^{k+1} - x^k, \bar{\nabla} V(x^k, x^{k+1}) \rangle \\ &= -\|\bar{\nabla} V(x^k, x^{k+1})\|^2 \\ &= -\left\| \frac{x^{k+1} - x^k}{\tau} \right\|^2. \end{aligned}$$

Gradient flow.

$$\dot{x}(t) = -\nabla V(x(t)).$$

$$\begin{aligned} \frac{d}{dt} V(x(t)) &= \langle \dot{x}(t), \nabla V(x(t)) \rangle \\ &= -\|\nabla V(x(t))\|^2 \\ &= -\|\dot{x}(t)\|^2. \end{aligned}$$

Convergence theorem for DG method

Theorem¹

Suppose V is C^1 -smooth, coercive, and bounded below, $0 < c < \tau < C$, and $x^{k+1} = x^k - \tau \bar{\nabla} V(x^k, x^{k+1})$. Then

- $\nabla V(x^k) \rightarrow 0$,
- $\|x^{k+1} - x^k\| \rightarrow 0$,
- (x^k) has an accumulation point x^* , and it satisfies $\nabla V(x^*) = 0$.



Figure: Inpainting with Itoh-Abe DG method¹

¹Grimm, McLachlan, McLaren, Quispel and Schönlieb (2017).

Examples of discrete gradients

- Gonzalez (midpoint) DG¹:

$$\bar{\nabla} V(x, y) = \nabla V\left(\frac{x+y}{2}\right) + \frac{V(y) - V(x) - \langle \nabla V\left(\frac{x+y}{2}\right), y-x \rangle}{\|x-y\|^2} (y-x).$$

- Mean value DG²:

$$\bar{\nabla} V(x, y) = \int_0^1 \nabla V((1-s)x + sy) ds.$$

- Itoh-Abe (coordinate increment) DG³:

$$\bar{\nabla} V(x, y) = \begin{pmatrix} \frac{V(y_1, x_2, \dots, x_n) - V(x_1, x_2, \dots, x_n)}{y_1 - x_1} \\ \frac{V(y_1, y_2, x_3, \dots, x_n) - V(y_1, x_2, x_3, \dots, x_n)}{y_2 - x_2} \\ \vdots \\ \frac{V(y_1, \dots, y_n) - V(y_1, y_2, \dots, y_{n-1}, x_n)}{y_n - x_n} \end{pmatrix}.$$

¹Gonzalez (1996). ²Celledoni, Grimm, et al. (2012). ³Itoh and Abe (1988).

Itoh-Abe discrete gradient (IADG) method

Applications

- Image inpainting with Euler's elastica (nonconvex).¹
- Successive-over-relaxation (SOR) and the Gauss-Seidel method, for linear systems $Ax = b$.²
 - Kaczmarz methods (by extension)



Figure: Inpainting with Itoh-Abe DG method¹.

¹ Ringholm, Lasić and Schönlieb (2017). ² Miyatake, Sogabe, Zhang (2017).

Rewrite IADG method as

$$x^{k+1} = x^k - \alpha d^k, \quad \text{s.t.} \quad \alpha = \tau_k \frac{V(x^k - \alpha d^k) - V(x^k)}{\alpha}, \quad (1)$$
$$d^k \in S^{n-1} := \{d \in \mathbb{R}^n : \|d\| = 1\}.$$

- Derivative-free gradient flow dissipative structure

$$V(x^{k+1}) - V(x^k) = -\frac{1}{\tau^k} \|x^{k+1} - x^k\|^2 = -\tau^k \left(\frac{V(x^{k+1}) - V(x^k)}{\|x^{k+1} - x^k\|} \right)^2$$

- Well-defined for nonsmooth functions; computationally tractable
- Descends along directions $(d^k)_{k \in \mathbb{N}}$
 - Standard IADG: Let (d^k) cycle through coordinates e^i
 - Can also randomise: Draw d^k randomly from $(e^i)_{i=1}^n$ or from S^{n-1}

- When the function is nonsmooth, nonconvex, and black-box
 - Bilevel optimisation of variational regularisation problems
 - Parameter optimisation of model simulations
 - 'Optimal camera placement to measure distances regarding static and dynamic obstacles'* [Hänel et al. (2012)]
- When gradients are computationally expensive
- When problem is poorly conditioned or stiff.

The DG method for nonsmooth, nonconvex optimisation

- The **Clarke subdifferential**¹ $\partial V(x)$ introduced by F. Clarke (1973).
- $\partial V(x) = \text{co} \left\{ p \mid \text{s.t. } x^k \rightarrow x, \nabla V(x^k) \rightarrow p \right\}$ (Convex hull of limiting gradients).
- Generalises gradient, and classical subdifferential for convex functions.
- Nice analytical properties for **locally Lipschitz continuous** functions. (**outer semicontinuity**; **mean value theorem**; **convex, compact**, non-empty sets)
- x is **Clarke stationary** when $0 \in \partial V(x)$.
- **Clarke directional derivative**:

$$V^o(x; d) := \limsup_{y \rightarrow x, \lambda \rightarrow 0} \frac{V(y + \lambda d) - V(y)}{\lambda}$$

- $0 \in \partial V(x) \iff V^o(x; d) \geq 0$ for all $d \in S^{n-1}$

¹Clarke (1990).

Want to prove accumulation points of $(x^k)_{k \in \mathbb{N}}$ are **Clarke stationary**.

- Local Lipschitz continuity of $V \Rightarrow$ **upper semicontinuity of $V^o(\cdot; \cdot)$** .
(closely related to outer semicontinuity of ∂V)
- $\|x^{k+1} - x^k\| \rightarrow 0, \quad \frac{V(x^{k+1}) - V(x^k)}{\|x^{k+1} - x^k\|} \rightarrow 0.$
- $(x^{k_j}, d^{k_j}) \rightarrow (x^*, d^*)$ and $\liminf_{j \rightarrow \infty} V^o(x^{k_j}; d^{k_j}) \geq 0 \Rightarrow \mathbf{V^o(x^*; d^*) \geq 0}.$
- Need to ensure \exists **dense subsequence d^{k_j}** corresponding to $x^{k_j} \rightarrow x^*$.
 - For **random** $(d^k)_{k \in \mathbb{N}}$, \iff full support of distribution of d^k on S^{n-1} .
 - For **deterministic** (d^k) , \iff “cyclical” density.

Theorem (Ehrhardt, Riis, Quispel, Schönlieb (2018))

Let $V : \mathbb{R}^n \rightarrow \mathbb{R}$ be a *locally Lipschitz continuous*, coercive function that is bounded below. Suppose (x^k) are the iterates from the generalised Itoh-Abe DG method with appropriate sequence of directions $(d_k)_{k \in \infty}$. Then

- The iterates converge to a *nonempty, connected*, compact set of *accumulation points*.
- All accumulation points are *Clarke stationary*.
- All accumulation points take the *same value* on V .

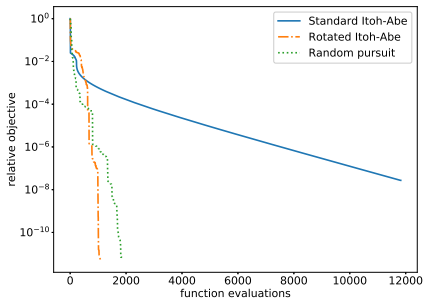
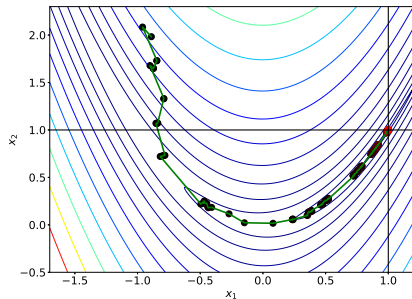
Other properties of DG method

When V is C^1 -smooth, the DG methods inherit properties from gradient descent/flow:

- **Convergence rates:**¹ $V(x^k) - V^* \rightarrow 0$.
 - $\mathcal{O}(1/k)$ if V is convex.
 - **Linearly** if V is Polyak–Łojasiewicz function/ strongly convex.
 - Itoh–Abe method has marginally better convergence rate than coordinate descent.
- **Kurdyka–Łojasiewicz inequality** $\implies (x^k)_{k \in \mathbb{N}}$ converges.
- Properties hold for **all time steps** $c < (\tau^k)_{k \in \mathbb{N}} < C$, $c, C > 0$.

¹Ehrhardt, Riis, Ringholm, Schönlieb (2018).

Rosenbrock function



Beyond gradient flow

Bregman distance and inverse scale space flow

- Let $J : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex, e.g.

$$J(x) = \gamma \|x\|_1 \quad \text{or} \quad \gamma \|x\|_1 + \|x\|^2/2 \quad \text{or} \quad \gamma \text{TV}(x).$$

- Define **Bregman distance** (notion of distance induced by J)

$$D_J^p(x, y) = J(y) - J(x) - \langle p, y - x \rangle \geq 0, \quad p \in \partial J(x).$$

- For inverse problem $b = Ax$, want to solve

$$\min_x J(x) \text{ s.t. } b = Ax \quad \text{or} \quad \min_x V(x) + \lambda J(x),$$

where $V(x) = \|Ax - b\|^2/2$.

- Consider **inverse scale space** (ISS) flow¹

$$\partial_t p(t) = -\partial V(x(t)), \quad p(t) \in \partial J(x(t)).$$

¹Burger, Gilboa, Osher, Xu (2006).

$$\partial_t p(t) = -\partial V(x(t)), \quad p(t) \in \partial J(x(t)).$$

- Backward Euler \rightarrow **Bregman iterations**¹:

$$\begin{aligned} p^{k+1} &= p^k - \tau^k \nabla V(x^{k+1}), \quad p^{k+1} \in \partial J(x^{k+1}) \\ \iff x^{k+1} &= \arg \min_x \tau^k V(x) + D_J^{p^k}(x^k, x). \end{aligned}$$

- Forward Euler \rightarrow **Linearised Bregman iterations**¹:

$$\begin{aligned} p^{k+1} &= p^k - \tau^k \nabla V(x^k), \quad p^{k+1} \in \partial J(x^{k+1}) \\ \iff x^{k+1} &= \arg \min_x \tau^k \left(V(x^k) + \langle \nabla V(x^k), x - x^k \rangle \right) + D_J^{p^k}(x^k, x). \end{aligned}$$

- DG method \rightarrow **Bregman DG method**:

$$p^{k+1} = p^k - \tau^k \bar{\nabla} V(x^k, x^{k+1}), \quad p^{k+1} \in \partial J(x^{k+1}).$$

¹Benning, Burger (2018).

$$p_i^{k+1} = p_i^k - \tau_i^k \frac{V(x_1^{k+1}, \dots, \mathbf{x}_i^{k+1}, x_{i+1}^k, \dots, x_n^k) - V(x_1^{k+1}, \dots, x_{i-1}^{k+1}, \mathbf{x}_i^k, \dots, x_n^k)}{\mathbf{x}_i^{k+1} - \mathbf{x}_i^k}.$$

- If V is continuous, and J is continuous and **strongly convex**, then updates are **well-defined**.
- If V is convex, updates are **unique**.
- Iterates $(x^k)_{k \in \mathbb{N}}$ converge to **Clarke stationary points** (under regularity assumption).
- Implications:
 - Can adapt pre-existing methods to incorporate bias (e.g. **sparsity**, variational regularisation problems)
 - Handles nonsmoothness better (e.g. $\|\cdot\|_1$ kinks)

Example: Bregman Itoh–Abe for linear systems (SOR)

(1/2)

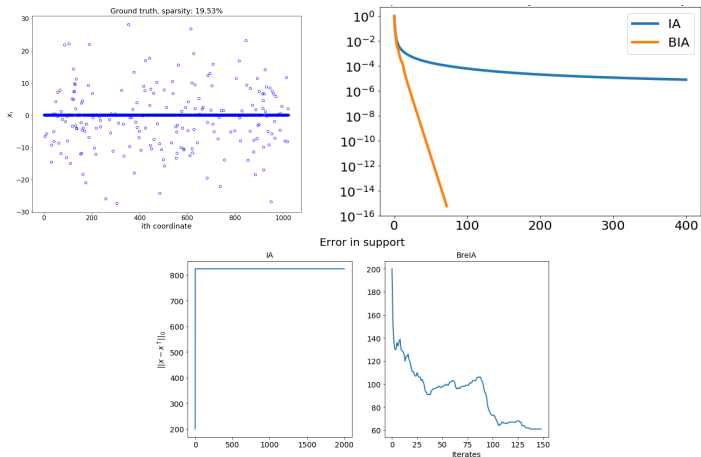


Figure: Top left: Ground truth. Top right: Normalised residual of data fidelity. Bottom: Error in support of iterate, $\text{supp}(x^k)$, to support of ground truth, $\text{supp}(x^*)$.

Example: Bregman Itoh–Abe for linear systems (SOR)

(2/2)

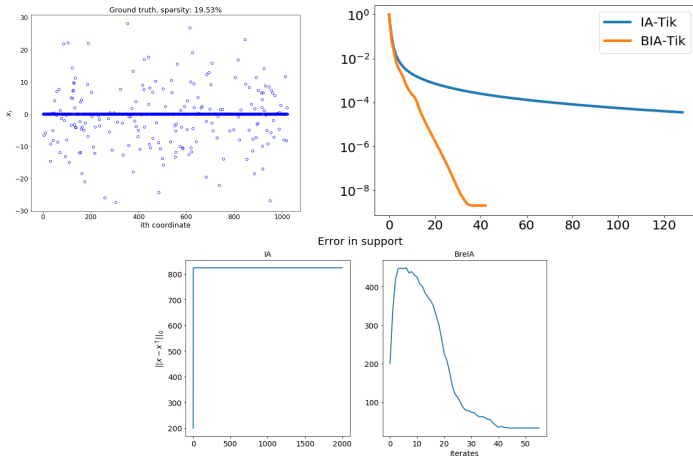


Figure: Top left: Ground truth. Top right: Normalised residual of data fidelity. Bottom: Error in support of iterate, $\text{supp}(x^k)$, to support of ground truth, $\text{supp}(x^*)$.

- Accelerate Itoh–Abe DG method.
- Apply discrete gradient methods to gradient flow under different metrics(e.g. optimal transport)

Thank you for your attention!

Relevant papers

- 1 Riis, Ehrhardt, Quispel, Schönlieb. *A geometric integration approach to nonsmooth, nonconvex optimisation*. (2018, preprint)
- 2 Ehrhardt, Riis, Ringholm, Schönlieb. *A geometric integration approach to smooth optimisation: Foundations of the discrete gradient method*. (2018, preprint)
- 3 Grimm, McLachlan, McLaren, Quispel, Schönlieb. *Discrete gradient methods for solving variational image regularisation models*. (2017)
- 4 Ringholm, Lazić, Schönlieb. *Variational image regularization with Euler's elastica using a discrete gradient scheme*. (2017)