

# Swarm-Based Gradient Descent Method for Non-Convex Optimization

Eitan Tadmor <sup>1,\*</sup>

## in honor of Albert Cohen

Nonlinear Approximation for High-Dimensional Problems  
Sorbonne University, July 2 2025



\* Joint works with J. Lu, A. Zenginoglu, Z. Ding, M. Guerra & Q. Li

# Gradient Descent method

- Objective function (or “loss function”)  $F(\cdot)$ ; state space  $\Omega \subset \mathbb{R}^d$ :

$$\text{Find } \underset{\mathbf{x} \in \Omega}{\operatorname{argmin}} F(\mathbf{x})$$

- Gradient Descent (GD) method (Cauchy<sup>2</sup> 1847)  $\min\{F(\mathbf{x}) : \nabla F(\mathbf{x}) = 0\}$

Proceed in gradient direction  $\nabla F^n = \nabla F(\mathbf{x}^n)$  discrete time steps  $t^{n+1} = t^n + h$

$$\mathbf{x}^{n+1} = \mathbf{x}^n - h \nabla F^n, \quad \mathbf{x}^n := \mathbf{x}(t^n)$$

- Gradient descent:  $\forall \lambda \in (0, 1), h \leq h_\lambda, F(\mathbf{x}^n - h \nabla F^n) \leq F(\mathbf{x}^n) - \lambda h |\nabla F^n|^2$

- Different protocols for time stepping  $h_\lambda = h(\mathbf{x}^n, \lambda)$ :

Backtracking, Adam ('momentum', Kingma & Ba (2017)), ...

AEGD ('energy', Liu (2021))...

---

<sup>2</sup> “I'll restrict myself here to outlining the principles underlying [my method], with the intention to come again over the same subject, in a paper to follow.”

A. Cauchy. Méthode générale pour la résolution des systèmes d'équations simultanées.  
C. R. Acad. Sci. Paris, 25:536-538, 1847.

# Swarm-based dynamics

- ☞ **Agents:** position  $\mathbf{x}_i^n = \mathbf{x}_i(t^n) \in \mathbb{R}^d$  with mass  $m_i^n = m_i(t^n) \in (0, 1]$
- ☞ **Communication:** adjust the masses according to relative heights  $\eta$ :

$$\left\{ \begin{array}{l} m_i^{n+1} = m_i^n - \eta_i^q m_i^n, \quad \eta_i := \frac{F(\mathbf{x}_i^n) - F_-(t^n)}{F_+(t^n) - F_-(t^n)}, \quad i \neq i_n \\ m_{i_n}^{n+1} = m_{i_n}^n + \sum_{i \neq i_n} \eta_i^q m_i^n, \quad i_n := \operatorname{argmin}_i F(\mathbf{x}_i^n) \quad F_-(t^n) = F(\mathbf{x}_{i_n}^n) \end{array} \right.$$

- Dynamic transfer of mass from high ground to lower ground
- Using  $\eta_i^q$  — the higher  $q$  the more tamed mass transfer

- ☞ **Protocol for time stepping:** dictated by relative masses  $\tilde{m}$ :

$$\mathbf{x}_i^{n+1} = \mathbf{x}_i^n - h(\mathbf{x}_i^n, \lambda \tilde{m}_i^{n+1}) \nabla F_i^n, \quad \tilde{m}_i^{n+1} = \frac{m_i^{n+1}}{m_+^{n+1}}, \quad m_+ := \max_i m_i$$

- Interplay between positions and masses

# Time stepping protocol — backtracking

- Protocol for choosing the step size (“learning rate”) — backtracking:

( $\star$ ) For “small enough”  $h$ :  $F(\mathbf{x}^n - h\nabla F^n) \leq F(\mathbf{x}^n) - \lambda h|\nabla F^n|^2$ ,  $\lambda \in (0, 1)$

Let  $h(\mathbf{x}^n, \lambda)$  – the largest for which ( $\star$ ) holds

· A key aspect:  $h(\cdot, \lambda\tilde{m})$  is decreasing function of the  $\tilde{m} = \frac{m}{m_+}$

$$F(\mathbf{x}_i^{n+1}) \leq F(\mathbf{x}_i^n) - \lambda\tilde{m}_i^{n+1} h_i^n |\nabla F_i^n|^2, \quad h_i^n := h(\mathbf{x}_i^n, \lambda\tilde{m}_i^{n+1})$$

☞ Dynamic distinction — ‘leaders’ and ‘explorers’...

heavy agents,  $m_i^{n+1} \sim m_+^{n+1} \rightsquigarrow$  small time steps; lead near critical point

light agents,  $m_i^{n+1} \ll m_+^{n+1} \rightsquigarrow$  large time steps; explore the landscape

- Each agent sheds an  $\eta$ -fraction of its mass – shifts to current minimizer

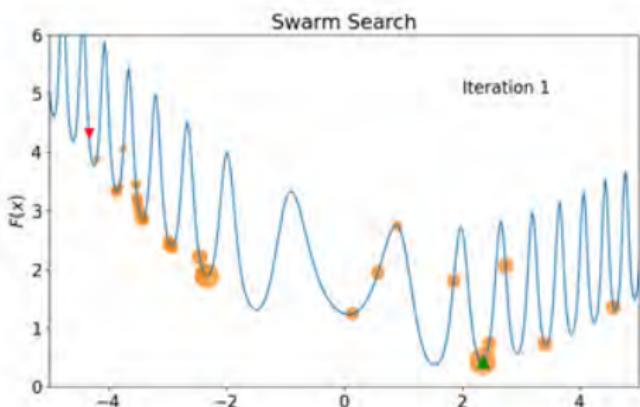
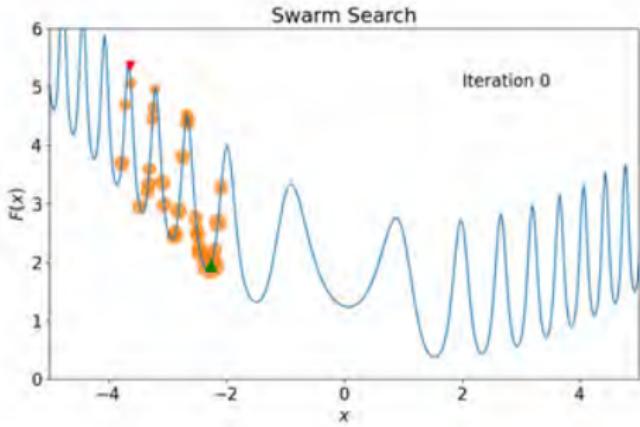


- Survival of the fittest: highest agent,  $\eta_{\max} = 1$ , is eliminated

- Conservation of mass  $\sum_i m_i^n = 1$ :

$\{m_i\}$  viewed as probabilities of agents to identify global minimum

# Why communication is important? survival of the fittest



# Alignment towards the minimal value

- Consensus Based Optimization (CBO)<sup>3</sup>:  
SDEs with steering towards the weighted min ('Laplace principle')
  - Biologically inspired methods —  
ant colony optimization<sup>4</sup>; artificial bee colony optimization<sup>5</sup>; firefly optimization<sup>6</sup>
- Wind-Driven Optimization (WDO)<sup>7</sup>: inspired by modeling pressure and wind
- Particle Swarm Optimization (PSO)<sup>8</sup>: explore the state space with randomized drifts towards the global best position
- Simulated Annealing (SA)<sup>9</sup>: exploration of state space is driven by noise terms that are 'cooled down' as time evolves

<sup>3</sup>Pinnau, Totzeck, Tse, & Martin, A CBO for global optimization, M3AS 27 (2017)  
Carrillo, Choi, Totzeck & Tse, An analytical framework for CBO, M3AS 28 (2018)

<sup>4</sup>Mohan & Baskaran, Ant colony optimization Expert Syst. Appl. 39 (2012)

<sup>5</sup>Karaboga et. al., Artificial bee colony (ABC) algorithm, Artif. Intell. Rev. 42 (2014)

<sup>6</sup>Yang, Firefly Algorithms for Multimodal Optimization. SAGA, V. 5792, (2009)

<sup>7</sup>Bayraktar et. al., The Wind Driven Optimization, IEEE Trans. Ant. Prop. 61 (2013)

<sup>8</sup>Kennedy & Eberhart, Particle Swarm Optimization. Proc. IEEE (1995)

Poli, Kennedy & Blackwell, Particle swarm optimization, Swarm Intell. 1 (2007)

<sup>9</sup>Kirkpatrick, Gelatt and Vecchi, Optimization by simulated annealing, Science (1983)

Chen et. al.. Accelerating nonconvex learning via replica exchange diffusion ICLR (2019)

SET tolerance parameters,  $\text{tolm}$ ,  $\text{tolmerge}$  and  $\text{tolres}$ ; adjustment parameter  $q > 0$   
INITIALIZE: #  $N$ ; positions  $\{\mathbf{x}_i^0\}$ ; masses  $m_i^0 = 1/N$ : optimal agent,  $i_0 = \operatorname{argmin}_i F(\mathbf{x}_i^0)$   
**for**  $n = 0, 1, 2, \dots$  **do**  
  Set  $F_{\min}^n = F(\mathbf{x}_{i_n}^n)$ ,  $F_{\max}^n = \max_i F(\mathbf{x}_i^n)$   
  **for**  $i = 1, \dots, N$  and  $i \neq i_n$  **do** % Mass transitions  
    **if**  $m_i^n < 1/N * \text{tolm}$  **then**  
      SET  $m_i^{n+1} = 0$   
      reduce the # of active agents:  $N \leftarrow N - 1$   
    **else**  $m_i^{n+1} = m_i^n - (\eta_i^n)^q m_i^n$  where  $\eta_i^n = \frac{F(\mathbf{x}_i^n) - F_{\min}^n}{F_{\max}^n - F_{\min}^n}$   
    **end if**  
  **end for**  
   $m_{i_n}^{n+1} = m_{i_n}^n + \sum_{i \neq i_n} (\eta_i^n)^q m_i^n$  % The total mass of the swarm is conserved  
  Compute  $m_+ = \max_i m_i^{n+1}$   
  **for**  $i = 1, \dots, N$  **do** % Gradient descent  
    Compute relative mass  $\tilde{m}_i^{n+1} := m_i^{n+1}/m_+$   
    Compute the step size  $h = h(\mathbf{x}_i^n, \lambda \tilde{m}_i^{n+1})$  according to algorithm 1 (backtracking).  
    MARCH:  $\mathbf{x}_i^{n+1} = \mathbf{x}_i^n - h \nabla F(\mathbf{x}_i^n)$   
  **end for**  
  MERGE agents if their distance  $< \text{tolmerge}$   
  SET the new optimal agent  $i_{n+1} = \operatorname{argmin}_i F(\mathbf{x}_i^{n+1})$ .  
  **if**  $\text{res} < \text{tolres}$  **then**  $\mathbf{x}_{\min} \leftarrow \mathbf{x}_{i_{n+1}}^{n+1}$  %  $\text{res} := |\mathbf{x}_{i_{n+1}}^{n+1} - \mathbf{x}_{i_n}^n|_2$   
  **end if**  
**end for**

## Quantifying the descent property

- Assume the Lip bound  $|\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})| \leq L|\mathbf{x} - \mathbf{y}|$

$$\begin{aligned} F(\mathbf{x}_i^n - h\nabla F_i^n) &\leq F(\mathbf{x}_i^n) - h|\nabla F_i^n|^2 + \frac{h^2}{2}L|\nabla F_i^n|^2 \\ &\leq F(\mathbf{x}_i^n) - \left(1 - \frac{h}{2}L\right)h|\nabla F_i^n|^2, \quad \nabla F_i^n := \nabla F(\mathbf{x}_i^n) \end{aligned}$$

hence set  $h_\lambda : \left(1 - \frac{1}{2}h_\lambda L\right) \geq \lambda \tilde{m}_i^{n+1}$

$$\forall h \leq h_\lambda : F(\mathbf{x}_i^n - h\nabla F_i^n) \leq F(\mathbf{x}_i^n) - \lambda \tilde{m}_i^{n+1} h |\nabla F_i^n|^2$$

- Certainly holds for  $h$  is small enough : OK if  $h_\lambda < \frac{2}{L}(1 - \lambda \tilde{m}_i^{n+1})$
- And this does not mean  $h_\lambda$  need to be small ...
- But no access to  $L$  ...

# Backtracking — implementation

---

## Algorithm 1 Backtracking Line Search

---

SET the descent parameter  $\lambda \in (0, 1)$ , and shrinkage parameter,  $\gamma \in (0, 1)$

Compute the relative masses  $\tilde{m}_i^{n+1} = \frac{m_i^{n+1}}{m_+^{n+1}}$

INITIALIZE large step size  $h = h_0 > \frac{2}{L}$ .

**while**  $F(\mathbf{x}_i^n - h \nabla F_i^n) \circlearrowleft F(\mathbf{x}_i^n) - \lambda \tilde{m}_i^{n+1} h |\nabla F_i^n|^2$  **do** % weighted descent  
     $h \leftarrow \gamma h$  % backtracking

**end while**

SET  $h \mapsto h_i^n := h(\mathbf{x}_i^n, \lambda \tilde{m}_i^{n+1}) \geq \frac{2\gamma}{L}(1 - \lambda \tilde{m}_i^{n+1})$

---

- Lower-bound:  $\frac{h_i^n}{\gamma} \geq \frac{2}{L}(1 - \lambda \tilde{m}_i^{n+1}), \quad |\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})| \leq L|\mathbf{x} - \mathbf{y}|$

- Descent property:  $F(\mathbf{x}_i^{n+1}) \leq F(\mathbf{x}_i^n) - \frac{2\gamma}{L}(1 - \lambda \tilde{m}_i^{n+1}) \lambda \tilde{m}_i^{n+1} |\nabla F_i^n|^2$ 
  - Noll & Rondepierre (2013): backtracking protocol with memory

Descent estimates:  $F(\mathbf{x}_i^{n+1}) \leq F(\mathbf{x}_i^n) - \frac{2\gamma}{L}(1 - \lambda \tilde{m}_i^{n+1})\lambda \tilde{m}_i^{n+1} |\nabla F_i^n|^2$

- Descent property: How small  $\tilde{m}_i^{n+1}$  can get for SBGD minimizers?
- Heaviest agent at  $\mathbf{X}_+^n$  w/relative mass  $\tilde{m} = 1$ ; compare w/minimizer  $\mathbf{X}_-^n$ :

$$\tilde{m}_{i_n}^{n+1} = 1 \rightsquigarrow F(\mathbf{X}_-^{n+1}) \leq F(\mathbf{X}_+^{n+1}) \leq F(\mathbf{X}_+^n) - \frac{2\gamma}{L}(1 - \lambda)\lambda |\nabla F(\mathbf{X}_+^n)|^2$$

A. Scenario (i): if  $F(\mathbf{X}_+^n) \leq F(\mathbf{X}_-^n) + \frac{\gamma}{L}(1 - \lambda)\lambda |\nabla F(\mathbf{X}_+^n)|^2$  then ...

$$F(\mathbf{X}_-^{n+1}) \leq F(\mathbf{X}_+^{n+1}) \leq F(\mathbf{X}_+^n) - \frac{2\gamma}{L}(1 - \lambda)\lambda |\nabla F_+^n|^2 \leq F(\mathbf{X}_-^n) - \frac{\gamma}{L}(1 - \lambda)\lambda |\nabla F_+^n|^2$$

B. Scenario (ii):  $F(\mathbf{X}_+^n) > F(\mathbf{X}_-^n) + \frac{\gamma}{L}(1 - \lambda)\lambda |\nabla F(\mathbf{X}_+^n)|^2$

The minimizer at  $\mathbf{X}_-^n$  must be different from the heaviest at  $\mathbf{X}_+^n$

hence must gain  $\eta$ -fraction mass from heaviest at  $\mathbf{X}_+^n$ :  $m_-^{n+1} > \eta_+^n m_+^n$

$$\rightsquigarrow \tilde{m}_-^{n+1} > \eta_+^n \geq \frac{1}{M^q} (F(\mathbf{X}_+^n) - F(\mathbf{X}_-^n))^q > \frac{\gamma}{ML}(1 - \lambda) |\nabla F_+^n|^2$$

$$F(\mathbf{X}_-^{n+1}) \leq F(\mathbf{X}_-^n) - \frac{2\gamma}{L}(1 - \lambda \tilde{m}_-^{n+1})\lambda \tilde{m}_-^{n+1} |\nabla F_-^n|^2$$

$$\leq F(\mathbf{X}_-^n) - \left( \frac{\gamma(1 - \lambda)\lambda}{ML} \right)^q |\nabla F_+^n|^{2q} \cdot \frac{\gamma}{L} \lambda |\nabla F_-^n|^2$$

☞ telescoping sum A+B ...

# Convergence. I (global)

- Recall the mass transfer:  $m_i^{n+1} = m_i^n - \eta_i^q m_i^n$ ,  $\eta_i := \frac{F(\mathbf{x}_i^n) - F_{\min}(t^n)}{F_{\max}(t^n) - F_{\min}(t^n)}$   
 $\rightsquigarrow F(\mathbf{X}_{-}^{n+1}) \leq F(\mathbf{X}_{-}^n) - \lambda h_- |\nabla F_{-}^n|^{2(q+1)}$

Theorem

$$\sum_{n=n_0}^{\infty} \min \left\{ |\nabla F(\mathbf{X}_+^n)|, |\nabla F(\mathbf{X}_-^n)| \right\}^{2(q+1)} < C \min_i F(\mathbf{x}_i^0)$$

☞ Limit set of minima of equi-height:  $\mathbf{X}_{-}^{n_\alpha} \xrightarrow{\alpha \rightarrow \infty} \mathbf{X}_\alpha^*$ ,  $\nabla F(\mathbf{X}_{-}^{n_\alpha}) \xrightarrow{\alpha \rightarrow \infty} 0$

- “Generic” descent property  $F(\mathbf{X}_{-}^{n+1}) \leq F(\mathbf{X}_{-}^n) - \lambda h_- |\nabla F_{-}^n|^2$

$$\sum_{n=n_0}^{\infty} \min \left\{ |\nabla F(\mathbf{X}_+^n)|, |\nabla F(\mathbf{X}_-^n)| \right\}^2 < C \min_i F(\mathbf{x}_i^0)$$

## Convergence. II (local)

- Lojasiewicz bound<sup>10a</sup>:  $\exists \mathcal{N}, \beta$  s.t.  $\mu|F(\mathbf{x}) - F(\mathbf{x}^*)| \leq |\nabla F(\mathbf{x})|^\beta$ ,  $\forall \mathbf{x} \in \mathcal{N}(\mathbf{x}^*)$ 
  - quantifies “Flatness” in terms of “ $\beta = \frac{m+1}{m} \in (1, 2]$ ” for SA  $F$ ’s<sup>10a</sup>
  - basin of convexity:  $m = 1 \rightsquigarrow \beta = 2$ ;
- Generic descent property<sup>10b</sup> at  $\mathbf{X}_-^{n_\alpha}$ :  $F_-^{n_\alpha+1} \leq F_-^{n_\alpha} - \lambda h_- |\nabla F_-^{n_\alpha}|^2 \rightsquigarrow$  $F_-^{n_\alpha+1} - F_\alpha^* \leq F_-^{n_\alpha} - F_\alpha^* - \lambda h_- (\mu(F_-^{n_\alpha} - F_\alpha^*))^{\frac{2}{\beta}}$ ,  $\mathbf{X}_-^{n_\alpha} \in \mathcal{N}(\mathbf{X}_\alpha^*)$ .

### Theorem

Analytic Semi-algebraic loss function  $F$  with minimal flatness  $\beta \in (1, 2]$ .  
SBGD minimizers  $\{\mathbf{X}_-^n\}_{n \geq 0}$ .  $\exists C = C(\gamma, \lambda, \mu)$ , such that

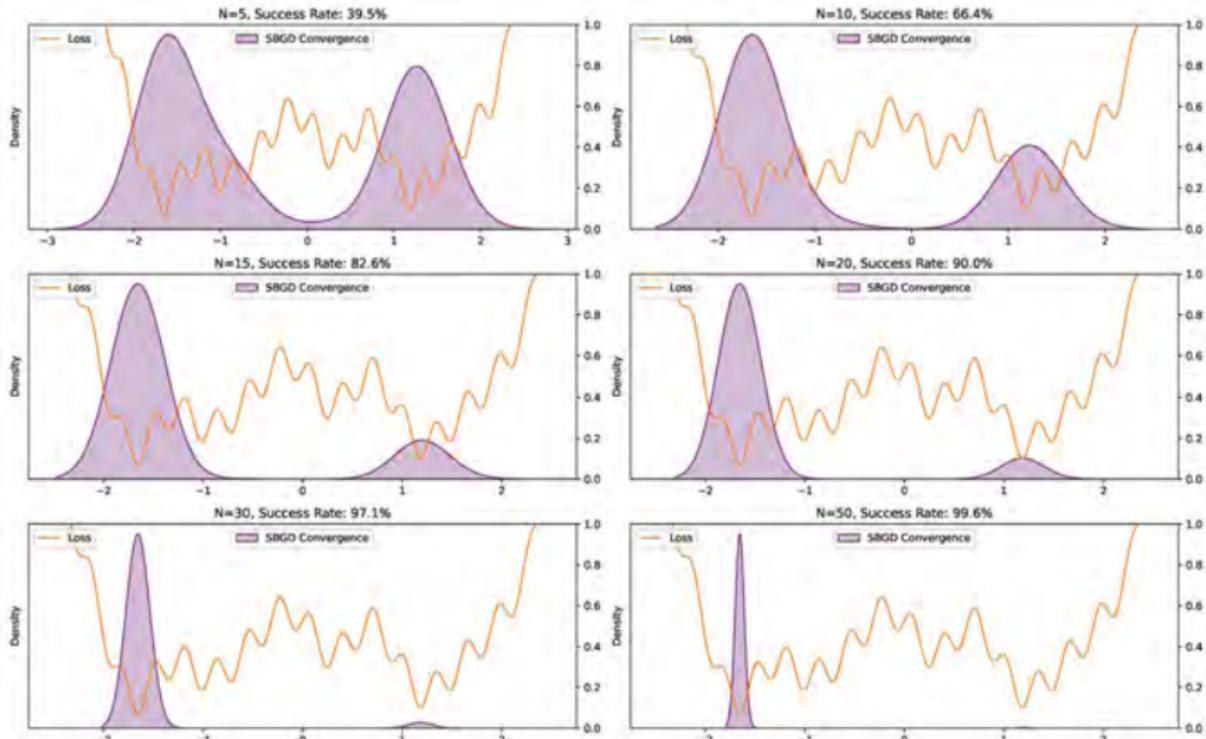
$$F(\mathbf{X}_-^{n_\alpha}) - F(\mathbf{X}_\alpha^*) \begin{cases} \leq \left(1 - \frac{2\mu\gamma\lambda(1-\lambda)}{L}\right)^{n_\alpha} (F_{\min}(\mathbf{x}^0) - F(\mathbf{x}^*)), & \beta = 2 \\ \lesssim \left(\frac{C}{n_\alpha}\right)^{\frac{\beta}{2-\beta}}, & \beta \in (1, 2) \end{cases}$$

<sup>10a</sup> Lojasiewicz IHES1965; J. Bolte et. al., on Kurdyka-Lojasiewicz ineq. TAMS 2010

<sup>10b</sup> Absil et. al., Convergence of descent iterations for analytic functions, SJOPT2005

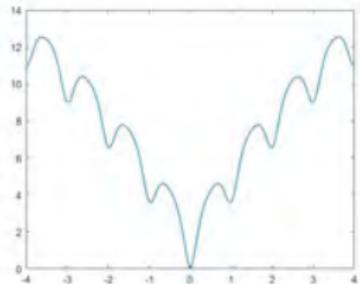
# Another convincing 1D example (Boffi-Slotine (2020))

$$F(x) = \frac{1}{C} \left( x^4 - 4x^2 + \frac{1}{5}x + \frac{2}{5} \left( 3 \sin(20x) - \frac{7}{2} \sin(2\pi x) + \cos\left(\frac{10\pi x}{3}\right) \right) \right)$$



# Simulations - one-dimensional

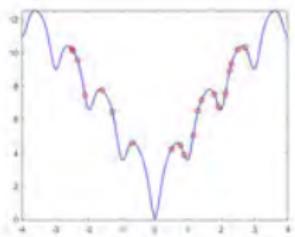
$$F(x) = -20e^{-0.2|x-B|} - e^{\cos(2\pi(x-B))} + 20 + e + C$$



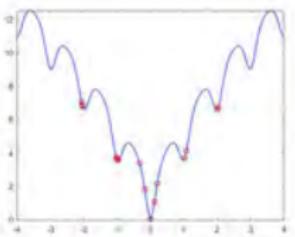
$x^* = B$		N=10	N=20	N=30	N=10	N=20	N=30
$B = 0$	success rate $\mathbb{E} x_{SOL} - x^* ^2$	100% $8.42e^{-10}$	100% $8.37e^{-10}$	100% $3.38e^{-10}$	100% $8.60e^{-10}$	100% $1.36e^{-9}$	100% $1.29e^{-9}$
$B = 5$	success rate $\mathbb{E} x_{SOL} - x^* ^2$	100% $8.41e^{-10}$	100% $7.58e^{-10}$	100% $5.01e^{-10}$	100% $8.51e^{-10}$	100% $1.25e^{-9}$	100% $1.21e^{-9}$
$B = 15$	success rate $\mathbb{E} x_{SOL} - x^* ^2$	98.5% $1.41e^{-2}$	100% $4.69e^{-3}$	100% $8.27e^{-10}$	46.5% $6.44e^{+1}$	75.0% $2.26e^{+1}$	85.5% $1.14e^{+1}$
$B = 25$	success rate $\mathbb{E} x_{SOL} - x^* ^2$	45.5% $1.48e^{+2}$	89.0% $1.49e^{+1}$	98.5% $3.28e^{-1}$	0% $4.86e^{+2}$	0% $4.50e^{+2}$	0% $4.38e^{+2}$

**Table:**  $m = 200$  runs of shifted 1D Ackley function. Left: SBGD. Right: GD( $\eta=0$ )(no communication). Backtracking descent parameter  $\lambda = 0.2$ ; shrinkage parameter  $\gamma = 0.9$

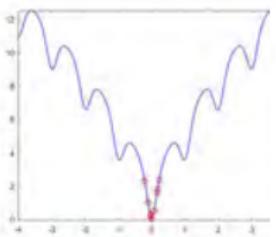
# Dynamics of SBGD agents; 1D Ackley simulation ( $B = C = 0$ )



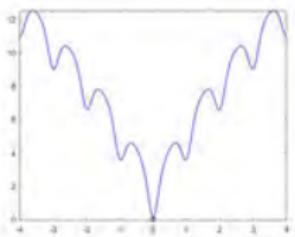
(a)  $n = 1, N = 20$



(b)  $n = 3, N = 17$



(c)  $n = 8, N = 9$

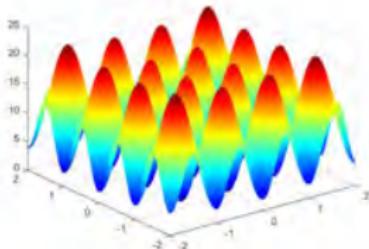


(d)  $n = 15, N = 1$

SBGD 1D Ackley

# Simulations – two-dimensional

$$F_{\text{Rastrigin}}(x) = \frac{1}{d} \sum_{i=1}^d \left\{ (\mathbf{x}_B)_i^2 - 10 \cos(2\pi(\mathbf{x}_B)_i) + 10 \right\} + C$$



N	10	20	30
SBGD <sub>1</sub>	52.1%	70.0%	75.8%
SBGD <sub>2</sub>	60.1%	84.3%	91.0%
GD( $h \equiv 0.004$ )	50.5%	70.0%	78.1%
GD( $\eta = 0$ )	51.0%	70.8%	79.3%
Adam(0.8)	29.6%	46.8%	65.5%
Adam(0.2)	40.9%	65.3%	79.4%

N	10	20	30
SBGD <sub>1</sub>	49.2%	67.0%	72.7%
SBGD <sub>2</sub>	46.7%	81.9%	89.6%
GD( $h \equiv 0.004$ )	0.0%	0.0%	0.0%
GD( $\eta = 0$ )	2.4%	4.3%	5.9%
Adam(0.8)	31.3%	49.2%	66.9%
Adam(0.2)	0.0%	0.0%	0.0%

**Table:**  $m = 1000$  runs for 2D Rastrigin function, success rates of SBGD compared with GD( $h$ ) GD( $\eta = 0$ ) and Adam methods (first component).

Left: initial data uniformly generated in  $[-3, 3]^2$ . Right: uniformly generated in  $[-3, -1]^2$ .

Backtracking parameters: descent parameter  $\lambda = 0.8$ ; shrinkage parameter  $\gamma = 0.9$

# Simulations – 20-dimensions

$\mathbf{x}_* = B$		N=50	N=100	N=200	N=50	N=100	N=200
$B = 0$	SBGD <sub>2</sub>	9.00e-07	2.02e-07	1.43e-07	4.03e-01	4.98e-01	3.91e-01
	GD( $\eta=0$ )	1.18e-01	6.90e-02	1.88e-02	2.96e-01	3.25e-01	4.35e-01
	Adam(0.5)	3.37e-03	4.03e-03	4.96e-03	3.75e-01	2.86e-01	4.14e-01
$B = 3$	SBGD <sub>2</sub>	1.65e-06	1.51e-06	7.96e-07	5.07	3.71	2.92
	GD( $\eta=0$ )	7.64	6.72	5.67	9.43	9.02	8.57
	Adam(0.5)	2.14e-01	1.22e-01	1.27e-01	9.02	8.63	8.15
$B = 5$	SBGD <sub>2</sub>	4.15	1.07	2.36e-01	3.53	2.45	1.29
	GD( $\eta=0$ )	17.99	17.59	17.11	18.02	17.64	17.18
	Adam(0.5)	11.01	9.27	7.52	17.13	16.73	16.3

**Table:**  $m = 1000$  runs.  $\mathbb{E}[\mathbf{x}_{SOL} - \mathbf{x}^*|$  (first component).

Left: 20-dimensional shifted Ackley. Right: for 20-dimensional shifted Rastrigin.

Backtracking parameters: descent parameter  $\lambda = 0.2$ ; shrinkage parameter  $\gamma = 0.9$

# Multi-D SBRD – “Where No One Has Gone Before”

- The need to enhance “space exploration” in high-D problems
- Time stepping:  $\mathbf{x}_i^{n+1} = \mathbf{x}_i^n - h\mathbf{p}_i^n$
- ☞ Allow general set of directions  $\{\mathbf{p}_i^n\}$  such that

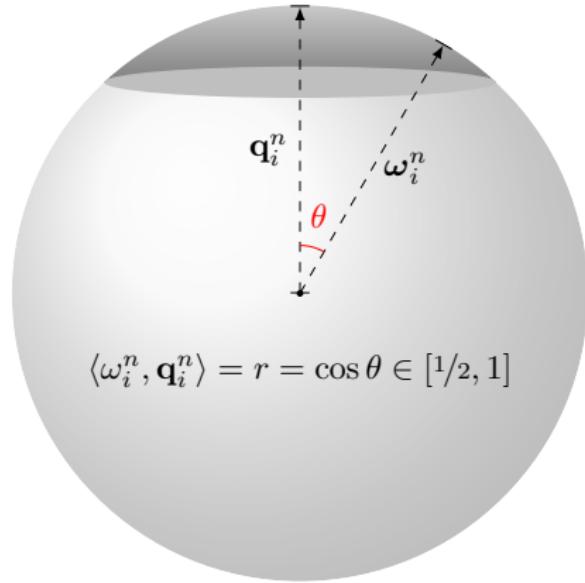
$$(*) \quad \langle \mathbf{p}_i^n, \nabla F(\mathbf{x}_i^n) \rangle \geq \frac{1 + \tilde{m}_i^{n+1}}{2} |\nabla F(\mathbf{x}_i^n)|^2 \geq \frac{1}{2} |\nabla F(\mathbf{x}_i^n)|^2$$

This secures the “one-half” descent property:

$$F(\mathbf{x}_i^{n+1}) \leq F(\mathbf{x}_i^n) - h \langle \mathbf{p}_i^n, \nabla F(\mathbf{x}_i^n) \rangle + \frac{h^2}{2} L |\nabla F(\mathbf{x}_i^n)|^2 \leq F(\mathbf{x}_i^n) - \frac{1}{2}(1 - Lh)h |\nabla F(\mathbf{x}_i^n)|^2$$

- Heterogeneity of orientations  $\leadsto$  effective exploration of the ambient space
- To secure (\*): choose a Random orientation  $\mathbf{p}_i^n = |\nabla F(\mathbf{x}_i^n)| \mathbf{\omega}_i^n$  such that

$$\langle \mathbf{\omega}_i^n, \mathbf{q}_i^n \rangle = \text{rand}, \quad \text{rand} \in \mathcal{N}\left(\frac{1 + \tilde{m}_i^{n+1}}{2}, 1\right), \quad \mathbf{q}_i^n := \frac{\nabla F(\mathbf{x}_i^n)}{|\nabla F(\mathbf{x}_i^n)|}$$



$$\langle \omega_i^n, \mathbf{q}_i^n \rangle = r = \cos \theta \in [1/2, 1]$$

**Figure:** Spherical cap of  $\mathbb{S}^{d-1}$ -sphere centered at the gradient orientation  
 $\mathbf{q}_i^n = \frac{\nabla F(\mathbf{x}_i^n)}{|\nabla F(\mathbf{x}_i^n)|}$ . “Opening”:  $\cos(\theta) = \frac{1 + \lambda \tilde{m}_i^{n+1}}{2}$  up to  $60^\circ$

## 2D Simulations — SBRD 2D Ackley (shift: $\mathbf{x}_B := \mathbf{x} - B\mathbf{1}$ )

- $F_{\text{Ackley}}(\mathbf{x}) = -20\exp\left\{\frac{-0.2}{\sqrt{d}}|\mathbf{x}_B|\right\} - \exp\left\{\frac{1}{d}\sum_i \cos(2\pi(\mathbf{x}_B)_i)\right\} + 20 + e + C$

SBRD 2D Ackley

## Algorithm 2 Random descent direction

```
if  $1 - \tilde{m}_i^{n+1} \leq tolrand$  then
    set  $\mathbf{p}_i^n = \nabla F(\mathbf{x}_i^n)$  % for heaviest agent use the gradient direction
end if
else set  $\mathbf{q}_i^n = \frac{\nabla F(\mathbf{x}_i^n)}{|\nabla F(\mathbf{x}_i^n)|}$ 
Choose random  $r$ ,  $\frac{1}{2}(1 + \tilde{m}_i^n) < \text{rand} < 1$ 
    % Random  $X$  in  $\mathbb{S}^{d-1}$ : spherical cap centered at the north pole  $\mathbf{z} := (0, \dots, 0, 1)$ 
for  $i = 1, \dots, d-1$  do %  $\frac{Y(i)}{|\mathbf{Y}|}$  is random in  $\mathbb{S}^{d-2}$  ball ( $= \mathbb{S}^{d-1}$  projected to plane)
     $Y(i)$  random  $\in \mathcal{N}(0, 1)$ 
end for
for  $i = 1, \dots, d-1$  do % from  $\mathbb{S}^{d-2}$  to cap 'up' in  $\mathbb{S}^{d-1}$ 
    Set  $X(i) = \sqrt{1 - \text{rand}^2} \frac{Y(i)}{|\mathbf{Y}|}$ 
end for
Set  $\mathbf{X} = (X(1), \dots, X(d-1), X(d))$ ,  $X(d) = \text{rand}$ 
% From  $\mathbf{X}$  to  $\omega_i^n$  - random orientation in a spherical cap centered at  $\mathbf{q}_i^n$ 
% Reflect:  $\omega_i^n = \mathbb{P}_i^n \mathbf{X}$ ,  $\mathbb{P}_i^n := \mathbb{I} - 2 \frac{(\mathbf{q}_i^n - \mathbf{z})(\mathbf{q}_i^n - \mathbf{z})^\top}{|\mathbf{q}_i^n - \mathbf{z}|^2}$ 
if  $1 - \mathbf{q}_i^n(d) \neq 0$  then
    Set  $\mathbf{v}_i^n = \mathbf{q}_i^n - \mathbf{z}$ 
    Set  $\omega_i^n = \mathbf{X} - 2 \frac{\langle \mathbf{v}_i^n, \mathbf{X} \rangle}{|\mathbf{v}_i^n|^2} \mathbf{v}_i^n$  % Note the simplification:  $|\mathbf{v}_i^n|^2 = 2(1 - \mathbf{q}_i^n(d))$ 
else  $\omega_i^n = \mathbf{X}$ 
end if
```

# Why randomization is important?

$d \backslash N$	5	10	25	50	100
4	27.2%	50.4%	91.2%	99.2%	100.0%
8	17.2%	53.0%	100.0%	100.0%	100.0%
12	2.4%	26.0%	96.0%	100.0%	100.0%
16	0.0%	0.0%	0.0%	1.0%	2.4%
20	0.0%	0.0%	0.0%	0.0%	0.0%

$d \backslash N$	5	10	25	50	100
4	24.0%	38.6%	77.8%	99.8%	100.0%
8	12.2%	26.6%	60.0%	92.8%	100.0%
12	2.2%	11.6%	56.2%	88.8%	99.6%
16	0.0%	0.4%	26.6%	60.2%	88.0%
20	0.0%	0.0%	0.4%	5.8%	21.6%

**Table:** Success rates of SBGD (top) vs. SBRD (bottom) for  $D$ -dimensional Ackley based on  $m = 500$  runs. Backtracking parameters  $\lambda = 0.2$  and  $\gamma = 0.9$

# Mass transfer — dependence on $q \gg 1$

$d \backslash N$	10	25		50		100		
$d$	8	4	8	4	8	4	8	4
<b>Ackley</b>								
14	<b>4.2%</b>	3.8%	<b>66.9%</b>	60.0%	100.0%	100.0%	100.0%	100.0%
16	0.1%	0.3%	38.4%	38.3%	<b>99.8%</b>	95.0%	100.0%	100.0%
18	0.0%	0.0%	14.6%	16.3%	<b>87.3%</b>	79.7%	100.0%	99.6%
20	0.0%	0.0%	1.0%	1.4%	<b>30.7%</b>	25.1%	<b>84.7%</b>	74.5%
<b>Rastrigin</b>								
2	<b>42.5%</b>	37.4%	99.0%	98.8%	100.0%	100.0%	100.0%	100.0%
3	6.2%	6.5%	<b>32.4%</b>	29.0%	<b>80.1%</b>	74.3%	<b>99.1%</b>	98.0%
4	0.8%	1.0%	4.7%	4.9%	<b>14.3%</b>	11.5%	<b>35.2%</b>	30.3%
5	0.2%	0.1%	1.1%	0.9%	<b>3.0%</b>	1.7%	<b>4.8%</b>	3.7%

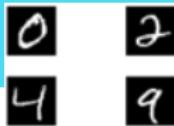
**Table:** Success rates of SBRD vs. SBGD for global optimization of the  $d$ -dimensional objective functions with mass transfer parameter  $q = 4$  and  $q = 8$ .

## 2D Simulations — SBRD 2D Rastrigin

- $F_{\text{Rastrigin}}(\mathbf{x}) = \frac{1}{d}|\mathbf{x}_B|^2 - \frac{10}{d} \sum_i \cos(2\pi(\mathbf{x}_B)_i) + 10 + C, \quad \mathbf{x}_B := \mathbf{x} - B\mathbf{1}$

SBRD 2D Rastrigin

# SBRD in $d = 1976^{11}$



- Training on MNIST:

- MNIST-4-4 Architecture ( $N = \text{batch size}$ ):

Layer type	Output shape	#Params.
1-3 Conv2d-1	[ $N, 4, 14, 14$ ]	16
MaxPool2d-2	[ $N, 4, 7, 7$ ]	0
Flatten-3	[ $N, 196$ ]	0
FullyConn-4	[ $N, 10$ ]	1960

- Loss function  $I(\mathbf{x}, \mathbf{y}) = -\frac{1}{N} \sum_{n=0}^{N-1} \sum_i y_i \cdot \log(s(x_i))$ ,  $s(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$

Trial	Optimizer	$N$	$\eta$	$h_0$	$\gamma$	$\lambda$	$q$
1-8 1	SGD	1	0.1	-	-	-	-
2	NSGD (multi-particle)	200	0.1	-	-	-	-
3	NGDBT (multi-backtracking)	200	-	1.0	0.25	0.2	-
4	SBRD	200	-	1.0	0.25	0.2	1.0

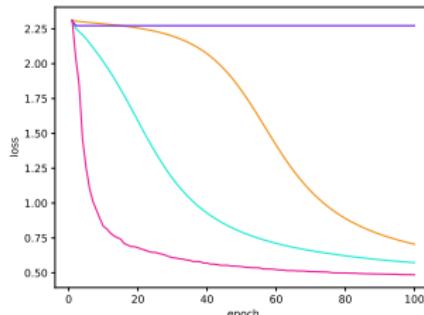


Figure: —, SGD; —, NSGD; —, NGDBT; —, SBRD

<sup>11</sup>with R. Leonard

# SBGD meets Simulated Annealing

- Simulated Annealing (SA) driven by Langevin diffusion:

$$d\mathbf{x}_i^t = -\nabla F(\mathbf{x}_i^t) dt + \sqrt{2\sigma_i^t} dW_i^t, \quad i = 1, 2, \dots, N,$$

- $\{W_i^t\}$  independent Brownian motions,  $i = 1, 2, \dots, N$ ;
- $\sigma_i^t$  is the scaled “temperature” or annealing-rate
- Key feature of one-particle SA — protocol for ‘cooling down’<sup>12</sup>  
Properly cooling down with annealing rate  $\sigma_i^t \propto c/\sqrt{\log t}$
- This is where the masses,  $\{m_i^t\}$ , come into play:

$$dm_i^t = -m_i^t \left( F(\mathbf{x}_i^t) - \bar{F}_N^t \right) dt, \quad i = 1, 2, \dots, N,$$

with provisional minimum:  $\bar{F}_N^t := \frac{\sum_{i=1}^N m_i^t F(\mathbf{x}_i^t)}{\sum_{i=1}^N m_i^t}$

- ☞ Let the swarm decide how to cool-down  $\sigma_i^t = \sigma(m_i^t) \downarrow$

---

<sup>12</sup>Gidas, ... convergence of the annealing algorithm, JSP (1985)  
Geman and Hwang, Diffusions for global optimization. SICOPT (1986)

# Why provisional minimum?

- Why provisional minimum?  $\bar{F}_N^t := \frac{\sum_{i=1}^N m_i^t F(\mathbf{x}_i^t)}{\sum_{i=1}^N m_i^t}$  why not  $F(\mathbf{x}^*) = \min_{\mathbf{x}} F(\mathbf{x})$ ?  
☞ Expected that for  $t_\epsilon \lesssim \frac{1}{\epsilon}$ ,  $N \gtrsim \frac{1}{\epsilon^2}$  there holds  $\mathbb{E} [\bar{F}_N^t - \min_{\mathbf{x}} F(\mathbf{x})] \leq \epsilon$
- Mean-field: empirical measure  $\mu_N^t(\mathbf{x}, m) = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_i^t}(\mathbf{x}) \otimes \delta_{m_i^t}(m) \rightarrow \mu_t(\mathbf{x}, m)$ :  
$$\partial_t \mu = \nabla_{\mathbf{x}} \cdot (\mu \nabla F) + (F(\mathbf{x}) - \bar{\mathcal{F}}_\mu^t) \partial_m(m\mu) + \sigma(m) \Delta_{\mathbf{x}} \mu \quad \bar{\mathcal{F}}_\mu^t := \frac{\iint m F(\mathbf{x}) d\mu^t(\mathbf{x}, m)}{\iint m d\mu^t(\mathbf{x}, m)}$$
- Indeed:  $\lim_{N \rightarrow \infty} \mathbb{E} [W_2(\mu^t, \mu_N^t)] = 0, \quad \mathbb{E} [\left| \bar{F}_N^t - \bar{\mathcal{F}}_\mu^t \right|] < \frac{C_t}{\sqrt{N}}$   
and for any  $\epsilon > 0$  there exists  $t_\epsilon \lesssim \frac{1}{\epsilon}$  such that  $\bar{\mathcal{F}}_\mu^t < F_* + \epsilon$

# Hydrodynamic description

- Density,  $\rho = \rho(t, \mathbf{x}) := \int \mu^t(\mathbf{x}, m) dm$  and Momentum,  $\rho v := \int m \mu^t(\mathbf{x}, m) dm$
- First two moments of the mean-field

$$\begin{cases} \rho_t + \nabla_{\mathbf{x}} \cdot (\rho \nabla F) = \Delta_{\mathbf{x}} \int \sigma(m) \mu^t(\mathbf{x}, m) dm, \\ (\rho v)_t + \nabla_{\mathbf{x}} \cdot (\rho v \nabla F) = (\bar{F}(t) - F(\mathbf{x})) \rho v + \Delta_{\mathbf{x}} \int m \sigma(m) \mu^t(\mathbf{x}, m) dm. \end{cases}$$

Provisional minimum  $\bar{F}(t) = \int F(\mathbf{x}) \rho v(t, \mathbf{x}) d\mathbf{x}$  (normalized mass/momentum= 1)

☞ Drift-diffusion equation for the velocity  $v$ ,

$$\rho_t + \nabla_{\mathbf{x}} \cdot (\rho \nabla F) = \Delta_{\mathbf{x}}(\sigma(t, \mathbf{x}) \rho(t, \mathbf{x})),$$

$$v_t + \nabla F \cdot \nabla_{\mathbf{x}} v = (\bar{F}(t) - F(\mathbf{x})) v + \frac{1}{\rho} \Delta_{\mathbf{x}}(\sigma(t, \mathbf{x}) \rho v(t, \mathbf{x})), \quad \mathbf{x} \in \text{supp}\{\rho(t, \cdot)\}.$$

- The decent estimate — expected  $\rho(t, \mathbf{x}) \xrightarrow{t \rightarrow \infty} \delta(\mathbf{x} - \mathbf{x}^*)$ ,  $v(t, \mathbf{x}) \xrightarrow{t \rightarrow \infty} \mathbb{1}(\mathbf{x}^*)$

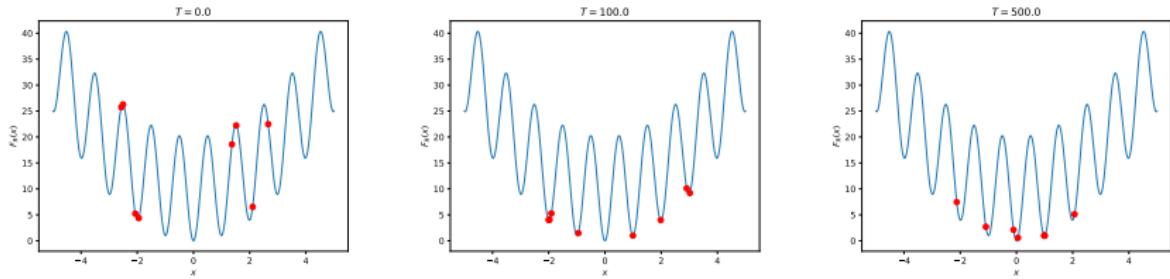
$$\frac{d}{dt} \int (F(\mathbf{x}) - F_*) \rho(t, \mathbf{x}) d\mathbf{x} = - \int |\nabla F(\mathbf{x})|^2 \rho(t, \mathbf{x}) d\mathbf{x} + \int \Delta_{\mathbf{x}} F(\mathbf{x}) \sigma(t, \mathbf{x}) \rho(t, \mathbf{x}) d\mathbf{x},$$

Communication —  $\sigma(\cdot) \geq 0$  is the key; observe:  $\sigma(t, \mathbf{x})$  should vanish as  $\mathbf{x} \rightarrow \mathbf{x}^*$

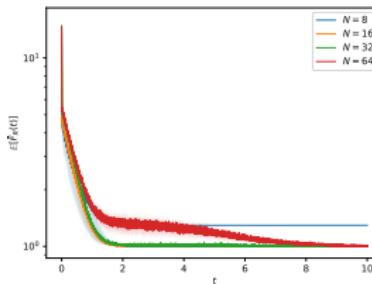
# SBGD-SA simulations: Rastrigin function for $d = 1$ :

$$\sigma(m) = \begin{cases} 2\exp(m/(m - 1/8)) & m < 1/8 \\ 0, & m \geq 1/8 \end{cases}; N = 8, h = 10^{-4}, T = 500$$

- Particles eventually find the global minimum after a long time



Convergence of  $\bar{F}_N^t$  in finding the global minimum is significantly faster



Mean-field convergence: 20 trajectories of  $\bar{F}_N^t$ ; bigger  $N \mapsto$  a lower value of  $F_*$



THANK YOU, Albert