# On Time-stepping Methods for Gradient-flow Optimization

#### Arash Sarshar

Computational Science Laboratory "Compute the Future!" Department of Computer Science, Virginia Polytechnic Institute and State University Blacksburg, VA 24060

CAIMS/SCMAI, June, 2022







Dr. Adrian Sandu Director of CSL

Email: sandu@cs.vt.edu Location: KWII 2222 Phone: (540) 231-2193 Fax: (540) 231-6075



Abhinab Bhattac

Email: abhinab93@y Location: KWII 2201 Research Interests: Data assimilation, numerical methods



Andrey Popov

Email: apppoy@vt.edu Location: KWII 2201 Research Interests: Large-scale dynamical systems, data assimilation, time integration



A. Sarshar CAIMS/SCMAI 2022 . [1/21] 🕐 csl.cs.vt.edu



(Stochastic) Gradient descent and gradient-flow



$$\arg \min_{\theta \in \mathbb{R}^{d}} L(\theta), \qquad (1)$$
  
$$\theta^{k+1} = \theta^{k} - \alpha \nabla_{\theta} L(\theta^{k}), \qquad (2)$$

$$\frac{\partial \theta}{\partial t} = -\nabla_{\theta} \mathbf{L}(\theta) \tag{3}$$



A. Sarshar CAIMS/SCMAI 2022 Introduction. [2/21]



# Inertial optimization methods

$$\underbrace{m(t)}_{mass} \frac{\partial^2 \theta}{\partial t^2} + \underbrace{\gamma(t)}_{damping} \frac{\partial \theta}{\partial t} + \underbrace{\nabla_{\theta} L(\theta)}_{force} = 0$$
(4)

- Adding inertia, and stiffness
- Heavy-ball (Polyak 1946)

B. T. Polyak, Some methods of speeding up the convergence of iteration methods, USSR Comput. Math. Math. Phys., 4 (1964), pp. 1–17

Nesterov acceleration (due to Su et al. 2016)

W. Su, S. Boyd, and E. J. Candes, A differential equation for modeling Nesterov's accelerated gradient method: Theory and insights, J. Mach. Learn. Res., 17 (2016), pp. 1–43.

How to enforce constraints on the optimization? (Solver + velocity shocks)



Introduction. [3/21]



#### Inertial optimization methods





A. Sarshar CAIMS/SCMAI 2022 Introduction. [4/21]



#### Inertial optimization methods





A. Sarshar CAIMS/SCMAI 2022 Introduction. [4/21] 🕐 csl.cs.vt.edu



# Momentum methods through Hamiltonian mechanics part I

$$\dot{\mathbf{x}}_t = \nabla_{\mathbf{z}} \mathcal{H}(\mathbf{x}, \mathbf{z}, t),$$

$$\dot{\mathbf{z}}_t = -\nabla_{\mathbf{x}} \mathcal{H}(\mathbf{x}, \mathbf{z}, t)$$
(5)
(6)

- Describe the dynamics in terms of pairs  $(\mathbf{x}, \mathbf{z})$  (position, momentum) and Hamiltonian function (Total Energy) such that  $\frac{\partial \mathcal{H}}{\partial t} = 0$ .
- Additive splitting into potential and kinetic components :  $\mathcal{H}(\mathbf{x}, \mathbf{z}) = f(\mathbf{x}) + \psi^*(\mathbf{z})$
- Norm of the averaged gradient  $||\nabla f||$  decays with  $\mathcal{O}(1/t)$
- Unconstrained, No convergence guarantees due to preservation of energy.
   Diakonikolas, Jelena, and Michael I. Jordan. 2021. "Generalized Momentum-Based Methods: A Hamiltonian Perspective." SIAM Journal on Optimization 31 (1): 915–44.



A. Sarshar CAIMS/SCMAI 2022 

# Momentum methods through Hamiltonian mechanics part I

$$\begin{aligned} \dot{\mathbf{x}}_{t} &= \nabla_{\mathbf{z}_{t}} \mathcal{H} \left( \mathbf{x}_{t}, \mathbf{z}_{t} \right) = \nabla \psi^{*} \left( \mathbf{z}_{t} \right) \\ \dot{\mathbf{z}}_{t} &= -\nabla_{\mathbf{x}_{t}} \mathcal{H} \left( \mathbf{x}_{t}, \mathbf{z}_{t} \right) = -\nabla f \left( \mathbf{x}_{t} \right) \end{aligned}$$
(5)

- Describe the dynamics in terms of pairs  $(\mathbf{x}, \mathbf{z})$  (position, momentum) and Hamiltonian function (Total Energy) such that  $\frac{\partial \mathcal{H}}{\partial t} = 0$ .
- Additive splitting into potential and kinetic components :  $\mathcal{H}(\mathbf{x}, \mathbf{z}) = f(\mathbf{x}) + \psi^*(\mathbf{z})$
- Norm of the averaged gradient ||
  abla f|| decays with  $\mathcal{O}(1/t)$
- Unconstrained, No convergence guarantees due to preservation of energy.
   Diakonikolas, Jelena, and Michael I. Jordan. 2021. "Generalized Momentum-Based Methods: A Hamiltonian Perspective." SIAM Journal on Optimization 31 (1): 915–44.



A. Sarshar CAIMS/SCMAI 2022 Introduction. [5/21]



# Momentum methods Hamiltonian mechanics part II

$$\mathcal{H}_{M}(\bar{\mathbf{x}}, \mathbf{z}, \tau) = h(\tau)f(\bar{\mathbf{x}}/\tau) + \psi^{*}(\mathbf{z}), \quad \bar{\mathbf{x}} = \tau \mathbf{x},$$
  
$$\dot{\mathbf{x}}_{t} = \frac{\dot{\alpha}_{t} \left(\nabla\psi^{*}\left(\mathbf{z}_{t}\right) - \mathbf{x}_{t}\right)}{\alpha_{t}},$$
  
$$\dot{\mathbf{z}}_{t} = -h\left(\alpha_{t}\right)\frac{\dot{\alpha}_{t}}{\alpha_{t}}\nabla f\left(\mathbf{x}_{t}\right).$$
  
(6)

- $h(\cdot)$  and  $\alpha_t$  are strictly increasing functions of t
- Covers most of the acceleration methods
- Hamiltonian system is useful in deriving convergence properties both for f and  $\mathbf{x} = \theta$ .
- Does not cover CS-favored methods such as Adam, RMSProp, etc.



Introduction. [6/21]



# Adaptive momentum optimization methods (Adam)

Discrete (original) form of Adam (Kingma and Ba 2014) :

$$\begin{cases} g_{k} = \nabla f(\theta_{k-1}) \\ m_{k} = \mu_{k} m_{k-1} + (1 - \mu_{k}) g_{k} \\ v_{k} = \nu_{k} v_{k-1} + (1 - \nu_{k}) g_{k}^{2} \\ \theta_{k} = \theta_{k-1} - sm_{k} / (\sqrt{v_{k}} + \varepsilon). \end{cases}$$
(7)

- m(t) is a decaying moving average of estimates of the first moment (the mean) of the gradient.
- v(t) is a decaying moving average of estimates of the second raw moment (the uncentered variance) of the gradient.



A. Sarshar CAIMS/SCMAI 2022 Introduction. [7/21]



### Adaptive momentum optimization methods (Adam)

Continuous form ( Belotto da Silva and Gazeau 2020 ) :

$$\begin{cases} \dot{\theta} = -m/\sqrt{v+\varepsilon} \\ \dot{m} = h(t,\lambda,\alpha_1) \left(\nabla f(\theta) - m\right) \\ \dot{v} = h(t,\lambda,\alpha_2) \left(\nabla f(\theta)^2 - v\right) \end{cases}$$
(7)

We recognize it as a partitioned ODE system:

$$\begin{bmatrix} \theta \\ V \end{bmatrix}' = \begin{bmatrix} 0 \\ F(\theta, V) \end{bmatrix} + \begin{bmatrix} G(V) \\ 0 \end{bmatrix}$$

with  $V = [m, v]^T$ 



A. Sarshar CAIMS/SCMAI 2022 Introduction. [7/21]



(8)

# Generalized-structure Additive Runge-Kutta methods I

Additively-partitioned ODE system:

$$y' = \sum_{m=1}^{N} f^{\{m\}}(y),$$

• One step of a GARK scheme (Sandu & Günther, 2013):

$$Y_{i}^{\{q\}} = y_{n} + h \sum_{m=1}^{N} \sum_{j=1}^{s^{\{m\}}} a_{i,j}^{\{q,m\}} f^{\{m\}} \left(Y_{j}^{\{m\}}\right),$$
  
$$i = 1, \dots, s^{\{q\}}, \quad q = 1, \dots, N,$$
  
$$y_{n+1} = y_{n} + h \sum_{q=1}^{N} \sum_{i=1}^{s^{\{q\}}} b_{i}^{\{q\}} f^{\{q\}} \left(Y_{i}^{\{q\}}\right).$$



A. Sarshar CAIMS/SCMAI 2022 GARK schemes. [8/21]



$\mathbf{A}^{\{1,1\}}$	$\mathbf{A}^{\{1,2\}}$	 $\mathbf{A}^{\{1,N\}}$
$\mathbf{A}^{\{2,1\}}$	$A^{\{2,2\}}$	 $\mathbf{A}^{\{2,N\}}$
÷	:	: .
<b>A</b> <sup>{<i>N</i>,1}</sup>	<b>A</b> <sup>{<i>N</i>,2}</sup>	 <b>A</b> { <i>N</i> , <i>N</i> }
$\mathbf{b}^{\{1\}}$	<b>b</b> <sup>{2}</sup>	 <b>b</b> { <i>N</i> }

- Each component function has a different argument.
- One set of RK coefficients for each component and stage (flexibility).
- The GARK Butcher tableau:



A. Sarshar CAIMS/SCMAI 2022 GARK schemes. [9/21]



Individual RKs plus coupling terms.

$$Y_{i}^{\{1\}} = y_{n} + h \sum_{j=1}^{s^{\{1\}}} a_{i,j}^{\{1,1\}} f^{\{1\}} \left(Y_{j}^{\{1\}}\right) + h \sum_{j=1}^{s^{\{2\}}} a_{i,j}^{\{1,2\}} f^{\{2\}} \left(Y_{j}^{\{2\}}\right),$$
  

$$Y_{i}^{\{2\}} = y_{n} + h \sum_{j=1}^{s^{\{2\}}} a_{i,j}^{\{2,2\}} f^{\{2\}} \left(Y_{j}^{\{2\}}\right) + h \sum_{j=1}^{s^{\{1\}}} a_{i,j}^{\{2,1\}} f^{\{1\}} \left(Y_{j}^{\{1\}}\right),$$
  

$$y_{n+1} = y_{n} + h \sum_{i=1}^{s^{\{1\}}} b_{i}^{\{1\}} f^{\{1\}} \left(Y_{i}^{\{1\}}\right) + h \sum_{i=1}^{s^{\{2\}}} b_{i}^{\{2\}} f^{\{2\}} \left(Y_{i}^{\{2\}}\right).$$



A. Sarshar CAIMS/SCMAI 2022 GARK schemes. [10/21]



# $\mathsf{GARK}$ methods for the momentum equation

$$\begin{bmatrix} \theta \\ V \end{bmatrix}' = \underbrace{\begin{bmatrix} 0 \\ F(\theta, V) \end{bmatrix}}_{f^{\{1\}}} + \underbrace{\begin{bmatrix} G(V) \\ 0 \end{bmatrix}}_{f^{\{2\}}},$$
(9)

$$Y_{i} = y_{n} + h \sum_{j=1}^{i} a_{i,j}^{\{1,1\}} F(Y_{j}) + h \sum_{j=1}^{i} a_{i,j}^{\{1,2\}} G(Z_{j}), \quad i = 1, \dots, s^{\{1\}},$$

$$Z_{i} = y_{n} + h \sum_{j=1}^{i} a_{i,j}^{\{2,1\}} F(Y_{j}) + h \sum_{j=1}^{s^{\{2\}}} a_{i,j}^{\{2,2\}} G(Z_{j}), \quad i = 1, \dots, s^{\{2\}},$$

$$y_{n+1} = y_{n} + h \sum_{i=1}^{s^{\{1\}}} b_{i}^{\{1\}} F(Y_{i}) + h \sum_{i=1}^{s^{\{2\}}} b_{i}^{\{2\}} G(Z_{i}),$$
(10)



A. Sarshar CAIMS/SCMAI 2022 GARK schemes. [11/21]



#### GARK methods for the momentum equation

$$\begin{bmatrix} \theta \\ V \end{bmatrix}' = \underbrace{\begin{bmatrix} 0 \\ F(\theta, V) \end{bmatrix}}_{f^{\{1\}}} + \underbrace{\begin{bmatrix} G(V) \\ 0 \end{bmatrix}}_{f^{\{2\}}},$$
(11)

$$Y_{i} = \begin{bmatrix} W^{Y} \\ V^{Y} \end{bmatrix}_{i} = y_{n} + h \sum_{j=1}^{i} a_{i,j}^{\{1,1\}} F(Y_{j}) + h \sum_{j=1}^{i} a_{i,j}^{\{1,2\}} G(Z_{j}), \quad i = 1, \dots, s^{\{1\}},$$

$$Z_{i} = \begin{bmatrix} W^{Z} \\ V^{Z} \end{bmatrix}_{i} = y_{n} + h \sum_{j=1}^{i} a_{i,j}^{\{2,1\}} F(Y_{j}) + h \sum_{j=1}^{s^{\{2\}}} a_{i,j}^{\{2,2\}} G(Z_{j}), \quad i = 1, \dots, s^{\{2\}}, \quad (12)$$

$$y_{n+1} = y_{n} + h \sum_{i=1}^{s^{\{1\}}} b_{i}^{\{1\}} F(Y_{i}) + h \sum_{i=1}^{s^{\{2\}}} b_{i}^{\{2\}} G(Z_{i}),$$



A. Sarshar CAIMS/SCMAI 2022 GARK schemes. [12/21]



$$W_{i}^{Y} = W_{n} + h \sum_{j=1}^{i} a_{i,j}^{\{1,2\}} g(V_{j}^{Z}), \quad i = 1, \dots, s^{\{1\}},$$

$$V_{i}^{Y} = V_{n} + h \sum_{j=1}^{i} a_{i,j}^{\{1,1\}} f\left([W_{j}^{Y}, V_{j}^{Y}]^{T}\right), \quad i = 1, \dots, s^{\{1\}},$$

$$V_{i}^{Z} = V_{n} + h \sum_{j=1}^{i} a_{i,j}^{\{2,1\}} f\left([W_{j}^{Y}, V_{j}^{Y}]^{T}\right), \quad i = 1, \dots, s^{\{2\}},$$

$$W_{n+1} = W_{n} + h \sum_{i=1}^{s^{\{2\}}} b_{i}^{\{2\}} g(V_{i}^{Z}),$$

$$V_{n+1} = W_{n} + h \sum_{i=1}^{s^{\{1\}}} b_{i}^{\{1\}} f\left([W_{i}^{Y}, V_{i}^{Y}]^{T}\right).$$
(13)

note that the term  $W^Z$  is not used in any calculation, and can thus be dropped. Notice that the term  $W_i^Y$  is quasi-implicit.



A. Sarshar CAIMS/SCMAI 2022 GARK schemes. [13/21]



(1 1)   - (1 1)	• (1.2)		$0 \frac{1}{2}$	0 $\frac{1}{2}$	0 0	0 0	$0 \frac{1}{2}$	0 0	
$\frac{\mathbf{c}^{\{1,1\}}}{\mathbf{c}^{\{2,1\}}} \frac{\mathbf{A}^{\{1,1\}}}{\mathbf{A}^{\{2,1\}}}$	$A^{\{1,2\}}$ $A^{\{2,2\}}$	:=	1	$\frac{1-\beta}{1-\beta}$	β	0	$\frac{\frac{1}{2}}{1}$	$\frac{1}{2}$	
$\mathbf{b}^{\{1\}}$	<b>b</b> <sup>{2}</sup>	_	1 4 3 4	$\frac{1}{4}$ $\frac{1}{4}$	0 $\frac{1}{2}$	0 0	$\frac{1}{4}$ $\frac{1}{2}$	$\frac{1}{4}$	
				$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{2}$	-



A. Sarshar CAIMS/SCMAI 2022 GARK schemes. [14/21]



# Vanilla Adam versus higher order explicit methods



$$X \xrightarrow{\theta} U \xrightarrow{\phi} \widetilde{X},$$
  
$$L(\widetilde{X}, X) = \left\| X - \widetilde{X} \right\|^2 + \alpha_1 \|\phi\|^2 + \alpha_2 \|\theta\|^2$$
$$\min_{\phi, \theta} L(\widetilde{X}, X) m$$





A. Sarshar CAIMS/SCMAI 2022 Numerical Experiments. [15/21]



### Vanilla Adam versus higher order explicit methods





A. Sarshar CAIMS/SCMAI 2022 Numerical Experiments. [16/21] csl.cs.vt.edu



# Vanilla Adam versus higher order explicit methods



Future directions: Partitioning the linear part

$$\begin{cases} \dot{\theta} = -m/\sqrt{v + \varepsilon} \\ \dot{m} = h(t, \lambda, \alpha_1) \left(\nabla f(\theta) - m\right) \\ \dot{v} = h(t, \lambda, \alpha_2) \left(\nabla f(\theta)^2 - v\right) \\ \left[ \begin{matrix} \theta \\ V \end{matrix} \right]' = \begin{bmatrix} 0 \\ f(W) \end{bmatrix} + \begin{bmatrix} g(V) \\ BV \end{bmatrix}$$

with  $V = [m, v]^T$ .



A. Sarshar CAIMS/SCMAI 2022 Future works. [18/21]



(14)

Future directions: Partitioning the linear part

$$\begin{bmatrix} \theta \\ V \end{bmatrix}' = \begin{bmatrix} 0 \\ f(W) \end{bmatrix} + \begin{bmatrix} g(V) \\ BV \end{bmatrix}$$
(15)  
$$W_i^{Y} = W_n + h \sum_{j=1} a_{i,j}^{\{1,2\}} g(V_j^{Z})$$
(16)  
$$V_i^{Z} = V_n + h \sum_{j=1} a_{i,j}^{\{2,1\}} f(W_j^{Y}) + h \sum_{j=1}^{i} a_{i,j}^{\{2,2\}} BV_j^{Z}$$
(17)  
$$W_{n+1} = W_n + h \sum_{i=1}^{s^{\{2\}}} b_i^{\{2\}} g(V_i^{Z}),$$
(18)  
$$V_{n+1} = W_n + h \sum_{i=1}^{s^{\{1\}}} b_i^{\{1\}} f(W_i^{Y}) + h \sum_{i=1}^{s^{\{2\}}} b_i^{\{2\}} BV_i^{Z}$$
(19)



A. Sarshar CAIMS/SCMAI 2022 Future works. [19/21]



# Summary and future directions

- ODE form of momentum methods for (stochastic) optimization
- Amenable to IMEX and ADI splittings
- Improved accuracy helps faster convergence to a suitable (local) minimum
- Indications that stability is a stronger requirement than order (Chebyshev explicit methods)
- Extensive testing is required for many machine learning tasks (Classification, Regression, Time-series forcasting, etc.)
- GARK framework allows offers new methods, what are the requirements for a successful one? (SDE analysis)
- arxiv.org/a/sarshar\_a\_1
- csl.cs.vt.edu



A. Sarshar CAIMS/SCMAI 2022 Future works. [20/21]





