

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

#### Objective-function-free optimization

#### Serge Gratton, Sadok Jerad, Alena Kopaničáková and Philippe Toint

INP - ANITI / UNamur

Louvain-La-Neuve, August 2024

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

# First: a brief publicity break :-)





▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

### The problem

Once more, the standard unconstrained nonconvex optimization problem

$$\min_{x\in \mathbf{R}^n} f(x)$$

where the objective function f is

- "sufficiently" smooth
- bounded below

Remarkable one can still say (hopefully) interesting things on this subject!

In this talk: focus on the deterministic case, but ...

# Why OFFO?



Our target: robust algorithms for noisy functions/inexact arithmetic

For convergence, standard methods (TR, AR) requires an error on function values which is the square (!) of that on the gradient (e.g. Bellavia et al, 22)



⇒ Design algorithms that do not evaluate the function

#### Adaptive gradient methods:

- Adagrad (Duchi et al, 2011)
- WNGrad (Wu, Ward, Bottou, 2018)
- Adam (Kingma, Ba, 2014)
- A trust-region method:
- Adatr (Grapiglia, 2022)

 $\Rightarrow$  Objective Function Free Optimization = OFFO

# ASTR1 an adaptive trust-region algorithm

Step 0: Initialization.  $x_0$  is given. Set k = 0. Step 1: Define the TR. Compute  $g_k = g(x_k)$  and define  $\Delta_{i,k} = \frac{|g_{i,k}|}{w_{i,k}}$ where  $w_{i,k} \ge \varsigma_i > 0$  are weights. Step 2: Hessian approximation. Select a symmetric  $B_k$ . Step 3: GCP. Define  $s_{i,k}^{L} = -\text{sgn}(g_{i,k})\Delta_{i,k}$  and  $s_{k}^{Q} = \gamma_{k}s_{k}^{L}$ with  $\gamma_k = \begin{cases} \min\left[1, \frac{|g_k' \, s_k^L|}{(s_k^L)^T B_k s_k^L}\right] & \text{if } (s_k^L)^T B_k s_k^L > 0, \\ 1 & \text{otherwise.} \end{cases}$ Step 3: Step. Compute a step  $s_k$  such that  $|s_{i,k}| \leq \Delta_{i,k}$  ( $\forall i$ ) and  $g_{l}^{T} s_{l} + \frac{1}{2} s_{l}^{T} B_{l} s_{l} < g_{l}^{T} s_{l}^{Q} + \frac{1}{2} (s_{l}^{Q})^{T} B_{l} s_{l}^{Q}$ Step 5: New iterate. Set  $x_{k+1} = x_k + s_k$ , increment k, and go to Step 1.

# ASTR1: comments

- ► the objective function is not evaluated ⇒ OFFO ... and thus the TR radius cannot depend on ared/prered.
- ► large weights ⇒ short steps
- $\triangleright$   $\gamma_k$  minimize the quadratic model between 0 and  $s_k^L$

Suppose that  $f \in C^1$ , has Lipschitz gradient with constant L and that  $||B_k|| \le \kappa_B$ . Then  $f(x_{k+1}) \le f(x_k) - \sum_{i=1}^n \frac{\operatorname{Smin} g_{i,j}^2}{2\kappa_B w_{i,j}} + \frac{1}{2}(\kappa_B + L) \sum_{i=1}^n \frac{g_{i,j}^2}{w_{i,j}^2}$ 

 $\Rightarrow$  descent for large enough weights  $w_{i,k}$ 



### ASTR1 with ADAGRAD-like weights (1) For given $\varsigma \in (0, 1]$ , $\vartheta \in (0, 1]$ and $\mu \in (0, 1)$ , define

$$w_{i,k} \in \left[\vartheta\left(\varsigma + \sum_{\ell=0}^{k} g_{i,\ell}^{2}\right)^{\mu}, \left(\varsigma + \sum_{\ell=0}^{k} g_{i,\ell}^{2}\right)^{\mu}\right]$$

For 
$$\vartheta = 1$$
 and  $\mu = rac{1}{2}$ ,  $w_{i,k} = \sqrt{\varsigma + \sum_{\ell=0}^{k} g_{i,\ell}^2}$  and

ASTR1 with 
$$\vartheta = 1$$
,  $\mu = \frac{1}{2}$  and  $B_k = 0$  is ADAGRAD

Suppose that  $f \in C^1$ , has Lipschitz gradient with constant L and is bounded below. Then ASTR1 with ADAGRAD-like weights,  $\mu \in (0, 1]$  and  $||B_k||$  uniformly bounded requires at most

iterations to produce an iterate k such that 
$$\operatorname{average}_{0,\ldots,k} \|g_{\ell}\|^2 \leq \epsilon$$
.

(日)

# More on ASTR1

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ● □ ● ● ●

- Extends known results (e.g., Wu, Ward, Bottou, 2018)
- Allows the use of curvature information in an ADAGRAD-like method (Barzilai-Borwein, LBFGS, quasi-Newton, ... true Hessian)

The above bound is essentially sharp.

Also possible with the "divergent" weights

$$w_{i,k} \in [v_{i,k}(k+1)^
u, v_{i,k}(k+1)^\mu]$$

for  $0 < \nu \leq \mu < 1$  and

$$v_{i,k} = \max_{0,\dots,k} |g_{i,\ell}|$$
 or  $v_{i,k} = \operatorname{average}_{0,\dots,k} |g_{i,\ell}|$ 

Slightly weaker (sharp) complexity result



## Some results on the small noiseless OPM problems

Method	$\pi_{\texttt{algo}}$	$ ho_{\texttt{algo}}$
adagbfgs3	0.75	69.75
sdba (using $f$ )	0.73	68.91
adagH	0.72	69.75
adagrad	0.69	73.11
maxg	0.66	66.39
adagbb	0.63	64.71
adam	0.54	30.25

Performance and reliability statistics for deterministic OFFO and steepest descent algorithms on small OPM problems ( $\epsilon = 10^{-6}$ )



#### The impact of noise

	$\rho_{\tt algo}/{\rm relative}$ noise level				
algo	0%	5%	15%	25%	50%
adagH	83.19	84.96	84.20	84.71	82.18
adagbfgs3	78.15	80.50	80.50	80.84	80.18
adagrad	77.31	80.50	80.25	80.17	80.17
adagbb	75.69	80.08	80.17	79.58	79.41
maxg	74.79	74.37	75.55	78.15	78.07
adam	40.34	35.55	36.30	44.03	45.80
sdba	81.51	30.92	31.85	34.87	29.58

Reliability of OFFO algorithms and steepest descent as a function of the level of relative Gaussian noise ( $\epsilon=10^{-3})$ 

OFFO and multilevel optimization: context

Statistical machine learning for solving PDE's

Key ingredients:

- an approximation set, typically a nonlinear Neural Network (NN) architecture
- a sampling technique of the space and time domains
- ► "training" = minimization of a loss encoding the PDE (or the underlying physics) ⇒ Physically Informed Neural Networks (PINNs) See [Raissi, Perdikaris, Karniadakis, 2017]

Single level convergence theory exists (e.g. [Shin, Darbon, Karniadakis, 2020]) and involves

- the universality property of NN,
- statistical sampling,
- ability of numerical optimizers (ADAM, SGD,...) to reach an approximate global optimum of nonconvex function

OFFO and multilevel optimization: idea(s) (1)

New algorithm:

- ► assume a hierarchy of (smooth) models h<sub>ℓ</sub>(x) (from fine to coarse)
- ensure first-order model coherence between levels by adding a linear term if necessary (cfr FAS))
- ► an OFFO trust-region based algorithm at each level *l*: (approx) minimize a quadratic model

$$g_{\ell,k}^T s + \frac{1}{2} s^T H_{\ell,k} s$$

in  $\mathcal{B}_{\ell,k} = \{s \text{ at level } \ell \mid ||s|| \leq \Delta_{\ell,k}\}$ 

• (no evaluation of  $f(x) \Rightarrow \text{accept all iterates}$ )

OFFO and multilevel optimization: idea(s) (2)



- either using a Taylor approximation of the model at level l
  (standard OFFO step), or
- ► use the OFFO algorithm to minimize the lower level model h<sub>ℓ-1</sub> with the TR at level ℓ (recursivity)
- only use lower level if lower-level gradient is not too small wrt upper-level one
- when using lower level, ensure the trust-region radius has a suitable size (a bit technical)

# OFFO and multilevel optimization: algorithm and results

Note:

- Detailed algorithm does not fit on slide, but...
- Subsumes most existing (first-order) OFFO methods in the single level case
- allows use of second-order information (LBFGS,... or even true Hessian)

Result:

• complexity is  $\mathcal{O}(\epsilon^{-2})$  (optimal for single level case)

 first-order version works well in experiments (stochastic context, momentum, resnets, ..., see (Gratton, Kopaničáková, T., 2023))

More can be said about PINN's (Gratton, Mercier, Riccetti, T., 2024)

## An example





|▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ | 画|| のへの

#### ANITI

### Towards second-order criticality

Use a trust-region mechanism for second-order criticality?

At  $x_k$ , let

$$T_{f,2}(x_k,d) = f(x_k) + g(x)_k^T d + \frac{1}{2} d^T H(x_k) d.$$

and the second-order criticality measure

$$\phi_{f,2}^{\delta}(x_k) = \max_{\|d\| \leq \delta} - \left(g(x_k)^T d + \frac{1}{2}d^T H(x_k)d\right) = \max_{\|d\| \leq \delta} \Delta q_k(d)$$

Define:

 $x_k$  is  $\epsilon$ -second-order critical if  $\phi_{f,2}^{\delta}(x_k) \leq \epsilon$ 

Idea: Use  $\phi_{f,2}^{\delta}(x_k)$  to define weights for the trust-region



(日) (日) (日) (日) (日) (日) (日) (日)

### Function decrease for ASTR2

Suppose that  $f \in C^2$  and has Lipschitz continuous gradient and Hessian. Then, if  $||g_k||^2 \ge \widehat{\phi}_k^3$ ,  $f_{k+1} \le f_k - \frac{||g_k||^2}{w_k^L} + \frac{L_1}{2} \frac{||g_k||^2}{(w_k^L)^2}$ while, if  $||g_k||^2 < \widehat{\phi}_k^3$ ,  $f_{k+1} \le f_k - \kappa \min\left[\frac{1}{2(1+L_1)}, \frac{1}{w_k^Q}, \frac{1}{(w_k^Q)^2}\right] \widehat{\phi}_k^3 + \frac{L_2}{6} \frac{\widehat{\phi}_k^3}{(w_k^Q)^3}$ .

 $\Rightarrow$  roles of  $w_k^L$  and  $w_k^Q$  complementary

# Complexity of ASTR2 for ADAGRAD-like weights

When using

$$\begin{split} w_{k}^{L} &\in \left[\vartheta \left(\varsigma + \sum_{\ell=0, \ell \in \mathcal{K}^{L}}^{k} \|g_{\ell}\|^{2}\right)^{\mu}, \left(\varsigma + \sum_{\ell=0, \ell \in \mathcal{K}^{L}}^{k} \|g_{\ell}\|^{2}\right)^{\mu}\right] \\ w_{k}^{Q} &\in \left[\vartheta \left(\varsigma + \sum_{\ell=0, \ell \in \mathcal{K}^{Q}}^{k} \widehat{\phi}_{k}^{3}\right)^{\mu}, \left(\varsigma + \sum_{\ell=0, \ell \in \mathcal{K}^{Q}}^{k} \widehat{\phi}_{k}^{3}\right)^{\mu}\right] \end{split}$$

Suppose that  $f \in C^2$  with Lipschitz gradient and Hessian and is bounded below. Then ASTR2 with the above weights and  $\mu \in (0, 1]$ requires at most  $\mathcal{O}(\epsilon^{-1})$  iterations to produce an iterate k such that  $\operatorname{average}_{0,\ldots,k} \|g_{\ell}\|^2 \leq \epsilon$  and  $\operatorname{average}_{0,\ldots,k} \widehat{\phi}_{\ell}^3 \leq \epsilon$ . [Essentially sharp!] ... and now for an OFFO regularization algorithm!

Consider now the more general

$$T_{f,p}(x,s)=f(x)+\sum_{i=1}^p\frac{1}{i!}\nabla^i_xf(x)[s]^i.$$

and the derived regularized model

$$m_k(s) = T_{f,p}(x_k,s) + \frac{\sigma_k}{(p+1)!} ||s||^{p+1}$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

We assume that  $\nabla_x^p f$  is globally Lipschitz.

# The OFFAR algorithm



(again using generic  $\kappa$ )

Step 0: Initialization:  $x_0$ ,  $\nu_0 > 0$ ,  $\epsilon$  and constants. Set k = 0. Step 1: Check for termination: Evaluate  $g_k = \nabla_x^1 f(x_k)$  and terminate if  $||g_k|| \le \epsilon$ . Else, evaluate  $\{\nabla_x^i f(x_k)\}_{i=2}^p$ . Step 2: Step calculation: If k = 0, set  $\sigma_0 = \mu_0 = \nu_0$ . Else set  $\mu_k = \frac{p! ||g_k||}{||s_{k-1}||^p} - \kappa \sigma_{k-1}$  and  $\sigma_k \in [\kappa \nu_k, \max(\nu_k, \mu_k)]$ .

Then compute a step  $s_k$  such that

$$m_k(s_k) < m_k(0)$$
 and  $\|\nabla^1_s \mathcal{T}_{f,p}(x_k,s_k)\| \leq \kappa \frac{\sigma_k}{p!} \|s_k\|^p$ .

Step 3: Updates. Set  $x_{k+1} = x_k + s_k$  and  $\nu_{k+1} = \nu_k + \nu_k ||s_k||^{p+1}$ . Increment k by one and go to Step 1.



# Complexity of OFFAR

- No objective function evaluation  $\Rightarrow$  OFFO
- The use of μ<sub>k</sub> is optional: one could simply set μ<sub>k</sub> = 0 without altering the theory. But it is important for performance.
- The definition of µ<sub>k</sub> promotes fast growth of the regularization parameter up the problem's Lispchitz constant
- The definition of σ<sub>k</sub> helps to limit this growth once the value of the Lipschitz constant has been reached.

▶ If 
$$p = 1$$
,  $\nu_{k+1} = \nu_k + \nu_k ||s_k||^2$ , recovering WNGrad (Wu, Ward,  
Bottou, 2018)

Suppose that  $f \in C^p$  with  $\nabla_x^p f$  Lipschitz gradient, is bounded below and is such that  $\min_{\|d\| \le 1} \nabla_x^i [d]^i \ge \kappa$  for i = 2, ..., p. Then OFFAR (with suitable constants) requires at most  $\mathcal{O}\left(\epsilon^{-\frac{p+1}{p}}\right)$  iterations to produce an iterate k such that  $\|g_k\| \le \epsilon$ .

### More on OFFAR

- Same rate as ARp using function values (Birgin et al, 2016)
- For p = 2, same rate as ARC/AR2 (Cartis, Gould, T. 2011). Optimal rate for second order methods

Optimal rates for exact *p*th order methods (Carmon et al. 2019).
 MOFFAR: If one requires that the step also satisfies

$$\max\left(0,-\lambda_{\min}[\nabla_{s}^{2}\mathcal{T}_{f,p}(x_{k},s_{k})]\right) \leq \frac{\kappa\sigma_{k}}{(p-1)!}\|s_{k}\|^{p-1}$$

Suppose that  $f \in C^p$  with  $\nabla_x^p f$  Lipschitz gradient, is bounded below and is such that  $\min_{\|d\| \le 1} \nabla_x^i [d]^i \ge \kappa$  for i = 2, ..., p. Then MOF-FAR (with suitable constants) requires at most  $\mathcal{O}\left(\epsilon^{-\frac{p+1}{p-1}}\right)$  iterations to produce an iterate k such that  $\|g_k\| \le \epsilon$  and  $\widehat{\phi}_k \le \epsilon$ .



#### Numerical illustration

For AR2 and two variants of OFFAR with p = 2, differing on how aggressively  $\mu_k$  forces growth in  $\sigma_k$  (b more aggressive than a)

	AR2	OFFAR2a	OFFAR2b
$\pi_{\texttt{algo}}$	0.99	0.78	0.83
$\rho_{\texttt{algo}}$	97.48	81.51	88.24

Performance and reliability statistics on the small OPM problems without noise

	5%	15%	25%	50%
AR2	40.67	30.84	24.54	6.81
OFFAR2a	80.76	75.38	70.76	56.30
OFFAR2b	85.97	80.67	72.69	47.98

Reliability statistics  $\rho_{algo}$  for 5%, 15%, 25% and 50% relative random Gaussian noise (averaged on 10 runs)



### Stochastic variants

Complexity bounds for first-order criticality (in expectation):

Algorithm	type	Compl.bound
ASTR1	trust-region	$\mathcal{O}\left(\epsilon^{-\frac{1}{2}(1-\mu)}\right)$
ASTR2	trust-region	??
STOFFAR	adat. regularization	$\mathcal{O}(\epsilon^{-3/2})$

STOFFAR = OFFAR +

$$\mathbb{E}_{k}\left[\left\|\nabla_{x}^{i}f(X_{k})-\overline{\nabla_{x}^{i}f}(X_{k})\right\|^{\frac{p+1}{p+1-i}}\right] \leq \kappa_{D}\sum_{i=1}^{m}\left\|S_{k-i}\right\|^{p+1} \quad (i=1,2)$$

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @

# ASTR1 on CIFAR-10 with cifar10-nv ( $\gamma = 10^{-5}$ )



▲ロト ▲樹 ト ▲ 臣 ト ▲ 臣 ト 一臣 - の Q ()~

# ASTR1 on CIFAR-10 with resnet18 ( $\gamma = 5.10^{-5}$ )



▲ロト ▲樹 ト ▲ 臣 ト ▲ 臣 ト 一臣 - の Q ()~

# **STOFFAR**



Figure: Loss function and number of samples for SUSY and w8a

#### Conclusions

Computing the value of f is not necessary for (theoretical) fast convergence

The use of curvature information is possible (and often beneficial) in standard OFFO adaptive methods

OFFO creates some interesting challenges in convergence theory!

Extension of ASTR1 to problems with convex constraints available!

Complexity of stochastic variants (ASTR1, OFFAR) also analyzed

Thank you for your interest and patience ... and good wind to Yurii!

#### Details in...

S. Gratton and S. Jerad and Ph. L. Toint, "Parametric Complexity Analysis for a Class of First-Order Adagrad-like Algorithms", to appear in Optimization Methods and Software, 2024.

S. Gratton and Ph. L. Toint, "OFFO minimization algorithms for second-order optimality and their complexity", Computational Optimization and Applications, vol. 84, pp. 573—607, 2022.

S. Gratton and S. Jerad and Ph. L. Toint, "Convergence properties of an Objective-Function-Free Optimization regularization algorithm, including an  $\mathcal{O}(\epsilon^{-3/2})$  complexity bound", SIAM Journal on Optimization, vol. 33(3), pp. 1621–1646, 2023.

S. Gratton, A. Kopaničáková and Ph. L. Toint, "Multilevel Objective-Function-Free Optimization with an Application to Neural Networks Training", SIAM Journal on Optimization, vol. 33(4), pp. 2772–2800, 2023.

S. Gratton and S. Jerad and Ph. L. Toint, "Complexity of Adagrad and other first-order methods for nonconvex optimization problems with bounds and convex constraints", arXiv:2406.15793,2024.

S. Gratton and S. Jerad and Ph. L. Toint, "A Stochastic Objective-Function-Free Adaptive Regularization Method with Optimal Complexity", arXiv:2407.08018, 2024.

S. Gratton, V. Mercier, E. Riccietti and Ph. L. Toint, "A Block-Coordinate Approach of Multi-level Optimization with an Application to Physics-Informed Neural Networks", Computational Optimization and Applications, to appear, 2024.