

# On the complexity of nonsmooth automatic differentiation

Jérôme Bolte , Ryan Boustany , Edouard Pauwels and Béatrice Pesquet-Popescu

ICLR 2023 - The Eleventh International Conference on Learning Representations



THALES



# Automatic Differentiation in Deep Learning

$$\min_{\theta \in \mathbb{R}^P} J(\theta) := \frac{1}{N} \sum_{i=1}^N \ell(f(x_i, \theta), y_i)$$

- $(x_i, y_i)_{i=1}^N$ : training set
- $f$ : can be composed of nonsmooth functions (e.g., ReLU, MaxPooling)
- $\theta \in \mathbb{R}^P$ : weight parameters
- $\ell$ : loss function

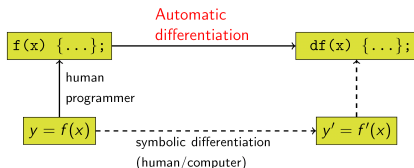


Figure: Illustration of how AD works

Two assumptions for training nonsmooth DNNs:

- ① Backpropagation outputs a gradient almost everywhere (theorem)
- ② The process is fast (empirical observation)

$f : \mathbb{R}^p \rightarrow \mathbb{R}$  differentiable function.

$P$ : program computing  $f$ .

$\text{backprop}(P)$ : program computing  $(f, \nabla f)$  using  $\text{backprop}$  AD.

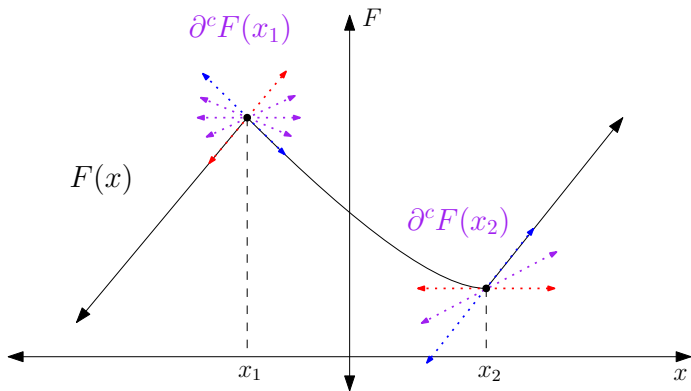
$\text{cost}(\cdot)$ : execution time to evaluate a program.

Theorem (Baur and Strassen, 1983)

For  $\text{rational}$  functions  $f$ :  $\text{cost}(\text{backprop}(P)) \leq 5 \times \text{cost}(P)$

**Motivation: generalize to nonsmooth functions.**

# Locally Lipschitz Functions and the Clarke Subgradient



$$\partial^c F(x) = \text{conv} \left\{ \lim_{k \rightarrow +\infty} \nabla F(x_k) : x_k \in \text{diff}_F, x_k \xrightarrow[k \rightarrow +\infty]{} x \right\}$$

# How does nonsmooth AD algorithm works ?

## Nonsmooth AD

$F : \mathbb{R}^p \rightarrow \mathbb{R}$  locally Lipschitz function in compositional form

$$F = g_1 \circ \dots \circ g_m.$$

- $d_i(x) = \nabla g_i(x)$  for **smooth**  $g_i$
- $d_i(x) \in \partial^c g_i(x)$  (when you hit a **nonsmooth** part)
- Ex :  $g_i = \text{relu}$  and take  $d_i(0) = \text{relu}'(0) = 0$  (Tensorflow, Pytorch)
- $\text{backprop}(P)$ : **chain rule** the  $d_i$ 's.

## Artifacts

- $\text{backprop } g_1(W) + \dots + \text{backprop } g_m(W) \notin \partial^c(g_1 + \dots + g_m)(W)$
- $\text{relu}_2 : t \rightarrow \text{relu}(-t) + t$  and  $\text{relu}'(0) = 0$
- $\text{zero} = \text{relu}_2 - \text{relu}$  and  $\text{zero}'(0) = 1$  (!)

# Conservative gradients (Bolte and Pauwels 2019)

## Main properties of conservative gradients

Let  $D_f : \mathbb{R}^p \rightrightarrows \mathbb{R}^p$  be a **conservative gradient** for  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  locally Lipschitz.

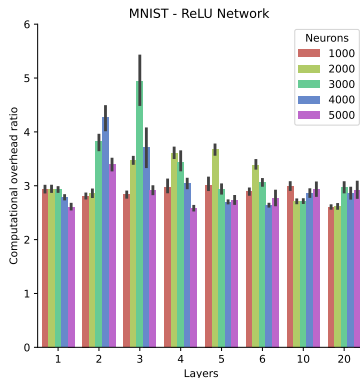
- Conservative gradients = gradients a.e
  - For all  $x \in \mathbb{R}^p$ ,  $\partial^c F(x) \subset \text{conv}(D_F(x))$ .
  - Compatible with **calculus rules** contrary to Clarke subdifferential.
- 
- **Most common** (virtually all, semialgebraic) functions using in DL admits **conservative gradients**.
  - Faithfully **model** what is computed by **backprop** (generated by Pytorch).
  - Preserving **convergence guaranties**.

# Nonsmooth Cheap Gradient Principle

## Theorem (Nonsmooth Cheap Gradient Principle, (Bolte et al. 2022))

Let  $P$  be a **program** that computes  $F = g_1 \circ \dots \circ g_m$  locally Lipschitz.

- 1 **backprop**( $P$ ) returns an element of a conservative gradient.
- 2  $\text{cost}(\text{backprop}(P)) \leq \boxed{\omega_b} \times \text{cost}(P)$ , where  $\omega_b$  is a constant.
- 3  $\boxed{\omega_b \approx 5}$  for ReLU networks.



# Comparison of conservative gradients with other nonsmooth AD

## Alternative AD methods

For a locally Lipschitz function  $F : \mathbb{R}^p \mapsto \mathbb{R}$ :

- $\omega := \text{cost}(p \times p \text{ matrix multiplication}) \approx p^{2.7}$  (best algorithm)
- $\text{cost}(p \text{ directional derivatives of } F) / \text{cost}(F) \geq p^{\omega-2}$

## Computational complexity comparison

$F$  relu network with matrix and vector entries in  $\{-1,0,1\}$ .

- ① Computing two distinct elements of  $\partial^c F$  is NP-hard.
- ② Computing two elements of conservative gradients is polynomial time solvable.



- Generalized subgradients cannot capture backpropagation (conservatives do).
- We extend the Baur-Strassen theorem (rational functions) to semi-algebraic functions (ubiquitous in ML) with a specific arithmetic model.
- We prove worst case lower bound in the nonsmooth context for directional derivatives and subgradient enumeration.

<https://arxiv.org/abs/2206.01730>

Thanks for your attention.