

# Garanties pour l'autodiff.

Travail en commun avec  
J. Bolte & E. Pauwels.

pas un talk sur des considérations pratiques  
(contrainte de temps)

Focus: Pourquoi l'autodiff d'un algorithme itératif fonctionnel ?

1) Algorithme itératif paramétrique  $\approx$   $x_{k+1}(\theta) = F(x_k(\theta), \theta) = F_{\theta}(x_k(\theta))$

↑  
opérateurs de "point fixe"

ex.  $F(x, \theta) = x - \theta \nabla \phi(x)$ . descente de gradient.

•  $F(x, \theta) = \text{prox}_{\theta \psi}(x - \theta \nabla \phi(x))$  Forward Backward

condition typique :  
 $F$   $\rho$ -Lip  $\rho < 1$

Fix  $F_{\theta} = \bar{x}(\theta)$

$x_{k+1}(\theta)$   
↓  
 $k \rightarrow +\infty$   
↓  
 $\bar{x}(\theta)$

cas de GD  
 $\phi \in C_{\rho}^{1,1}(\mathbb{R}^d)$   
 $\bar{x} = \text{argmin } \phi$   
 $x_k \rightarrow \bar{x}$ .

## 2) Autodiff.

Si  $F$  lisse ( $C^1$ ), modèle simple  
de différentiation automatique d'algo

→ règle de la chaîne appliquée  
à  $\theta$ . "piggy back"

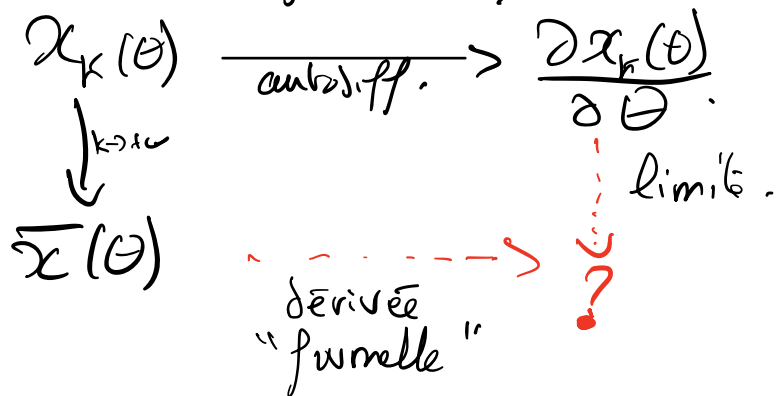
$$\frac{\partial x_k(\theta)}{\partial \theta} = A \frac{\partial x_k(\theta)}{\partial \theta} + B$$

$\uparrow$   $\uparrow$

$[\partial_1 F(x_k(\theta), \theta), \partial_2 F(x_k(\theta), \theta)]$

dérivée totale.

( $\approx$  backward de PyTorch)



→  $\bar{x}$  candidat :  $\frac{\partial \bar{x}(\theta)}{\partial \theta}$   
• naturel.

Thm (Gilbert 1992) informel.  
 Si  $\|A\| < 1$  alors  $\bar{x}$  est  $C^1$  et

$$\frac{\partial \bar{x}(\theta)}{\partial \theta} \stackrel{\text{linéarité}}{\underset{\text{K-S}}{\text{K-S}}} \frac{\partial \bar{\mathcal{L}}(\theta)}{\partial \theta}$$

Note 1  $\|A\| < 1$  est suffisant mais pas nécessaire (Devets, V. GORT '23).

Note 2  $\bar{\mathcal{L}}$  peut être également obtenu par différentiation implicite.

$$\frac{\partial \bar{x}}{\partial \theta}(\theta) = (\text{Id} - \partial_x F(\bar{x}(\theta), \theta))^{-1} \partial_\theta F(\bar{x}(\theta), \theta)$$

3) L'éléphant dans la salle :  
dynamique non-lisse.

ex A  $\min_x \phi(x)$        $F(x, \theta) = x - \theta \nabla \phi(x)$

pb d'optim

• Si  $\phi \in C_p^2(\mathbb{R}^d)$ .

→  $F$  est  $C^1$  😊  
 (Neuron?)

ex B  $F$  est un DEW  
 avec ReLU

• Si  $\phi \in C_p^{1,1}(\mathbb{R}^d)$  cas typique  
 →  $F$  pas  $C^1$  ☹️

Solution Jacobien conservatif.  
(Bolte, Pauvel, Math Prog 21)

Def Si:  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  continue loc. Lipschitz  
on dit que  $J: \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  est un

app. multivoque  
gradient conservatif pour la fonction  
 $f$  (dit chemin dérivable) si:

- $J$  a un graph fermé, loc. borné, non-vide partout.
- $J$  satisfait la règle de la chaîne  
sur les chemins:

$\forall \gamma: [0,1] \rightarrow \mathbb{R}^d$  a.c. wrt Lebesgue

$$\frac{d}{dt} f(\gamma(t)) = J(\gamma(t)) \cdot \dot{\gamma}(t)$$

Propriétés importantes:

- $\partial^c f(x) \subseteq J(x)$   
 $\uparrow$  Clarke.

- $J(x) = \left\{ \frac{\partial f}{\partial x}(x) \right\}$  P.P.  $\uparrow$  Rademacher.

- compatible avec + et  $\circ$ .

On peut définir de même pour les fonctions de plusieurs arguments  $F(x, \theta)$ .

4) Piggyback non-lisse

$$\frac{\partial x_{k+1}(\theta)}{\partial \theta} = A \frac{\partial x_k(\theta)}{\partial \theta} + B$$

$$\downarrow \qquad \qquad \qquad \downarrow$$

$$J_{x_{k+1}}(\theta) = A \underbrace{J_{x_k}(\theta)}_{\text{requel??}} + B$$

$$\hookrightarrow J_{x_{k+1}}(\theta) = \left\{ A J + B \mid \begin{array}{l} [A, B] \in J_F(x_k(\theta), \theta) \\ J \in J_{x_k}(\theta) \end{array} \right\}$$

PyTorch

$\hookrightarrow$  selection

$$J_{x_{k+1}} = A_k J_k + B_k$$

avec  $[A_k, B_k] \in J_F(x_k(\theta), \theta)$

# 5) Une règle de la chaire infinie

Hypothèse  $J_F$  "contracte"

→ **Assumption 1 (The conservative Jacobian of the iteration mapping is a contraction)**  
 $F$  is locally Lipschitz, path differentiable, jointly in  $(x, \theta)$ , and  $J_F$  is a conservative Jacobian for  $F$ . There exists  $0 \leq \rho < 1$ , such that for any  $(x, \theta) \in \mathbb{R}^p \times \mathbb{R}^m$  and any pair  $[A, B] \in J_F(x, \theta)$ , with  $A \in \mathbb{R}^{p \times p}$  and  $B \in \mathbb{R}^{p \times m}$ , the operator norm of  $A$  is at most  $\rho$ .

$$J_{\frac{pb}{x}} : \theta \mapsto \text{Fix} \left[ J_F(\text{Fix}[F_\theta], \theta) \right]$$

$$\uparrow \text{Fix} \left[ J_F(\text{Fix}[F_\theta], \theta) \right]$$

au sens des appli  
multi-valees.

au sens des  
applications  
univoques

Thm  $J_{\frac{pb}{x}}$  est bien définie.  
BPR 22

↗  
 $\simeq$  Banach-Dicard point fixe  
 par les récursions sur les ensembles.

Cor ↗  $\lim_{k \rightarrow +\infty} \text{gap} \left( J_{x_k}(\theta), J_{\frac{pb}{x}}(\theta) \right) = 0$

$$\hookrightarrow \lim_{k \rightarrow +\infty} \text{gap} \left( \prod_{\ell=0}^k \mathcal{J}_F(x_\ell(\theta), \theta), \mathcal{J}_{\bar{c}}(\theta) \right) = 0.$$

Version "unrolled"

2) Pour presque tout  $\theta$ ,

$$\frac{\partial x_k(\theta)}{\partial \theta} \xrightarrow{k \rightarrow +\infty} \frac{\partial \bar{c}(\theta)}{\partial \theta} !$$

Note différentiation implicite pas compatible.

Note Si  $F$  a une structure "semi-afj" alors version qualitative (convergence linéaire).