

## Differentiating Solution Mappings of Non-smooth Optimization Problems



Peter Ochs  
Mathematical Optimization for Data Science  
Saarland University

— 04.09.2023 —

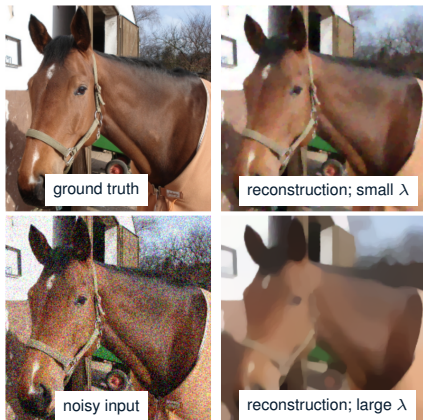


joint work: Sheheryar Mehmood

# Variational Problems with Regularization

$$\min_x E_\lambda(x), \quad E_\lambda(x) = \underbrace{D(x)}_{\text{data term}} + \lambda \underbrace{R(x)}_{\text{regularization term}} .$$

- ◆  $\lambda > 0$  is a regularization parameter.



**Regularization weight has a huge impact on the result!**

## What about learning the whole regularization term?

$$R(x, \nu, \vartheta) := \sum_{i,j} \left( \sum_{k=1}^m \nu_k \rho \left( \sum_{l=1}^L \vartheta_{kl} (K_l x)_{ij} \right) \right)$$

- ◆  $K_l$  are predefined basis filters (e.g. DCT filter)
- ◆  $\rho$  is a potential function (convex)

*(This regularizer reflects a 1-hidden-layer neural network.)*

## What about learning the whole regularization term?

$$R(x, \nu, \vartheta) := \sum_{i,j} \left( \sum_{k=1}^m \nu_k \rho \left( \sum_{l=1}^L \vartheta_{kl} (K_l x)_{ij} \right) \right)$$

- ◆  $K_l$  are predefined basis filters (e.g. DCT filter)
- ◆  $\rho$  is a potential function (convex)

*(This regularizer reflects a 1-hidden-layer neural network.)*

## How to find the best weights $\nu, \vartheta$ ?

- ◆ hand-tuning and grid search are not feasible
- ◆ sampling and regression of loss function using Gaussian processes or Random Fields (up to  $\approx 200$  parameters)
- ◆ **gradient based bi-level optimization** (several 100 000 parameters)

## Similar Problems:

- ◆ (Hyper)Parameter Learning [Domke '12], [Kunisch, Pock '13], [O. et al '16], ...
- ◆ Implicit Models [Amos, Kolter '17], [Agrawal et al. '19], [Bai, Kolter, Koltun '19], [Chen et al. '21], ...
- ◆ Meta Learning [Hospedales et al. '21], ...
- ◆ Learning to Optimize [Chen et al. '21], ...
- ◆ Sensitivity Analysis.

$$\min_{\theta \in \mathbb{R}^P} \mathcal{L}(x^*(\theta), \theta) \quad (\text{upper level})$$

$$s.t. \quad x^*(\theta) \in \arg \min_{x \in \mathbb{R}^N} E(x, \theta) \quad (\text{lower level})$$

- ◆  $\theta \in \mathbb{R}^P$ : optimization variable **parameter (vector)**.
- ◆  $\mathcal{L}: \mathbb{R}^N \times \mathbb{R}^P \rightarrow \mathbb{R}$ : **smooth loss function**.
- ◆  $E: \mathbb{R}^N \times \mathbb{R}^P \rightarrow \overline{\mathbb{R}}$ : parametric (energy) **minimization** problem.  
(convex for each  $\theta \in \mathbb{R}^P$ )
- ◆  $x^*: \mathbb{R}^P \rightarrow \mathbb{R}^N$  is a selection of the **solution mapping** of  $E$ .

$$\min_{\theta \in \mathbb{R}^P} \mathcal{L}(x^*(\theta), \theta) \quad (\text{upper level})$$

$$s.t. \quad x^*(\theta) \in \arg \min_{x \in \mathbb{R}^N} E(x, \theta) \quad (\text{lower level})$$

- ◆  $\theta \in \mathbb{R}^P$ : optimization variable **parameter (vector)**.
- ◆  $\mathcal{L}: \mathbb{R}^N \times \mathbb{R}^P \rightarrow \mathbb{R}$ : **smooth loss function**.
- ◆  $E: \mathbb{R}^N \times \mathbb{R}^P \rightarrow \overline{\mathbb{R}}$ : parametric (energy) **minimization** problem.  
(convex for each  $\theta \in \mathbb{R}^P$ )
- ◆  $x^*: \mathbb{R}^P \rightarrow \mathbb{R}^N$  is a selection of the **solution mapping** of  $E$ .

**Gradient based optimization requires  $\nabla_{\theta} \mathcal{L}(x^*(\theta^{(k)}), \theta^{(k)})$ ,**

**i.e., in particular, the sensitivity of the solution mapping  $\frac{\partial x^*}{\partial \theta}$ .**

# Bilevel optimization / parameter learning

$$\min_{\theta \in \mathbb{R}^P} \mathcal{L}(x^*(\theta), \theta) \quad (\text{upper level})$$

$$\text{s.t. } x^*(\theta) \in \arg \min_{x \in \mathbb{R}^N} E(x, \theta) \quad (\text{lower level})$$

- ◆  $\theta \in \mathbb{R}^P$ : optimization variable **parameter (vector)**.
- ◆  $\mathcal{L}: \mathbb{R}^N \times \mathbb{R}^P \rightarrow \mathbb{R}$ : **smooth loss function**.
- ◆  $E: \mathbb{R}^N \times \mathbb{R}^P \rightarrow \bar{\mathbb{R}}$ : parametric (energy) **minimization** problem.  
(convex for each  $\theta \in \mathbb{R}^P$ )
- ◆  $x^*: \mathbb{R}^P \rightarrow \mathbb{R}^N$  is a selection of the **solution mapping** of  $E$ .

**Gradient based optimization requires  $\nabla_{\theta} \mathcal{L}(x^*(\theta^{(k)}), \theta^{(k)})$ ,**

**i.e., in particular, the sensitivity of the solution mapping  $\frac{\partial x^*}{\partial \theta}$ .**

## Outline:

(I) smooth setting

(II) partly smooth setting

(III) non-smooth setting

$E$

smooth

structured non-smooth

non-smooth

**type of bilevel algorithm**

gradient

gradient

generalized derivative



## (I) Smooth Setting: Strategies for differentiation

In the following:

$$x^*(\theta) \in \arg \min_{x \in \mathbb{R}^N} E(x, \theta) \quad (E \text{ smooth})$$

Goal: Compute

$$\frac{\partial x^*}{\partial \theta}(\theta)$$

## (I) Smooth Setting: Strategies for differentiation

In the following:  $x^*(\theta) \in \arg \min_{x \in \mathbb{R}^N} E(x, \theta)$  ( $E$  smooth)

Goal: Compute  $\frac{\partial x^*}{\partial \theta}(\theta)$

### Strategies:

- ◆ Implicit differentiation. (requires  $E$  strictly convex)
- ◆ Unrolling an algorithm. (*use automatic differentiation*)
- ◆ Unrolling of a fixed point equation.

### Warning:

- ◆ Even for smooth  $E$ , the solution mapping  $x^*(\theta)$  is often non-smooth.
- ◆ Solution mapping multivalued, when  $\arg \min_{x \in \mathbb{R}^N} E(x, \theta)$  is not unique.

## (I) Implicit Differentiation (ID)

- ◆ The optimality condition is  $\nabla_x E(x, \theta) = 0$ .
- ◆ This implicitly defines  $x^*(\theta)$  (**implicit function theorem**).
- ◆ Let  $(x^*, \theta)$  be such that  $\nabla_x E(x^*, \theta) = 0$ , then, if [...] we have

$$\frac{\partial x^*}{\partial \theta}(\theta) = -\left(\frac{\partial^2 E}{\partial x^2}(x^*, \theta)\right)^{-1} \frac{\partial^2 E}{\partial \theta \partial x}(x^*, \theta).$$

## (I) Implicit Differentiation (ID)

- ◆ The optimality condition is  $\nabla_x E(x, \theta) = 0$ .
- ◆ This implicitly defines  $x^*(\theta)$  (**implicit function theorem**).
- ◆ Let  $(x^*, \theta)$  be such that  $\nabla_x E(x^*, \theta) = 0$ , then, if [...] we have

$$\frac{\partial x^*}{\partial \theta}(\theta) = -\left(\frac{\partial^2 E}{\partial x^2}(x^*, \theta)\right)^{-1} \frac{\partial^2 E}{\partial \theta \partial x}(x^*, \theta).$$

### Disadvantages:

- ◆ Requires twice differentiability of  $E$ .
- ◆ Requires several approximations:  $x^*$  and  $\left(\frac{\partial^2 E}{\partial x^2}\right)^{-1}$  (solve linear system).
- ◆ Unstable for badly conditioned  $\frac{\partial^2 E}{\partial x^2}$ .
- ◆ Requires estimation (and storing)  $\frac{\partial^2 E}{\partial x^2}$ .

## (I) Unrolling / Automatic Differentiation (AD)

- ◆ Approximate by a fixed  $n \in \mathbb{N}$ : (where  $x^{(0)}$  is some fixed initialization)

$$x^*(\theta) \approx x^{(n+1)} := \mathcal{A}^{(n+1)}(x^{(0)}, \theta) = \mathcal{A} \circ \dots \circ \mathcal{A}(x^{(0)}, \theta)$$

- ◆ Evaluate the chain rule by reverse mode AD (backpropagation):

$$\frac{\partial \mathcal{A}^{(n+1)}}{\partial \theta}(x^{(0)}, \theta) = \frac{\partial \mathcal{A}}{\partial x}(x^{(n)}, \theta) \frac{\partial \mathcal{A}^{(n)}}{\partial \theta}(x^{(0)}, \theta) + \frac{\partial \mathcal{A}}{\partial \theta}(x^{(n)}, \theta) = \dots \text{unroll} \dots$$

## (I) Unrolling / Automatic Differentiation (AD)

- ◆ Approximate by a fixed  $n \in \mathbb{N}$ : (where  $x^{(0)}$  is some fixed initialization)

$$x^*(\theta) \approx x^{(n+1)} := \mathcal{A}^{(n+1)}(x^{(0)}, \theta) = \mathcal{A} \circ \dots \circ \mathcal{A}(x^{(0)}, \theta)$$

- ◆ **Algorithmic formulation:** Set  $J_\theta^{(1)} := \frac{\partial \mathcal{A}}{\partial \theta}(x^{(0)}, \theta)$  and for  $n = 1, 2, \dots$ :

$$J_\theta^{(n+1)} = \frac{\partial \mathcal{A}}{\partial x}(x^{(n)}, \theta) J_\theta^{(n)} + \frac{\partial \mathcal{A}}{\partial \theta}(x^{(n)}, \theta)$$

## (I) Unrolling / Automatic Differentiation (AD)

- ◆ Approximate by a fixed  $n \in \mathbb{N}$ : (where  $x^{(0)}$  is some fixed initialization)

$$x^*(\theta) \approx x^{(n+1)} := \mathcal{A}^{(n+1)}(x^{(0)}, \theta) = \mathcal{A} \circ \dots \circ \mathcal{A}(x^{(0)}, \theta)$$

- ◆ **Algorithmic formulation:** Set  $J_\theta^{(1)} := \frac{\partial \mathcal{A}}{\partial \theta}(x^{(0)}, \theta)$  and for  $n = 1, 2, \dots$ :

$$J_\theta^{(n+1)} = \frac{\partial \mathcal{A}}{\partial x}(x^{(n)}, \theta) J_\theta^{(n)} + \frac{\partial \mathcal{A}}{\partial \theta}(x^{(n)}, \theta)$$

### Advantages:

- ◆ Output is **unambiguous** (unique), even in case  $\operatorname{argmin} E$  is multi-valued.
- ◆ After the algorithm  $\mathcal{A}$  and  $n \in \mathbb{N}$  are fixed, the approach is exact.
- ◆ All iterations depend on the same parameter  $\rightsquigarrow$  possibly truncate backprop.
- ◆ Easy implementation using standard AD packages.
- ◆ For  $E$  strongly convex, we have convergence rates for

$$J_\theta^{(n+1)} \xrightarrow{n \rightarrow \infty} \frac{\partial x^*}{\partial \theta}(\theta)$$

**Accelerated convergence for accelerated algorithms!** [Mehmoed, O. '20]

- ◆ **Disadvantages:** Store all intermediate iterates (in reverse mode).

## (I) Fixed Point Automatic Differentiation (FPAD)

- ◆ Replace  $\arg \min E(x, \theta)$  by a fixed point equation:

$$x^*(\theta) = \mathcal{A}(x^*(\theta), \theta)$$

- ◆ Imitate  $n \in \mathbb{N}$  iterations of  $\mathcal{A}$  starting at  $x^{(0)} := x^*$ :

$$\mathcal{A}^{(n+1)}(x^*, \theta) = \mathcal{A} \circ \dots \circ \mathcal{A}(x^*, \theta)$$



## (I) Fixed Point Automatic Differentiation (FPAD)

- ◆ Replace  $\arg \min E(x, \theta)$  by a fixed point equation:

$$x^*(\theta) = \mathcal{A}(x^*(\theta), \theta)$$

- ◆ Imitate  $n \in \mathbb{N}$  iterations of  $\mathcal{A}$  starting at  $x^{(0)} := x^*$ :

$$\mathcal{A}^{(n+1)}(x^*, \theta) = \mathcal{A} \circ \dots \circ \mathcal{A}(x^*, \theta)$$

- ◆ Evaluate the chain rule by reverse mode AD (backpropagation): For  $n = 1, 2, \dots$ :

$$\hat{J}_\theta^{(n+1)} = \frac{\partial \mathcal{A}}{\partial x}(x^*, \theta) \hat{J}_\theta^{(n)} + \frac{\partial \mathcal{A}}{\partial \theta}(x^*, \theta)$$

## (I) Fixed Point Automatic Differentiation (FPAD)

- ◆ Replace  $\arg \min E(x, \theta)$  by a fixed point equation:

$$x^*(\theta) = \mathcal{A}(x^*(\theta), \theta)$$

- ◆ Imitate  $n \in \mathbb{N}$  iterations of  $\mathcal{A}$  starting at  $x^{(0)} := x^*$ :

$$\mathcal{A}^{(n+1)}(x^*, \theta) = \mathcal{A} \circ \dots \circ \mathcal{A}(x^*, \theta)$$

- ◆ Evaluate the chain rule by reverse mode AD (backpropagation): For  $n = 1, 2, \dots$ :

$$\hat{J}_\theta^{(n+1)} = \frac{\partial \mathcal{A}}{\partial x}(x^*, \theta) \hat{J}_\theta^{(n)} + \frac{\partial \mathcal{A}}{\partial \theta}(x^*, \theta) \quad \xrightarrow{n \rightarrow \infty} \quad \frac{\partial x^*}{\partial \theta}(\theta)$$

$\rightsquigarrow$  **faster convergence rate of derivative sequence** [Mehmoor, O. '20].

(known in AD community [Gilbert '92, Christianson '94];  
connection to ID via Von Neumann series.)

# (I) Fixed Point Automatic Differentiation (FPAD)

- ◆ Replace  $\arg \min E(x, \theta)$  by a fixed point equation:

$$x^*(\theta) = \mathcal{A}(x^*(\theta), \theta)$$

- ◆ Imitate  $n \in \mathbb{N}$  iterations of  $\mathcal{A}$  starting at  $x^{(0)} := x^*$ :

$$\mathcal{A}^{(n+1)}(x^*, \theta) = \mathcal{A} \circ \dots \circ \mathcal{A}(x^*, \theta)$$

- ◆ Evaluate the chain rule by reverse mode AD (backpropagation): For  $n = 1, 2, \dots$ :

$$\hat{J}_\theta^{(n+1)} = \frac{\partial \mathcal{A}}{\partial x}(x^*, \theta) \hat{J}_\theta^{(n)} + \frac{\partial \mathcal{A}}{\partial \theta}(x^*, \theta) \quad \xrightarrow{n \rightarrow \infty} \quad \frac{\partial x^*}{\partial \theta}(\theta)$$

↪ **faster convergence rate of derivative sequence** [Mehmoor, O. '20].

(known in AD community [Gilbert '92, Christianson '94];  
connection to ID via Von Neumann series.)

- ◆ Equivalent to **AD with intermediate variables replaced by the optimum**: [O. et al. '16]

$$J_\theta^{(n+1)} = \frac{\partial \mathcal{A}}{\partial x}(x^{(n)}, \theta) J_\theta^{(n)} + \frac{\partial \mathcal{A}}{\partial \theta}(x^{(n)}, \theta)$$

# (I) Fixed Point Automatic Differentiation (FPAD)

- ◆ Replace  $\arg \min E(x, \theta)$  by a fixed point equation:

$$x^*(\theta) = \mathcal{A}(x^*(\theta), \theta)$$

- ◆ Imitate  $n \in \mathbb{N}$  iterations of  $\mathcal{A}$  starting at  $x^{(0)} := x^*$ :

$$\mathcal{A}^{(n+1)}(x^*, \theta) = \mathcal{A} \circ \dots \circ \mathcal{A}(x^*, \theta)$$

- ◆ Evaluate the chain rule by reverse mode AD (backpropagation): For  $n = 1, 2, \dots$ :

$$\hat{J}_\theta^{(n+1)} = \frac{\partial \mathcal{A}}{\partial x}(x^*, \theta) \hat{J}_\theta^{(n)} + \frac{\partial \mathcal{A}}{\partial \theta}(x^*, \theta) \quad \xrightarrow{n \rightarrow \infty} \quad \frac{\partial x^*}{\partial \theta}(\theta)$$

↪ **faster convergence rate of derivative sequence** [Mehmood, O. '20].

(known in AD community [Gilbert '92, Christianson '94];  
connection to ID via Von Neumann series.)

- ◆ Equivalent to **AD with intermediate variables replaced by the optimum**: [O. et al. '16]

$$J_\theta^{(n+1)} = \frac{\partial \mathcal{A}}{\partial x}(x^{(n)}, \theta) J_\theta^{(n)} + \frac{\partial \mathcal{A}}{\partial \theta}(x^{(n)}, \theta)$$

## Key advantages:

- ◆ Requires to store only  $x^*$ .
- ◆ Easy implementation using standard AD packages.

## (II) Structured Non-Smooth Setting: Strategies for differentiation

In the following:

$$x^*(\theta) \in \arg \min_{x \in \mathbb{R}^N} E(x, \theta) \quad (E \text{ structured non-smooth})$$

Goal: Compute

$$\frac{\partial x^*}{\partial \theta}(\theta) \quad !?$$

## (II) Structured Non-Smooth Setting: Strategies for differentiation

In the following:  $x^*(\theta) \in \arg \min_{x \in \mathbb{R}^N} E(x, \theta)$  ( $E$  structured non-smooth)

Goal: Compute  $\frac{\partial x^*}{\partial \theta}(\theta)$  !?

Strategies that compute (classic) gradients:

◆ Strategies above after smoothing  $E$  (ID, AD, FPAD):

↪ often instable or requires significant smoothing

↪ no approximation bounds

## (II) Structured Non-Smooth Setting: Strategies for differentiation

In the following:  $x^*(\theta) \in \arg \min_{x \in \mathbb{R}^N} E(x, \theta)$  ( $E$  structured non-smooth)

Goal: Compute  $\frac{\partial x^*}{\partial \theta}(\theta)$  !?

### Strategies that compute (classic) gradients:

◆ Strategies above after smoothing  $E$  (ID, AD, FPAD):

- ↪ often instable or requires significant smoothing
- ↪ no approximation bounds

◆ **Unrolling a “smooth algorithm”** that solves the non-smooth problem. [O. et al. '16]

- $\mathcal{A}^{(n+1)}(x^{(0)}, \theta) \rightarrow x^*(\theta)$  for  $n \rightarrow \infty$ , where  $\mathcal{A}$  is a smooth mapping.
- **Idea:** Assert that iterates lie in interior of constraint set.  
(e.g. *Bregman Proximal Gradient Method*)

- ↪ **Smoothing is controlled by the number of iterations**  $n \in \mathbb{N}$ .
- ↪ Convergence of derivative sequence was **not** studied.  
Limit requires generalized notion of derivatives.

## (II) Structured Non-Smooth Setting: Strategies for differentiation

In the following:  $x^*(\theta) \in \arg \min_{x \in \mathbb{R}^N} E(x, \theta)$  ( $E$  structured non-smooth)

Goal: Compute  $\frac{\partial x^*}{\partial \theta}(\theta)$  !?

### Strategies that compute (classic) gradients:

◆ Strategies above after smoothing  $E$  (ID, AD, FPAD):

↪ often instable or requires significant smoothing

↪ no approximation bounds

◆ **Unrolling a “smooth algorithm”** that solves the non-smooth problem. [O. et al. '16]

●  $\mathcal{A}^{(n+1)}(x^{(0)}, \theta) \rightarrow x^*(\theta)$  for  $n \rightarrow \infty$ , where  $\mathcal{A}$  is a smooth mapping.

● **Idea:** Assert that iterates lie in interior of constraint set.

(e.g. *Bregman Proximal Gradient Method*)

↪ **Smoothing is controlled by the number of iterations**  $n \in \mathbb{N}$ .

↪ Convergence of derivative sequence was **not** studied.

Limit requires generalized notion of derivatives.

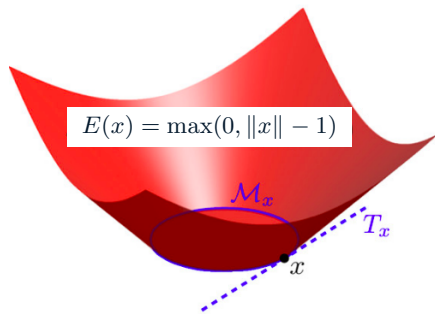
● **FPAD variant** requires  $x^*$  to lie in the interior of the constraint set.



## (II) Differentiation Strategies for Partly Smooth Functions

**Definition** of **partial smoothness** for convex  $E$  from [Liang, Fadili, Peyré '14] (original [Lewis '02]):  
 $E$  with  $\partial E(x) \neq \emptyset$  is **partly smooth at  $x$  relative to  $\mathcal{M} \ni x$** , if

- ◆ (Smoothness)  $\mathcal{M}$  is a  $C^2$ -Manifold around  $x$  and  $E|_{\mathcal{M}} \in C^2$ ,
- ◆ (Sharpness) the normal space  $\mathcal{N}_{\mathcal{M}}$  is  $\text{par}(\partial E(x))$ , and
- ◆ (Continuity)  $\partial E$  is continuous at  $x$  relative to  $\mathcal{M}$ .



**Examples:**  $\ell_1$ -norm,  $\ell_{2,1}$ -norm,  $\ell_\infty$ -norm, nuclear norm, TV-norm, ...

from [G. Peyré, talk "Low Complexity Regularization of Inverse Problems", 2014]

## (III) Implicit Differentiation (ID) under Partial Smoothness

Consider:  $x^*(\theta) \in \arg \min_{x \in \mathbb{R}^N} E(x, \theta)$

**Theorem:** [ID under Partial Smoothness [Lewis '02], [Vaiter et al. '17]]

**Assumptions:**

- ◆  $E$  is partly smooth relative to  $\mathcal{M} \times \Omega$  where  $\Omega \subset \mathbb{R}^P$  is open, and
- ◆ for some  $(x^*, \theta^*)$ , the following **restricted positive definiteness** holds

$$\forall v \in \mathcal{T}_{x^*} \mathcal{M}: \quad \langle v, \nabla_{\mathcal{M}}^2 E(x^*, \theta^*) v \rangle > 0$$

and **non-degeneracy** holds

$$0 \in \text{ri}(\partial_x E(x^*, \theta^*)).$$

**Then the manifold version of the Implicit Differentiation (ID) holds:**

- ◆ There exists a neighborhood  $\Theta$  of  $\theta^*$  and a  $C^1$  mapping  $\psi: \Theta \rightarrow \mathcal{M}$  such that
- ◆  $\Theta \ni \theta \mapsto \psi(\theta) = \operatorname{argmin}_x E(x, \theta)$ , and
- ◆  $\Theta \ni \theta \mapsto \frac{\partial \psi}{\partial \theta}(\theta) = -\nabla_{\mathcal{M}}^2 E(\psi(\theta), \theta)^\dagger \frac{\partial}{\partial \theta} \nabla_{\mathcal{M}} E(\psi(\theta), \theta)$ .

## (II) Unrolling / Automatic Differentiation (AD) under Partial Smoothness

- ◆ **Idea:** Many algorithms have the **finite identification property** (see papers by [\[J. Liang\]](#)):

$$\exists n_0 \in \mathbb{N}: \quad x^{(n)} \in \mathcal{M} \text{ for all } n \geq n_0,$$

hence the update mapping  $\mathcal{A}$  becomes smooth eventually.

**Examples:** Forward–backward Splitting (Proximal Gradient Descent), FISTA, ...

## (II) Unrolling / Automatic Differentiation (AD) under Partial Smoothness

- ◆ **Idea:** Many algorithms have the **finite identification property** (see papers by [J. Liang]):

$$\exists n_0 \in \mathbb{N}: \quad x^{(n)} \in \mathcal{M} \text{ for all } n \geq n_0,$$

hence the update mapping  $\mathcal{A}$  becomes smooth eventually.

**Examples:** Forward–backward Splitting (Proximal Gradient Descent), FISTA, ...

**Theorem:** [AD for Forward–backward Splitting [Liang et al. '14], [Mehmood, O. '22]]

**Assumptions:**

- ◆ Partial smoothness, restricted positive definiteness, the non-degeneracy assumption, and convergence of the iterates.

**Then,**

- ◆ the rate of convergence for  $(x^{(n)})_{n \in \mathbb{N}}$  is actually **locally linear**.
- ◆ If, additionally,  $x^{(0)}$  is close enough to  $x^*$ , then

$$\frac{\partial \mathcal{A}^{(n+1)}}{\partial \theta}(x^{(0)}, \theta) =: J_{\theta}^{(n)} \xrightarrow{n \rightarrow \infty} \frac{\partial \psi}{\partial \theta}(\theta).$$

↪ **Crucial issue:** Derivatives are not defined for  $x^{(n)} \notin \mathcal{M}$ .

- ◆ Linear convergence, if all involved derivatives are locally Lipschitz.

## (III) Fixed Point Automatic Differentiation (FPAD)

### Remedy of Differentiation Issue:

- ◆ Since  $x^* \in \mathcal{M}$ , by definition, all iteration mappings in FPAD are well defined.

### Remedy of Differentiation Issue:

- ◆ Since  $x^* \in \mathcal{M}$ , by definition, all iteration mappings in FPAD are well defined.

**Theorem:** [Convergence of FPAD for Forward–backward Splitting [Mehmood, O. '22]]

#### Assumptions:

- ◆ Partial smoothness, restricted positive definiteness and the non-degeneracy.

Then,

- ◆ the derivative sequence converges linearly

$$\frac{\partial \mathcal{A}^{(n+1)}}{\partial \theta}(x^*, \theta) =: \hat{J}_\theta^{(n)} \xrightarrow{n \rightarrow \infty} \frac{\partial \psi}{\partial \theta}(\theta).$$

## (III) Non-Smooth Setting: Strategies for Generalized Differentiation

In the following:

$$x^*(\theta) \in \arg \min_{x \in \mathbb{R}^N} E(x, \theta) \quad (E \text{ non-smooth})$$

Goal: Compute

$$\frac{\partial x^*}{\partial \theta}(\theta) \quad !?$$

## (III) Non-Smooth Setting: Strategies for Generalized Differentiation

In the following:  $x^*(\theta) \in \arg \min_{x \in \mathbb{R}^N} E(x, \theta)$  ( $E$  non-smooth)

Goal: Compute  $\frac{\partial x^*}{\partial \theta}(\theta)$  !?

### Strategies:

- ◆ Weak differentiation of iterative algorithms. [Deledalle et al. '14]  
(using Rademacher Theorem; show update map is Lipschitz)
  - ↪ yields only weak derivatives (not defined pointwise)
  - ↪ convergence of derivative sequence unknown
- ◆ Using generalized non-smooth derivatives. [Bolte et al. '21]
  - ↪ *Some details on next slides.*



**Definition:** [Bolte, Pauwels '19]

A locally Lipschitz function  $F: \mathbb{R}^N \rightarrow \mathbb{R}$  is **path differentiable**, if

- ◆ there exists a nowhere empty-valued mapping  $D_F: \mathbb{R}^N \rightrightarrows \mathbb{R}^N$
- ◆ that has a closed graph and is locally bounded, such that
- ◆ for any absolutely continuous curve  $\gamma: [0, 1] \rightarrow \mathbb{R}^N$

$$\frac{d}{dt} F \circ \gamma(t) = A \dot{\gamma}(t), \quad \forall A \in D_F(\gamma(t)), \quad \text{a.e. for } t \in [0, 1].$$

- ◆ Call  $D_F$  **conservative Jacobian** of  $F$ .

**Construction idea:**

- ↪ By definition, **path differentiable functions admit a chain rule**.
- ↪ Formalization of backpropagation used in AD packages for neural networks.
- ↪ Straightforward generalization to locally Lipschitz mappings.

### (III) Implicit Path Differentiation (ID)

**Theorem:** [Implicit Path Differentiation (ID) [Bolte et al. '21], [Bolte, Pauwels, Silveti-Falls '23]]

**Assumptions:**

◆  $F: \mathbb{R}^N \times \mathbb{R}^P \rightarrow \mathbb{R}^N$  is path differentiable.

◆  $(x^*, \theta^*)$  satisfies the optimality condition

$$F(x^*, \theta^*) = 0.$$

◆  $D_F(x^*, \theta^*)$  is convex and  $\forall [A, B] \in D_F(x^*, \theta^*)$ , we have that  $A$  is invertible.

**Then the path differentiable version of the Implicit Differentiation (ID) holds:**

◆  $\exists \Theta$  neighborhood of  $\theta^*$  and path differentiable  $\Psi: \Theta \rightarrow \mathbb{R}^N$  such that

$$\forall \theta \in \Theta: F(\Psi(\theta), \theta) = 0,$$

◆ and  $\Theta \ni \theta \mapsto D_\Psi(\theta) = \{-A^{-1}B \mid [A, B] \in D_F(\Psi(\theta), \theta)\}$ .

**Key Disadvantage:**

◆ Invertibility condition for all  $A$  above; Sufficient condition by strong monotonicity of  $F$ .

### (III) Unrolling / Automatic Differentiation (AD)

#### Idea:

- ◆ Compositions of path differentiable functions are path differentiable.
- ◆ Path differentiable functions admit a chain rule.

#### Unrolling:

- ◆ Evaluate the chain rule. Set  $\mathcal{J}_\theta^{(1)} := \{J_\theta^{(1)} \mid [A_x^{(1)}, J_\theta^{(1)}] = D_{\mathcal{A}}(x^{(0)}, \theta)\}$  and for  $n = 1, 2, \dots$ :

$$\mathcal{J}_\theta^{(n+1)} = \{A_x^{(n)} J_\theta^{(n)} + A_\theta^{(n)} \mid [A_x^{(n)}, A_\theta^{(n)}] \in D_{\mathcal{A}}(x^{(n)}, \theta), J_\theta^{(n)} \in \mathcal{J}_\theta^{(n)}\}.$$

**Theorem:** [Convergence of AD under path differentiability [Bolte, Pauwels, Vaiter '22]]

#### Assumption:

- ◆  $\mathcal{A}$  is locally Lipschitz with unique fixed point  $x^*(\theta)$ , path differentiable in  $(x, \theta)$ ,
- ◆ and  $\exists \rho \in [0, 1)$ :  $\forall (x, \theta)$  and all  $[A, B] \in D_{\mathcal{A}}(x, \theta)$ , we have  $\|A\| \leq \rho$ .

#### Then:

- ◆ The gap between  $\mathcal{J}_\theta^{(n+1)}$  and  $\text{fix}(D_{\mathcal{A}}(x^*(\theta), \theta))$  vanishes as  $n \rightarrow \infty$ .
- ◆ Under additional assumptions on a Lipschitz gradient selection structure, the convergence is linear.

## Bi-level optimization / parameter learning

$$\min_{\theta \in \mathbb{R}^P} \mathcal{L}(x^*(\theta), \theta) \quad (\text{upper level})$$

$$\text{s.t. } x^*(\theta) \in \arg \min_{x \in \mathbb{R}^N} E(x, \theta) \quad (\text{lower level})$$

**Goal:** Compute  $\frac{\partial x^*}{\partial \theta}(\theta)$  **!?**

Outline:	$E$	type of bilevel algorithm
(I) smooth setting	smooth	gradient
(II) partly smooth setting	structured non-smooth	gradient
(III) non-smooth setting	non-smooth	generalized derivative

### For each part:

- ◆ Implicit differentiation
- ◆ Automatic Differentiation / Backpropagation
- ◆ Fixed Point Automatic Differentiation
- ◆ Generalized derivatives of  $x^*(\theta)$  requires algorithms for non-smooth bi-level opt.