# A New Look at First Order Methods Lifting the Lipschitz Gradient Continuity Restriction

Marc Teboulle

School of Mathematical Sciences Tel Aviv University

Joint work with H. Bauschke and J. Bolte

Optimization and Discrete Geometry: Theory and Practice 24-26 April, 2018, Tel Aviv University



#### Recall: The Basic Pillar underlying FOM

 $\inf \{\Phi(x) := f(x) + g(x) : x \in \mathbb{R}^d\}, f, g \text{ convex, with } g \in C^1.$ 

Captures many applied problems, and the source for fundamental FOM.

Usual key assumption: g admits L-Lipschitz continuous gradient on  $\mathbb{R}^d$ 



## Recall: The Basic Pillar underlying FOM

$$\inf \{\Phi(x) := f(x) + g(x) : x \in \mathbb{R}^d\}, f, g \text{ convex}, \text{ with } g \in C^1.$$

Captures many applied problems, and the source for fundamental FOM.

Usual key assumption: g admits L-Lipschitz continuous gradient on  $\mathbb{R}^d$ 

A simple, yet key consequence of this, is the so-called descent Lemma:

$$g(x) \leq g(y) + \langle 
abla g(y), x - y 
angle + rac{L}{2} \|x - y\|^2, \ orall x, y \in \mathbb{R}^d.$$

This inequality naturally provides

- 1. An upper quadratic approximation of g
- 2. A crucial pillar in the analysis of current FOM.

However, in many contexts and applications:

- $\ominus$  the differentiable function g <u>does not</u> have a L-smooth gradient
- ⊖ Hence precludes direct use of basic FOM methodology and schemes.



# FOM Beyond Lipschitz Gradient Continuity

#### Goals/Outline:

- Circumvent the longstanding question of Lipschitz Gradient continuity imposed on FOM.
- Derive FOM "free" from this smoothness assumption, with guaranteed complexity estimates and convergence results.
- Apply our results to a broad class of important problems lacking smooth gradients.



Marc Teboulle (Tel Aviv University)

A New Look at First Order Methods Lifting the Lipschitz Gradient Continuity Restriction

Consider the descent Lemma for the smooth  $g \in C_L^{1,1}$  on  $\mathbb{R}^d$ :

$$g(x) \leq g(y) + \langle x - y, \nabla g(y) \rangle + \frac{L}{2} ||x - y||^2, \ \forall x, y \in \mathbb{R}^d.$$



Consider the descent Lemma for the smooth  $g \in C^{1,1}_L$  on  $\mathbb{R}^d$ :

$$g(x) \leq g(y) + \langle x - y, \nabla g(y) \rangle + \frac{L}{2} ||x - y||^2, \ \forall x, y \in \mathbb{R}^d.$$

Simple algebra shows that it can be equivalently written as:

$$\left(rac{L}{2}\|x\|^2 - g(x)
ight) - \left(rac{L}{2}\|y\|^2 - g(y)
ight) \geq \langle Ly - 
abla g(y), x - y 
angle \quad orall x, y \in \mathbb{R}^d$$

Consider the descent Lemma for the smooth  $g \in C_L^{1,1}$  on  $\mathbb{R}^d$ :

$$g(x) \leq g(y) + \langle x - y, \nabla g(y) \rangle + \frac{L}{2} ||x - y||^2, \ \forall x, y \in \mathbb{R}^d.$$

Simple algebra shows that it can be equivalently written as:

$$\left(rac{L}{2}\|x\|^2 - g(x)
ight) - \left(rac{L}{2}\|y\|^2 - g(y)
ight) \ge \langle Ly - 
abla g(y), x - y 
angle \quad orall x, y \in \mathbb{R}^d$$

Nothing else but the gradient inequality for the convex  $\frac{L}{2}||x||^2 - g(x)$  !

Thus, for a given smooth function g on  $\mathbb{R}^d$ 

$$\text{Descent Lemma} \quad \Longleftrightarrow \quad \frac{\mathsf{L}}{2} \| x \|^2 - g(x) \text{ is convex on } \mathbb{R}^d.$$



Consider the descent Lemma for the smooth  $g \in C_L^{1,1}$  on  $\mathbb{R}^d$ :

$$g(x) \leq g(y) + \langle x - y, \nabla g(y) \rangle + \frac{L}{2} ||x - y||^2, \ \forall x, y \in \mathbb{R}^d.$$

Simple algebra shows that it can be equivalently written as:

$$\left(rac{L}{2}\|x\|^2 - g(x)
ight) - \left(rac{L}{2}\|y\|^2 - g(y)
ight) \geq \langle Ly - 
abla g(y), x - y 
angle \quad orall x, y \in \mathbb{R}^d$$

Nothing else but the gradient inequality for the convex  $\frac{L}{2} ||x||^2 - g(x)$  !

Thus, for a given smooth function g on  $\mathbb{R}^d$ 

Descent Lemma 
$$\iff \frac{L}{2} \|\mathbf{x}\|^2 - \mathbf{g}(\mathbf{x})$$
 is convex on  $\mathbb{R}^d$ .

**Capture the Geometry of Constraint/Objective** Naturally suggests to replace the *squared norm* with a general convex function  $h(\cdot)$  that captures the geometry of the constraint/objective.



# A Lipschitz-Like Convexity Condition

Following our basic observation: Replace the  $\|\cdot\|^2$  with a convex *h*.

- ▶ Trade *L*-smooth gradient of g on  $\mathbb{R}^d$  with
- ▶ Convexity condition on couple (g, h), dom  $g \supset$  dom  $h, g \in C^1($ int dom h).

# A Lipschitz-like/Convexity Condition (LC) $\exists L > 0$ with Lh - g convex on int dom h,

- ► Condition (LC) ⇔ New descent Lemma we seek for.
- It also naturally leads to the well-known Bregman distance.

# A Descent Lemma without Lipschitz Gradient Continuity

Lemma (Descent lemma without Lipschitz Gradient Continuity) The condition (LC): Lh - g convex on int dom h is equivalent to

 $g(x) \leq g(y) + \langle \nabla g(y), x - y \rangle + LD_h(x, y), \ \forall (x, y) \in dom \ h \times int \ dom \ h$ 

 $D_h$  stands for the Bregman Distance associated to a convex h:

$$D_h(x,y) := h(x) - h(y) - \langle \nabla h(y), x - y \rangle, \ \forall x \in \text{dom } h, y \in \text{int dom } h.$$

**Proof of Descent Lemma.**  $D_{Lh-g}(x, y) \ge 0$  for the convex function Lh - g!



# A Descent Lemma without Lipschitz Gradient Continuity

Lemma (Descent lemma without Lipschitz Gradient Continuity) The condition (LC): Lh - g convex on int dom h is equivalent to

 $g(x) \leq g(y) + \langle \nabla g(y), x - y \rangle + LD_h(x, y), \ \forall (x, y) \in dom \ h \times int \ dom \ h$ 

 $D_h$  stands for the Bregman Distance associated to a convex h:

$$D_h(x,y) := h(x) - h(y) - \langle \nabla h(y), x - y \rangle, \ \forall x \in \operatorname{dom} h, y \in \operatorname{int} \operatorname{dom} h.$$

**Proof of Descent Lemma.**  $D_{Lh-g}(x, y) \ge 0$  for the convex function Lh - g!

**Distance-Like Properties** - For all  $(x, y) \in \text{dom } h \times \text{int dom } h$ 

- $x \to D_h(x, y)$  is convex with *h* convex.
- $D_h(x, y) \ge 0$  and " = 0" iff x = y.(h strictly convex).
- ▶ However, note that *D<sub>h</sub>* is in general not symmetric!

The use of Bregman distances in optimization started with Bregman (67). For initial works and main results on *Proximal Bregman Algorithms:* [Censor-Zenios (92), T. (92), Chen-T. (93), Eckstein (93), Bauschke-Borwein (97).]

Marc Teboulle (Tel Aviv University)

A New Look at First Order Methods Lifting the Lipschitz Gradient Continuity Restriction

# Some Useful Examples for Bregman Distances D<sub>h</sub>

Each example is a one dimensional convex h. The corresponding function  $\tilde{h}$  and Bregman distance in  $\mathbb{R}^d$  simply use the formulae

$$\widetilde{h}(x) = \sum_{j=1}^n h(x_j) ext{ and } D_{\widetilde{h}}(x,y) = \sum_{j=1}^n D_h(x_j,y_j).$$

Name	h	dom <i>h</i>
Energy	$\frac{1}{2}x^2$	$\mathbb{R}$
Boltzmann-Shannon entropy	$x \log x$	$[0,\infty)$
Burg's entropy	$-\log x$	$(0,\infty)$
Fermi-Dirac entropy	$x\log x + (1-x)\log(1-x)$	[0, 1]
Hellinger	$-(1-x^2)^{1/2}$	[-1, 1]
Fractional Power	$(px-x^p)/(1-p), p\in(0,1)$	$[0,\infty)$

► Other possible/useful kernels h include: Nonseparable Bregman, e.g., any convex h on ℝ<sup>d</sup> as well as for handling matrix problems: PSD matrices, cone constraints, etc.., [details in Auslender and T. (2005)].



# The Convex Model and Blanket Assumption

Our aim is to solve the composite convex problem

$$v(\mathcal{P}) = \inf\{\Phi(x) := f(x) + g(x) \mid x \in \overline{\mathrm{dom}} h\},\$$

where  $\overline{\operatorname{dom}} h$  denotes the closure of dom h,

Under the following standard assumption.

The "Hidden h " (in unconstrained case) will adapt to Nonlinear Geometry of  ${\mathcal P}$ 

#### Blanket Assumption as Usual:

- (i)  $f: \mathbb{R}^d \to (-\infty, \infty]$  is proper lower semicontinuous (lsc) convex,
- (ii)  $h: \mathbb{R}^d \to (-\infty, \infty]$  is proper, lsc convex.
- (iii)  $g: \mathbb{R}^d \to (-\infty, \infty]$  is proper lsc convex with dom  $g \supset \text{dom } h$  and  $g \in C^1(\text{int dom } h)$
- (iv) dom  $f \cap$  int dom  $h \neq \emptyset$ ,

$$(\mathsf{v}) -\infty < \mathsf{v}(\mathcal{P}) = \inf\{\Phi(x) : x \in \overline{\mathrm{dom}} \ h\} = \inf\{\Phi(x) : x \in \mathrm{dom} \ h\}.$$

Algorithm NoLips for  $\inf \{f(x) + g(x) : x \in C \equiv \overline{\operatorname{dom}} h\}$ 

Main Algorithmic Operator- [Reduces to classical prox-grad, when h quadratic]

$$\mathsf{T}_\lambda(\mathsf{x}) := \mathsf{argmin}\left\{\mathsf{f}(\mathsf{u}) + \mathsf{g}(\mathsf{x}) + \langle \nabla \mathsf{g}(\mathsf{x}), \mathsf{u} - \mathsf{x} \rangle + \frac{1}{\lambda}\mathsf{D}_\mathsf{h}(\mathsf{u},\mathsf{x}) : \mathsf{u} \in \overline{\mathsf{dom}}\,\mathsf{h}\right\} \ (\lambda > \mathbf{0}).$$

**NoLips Main Iteration:**  $x \in \text{int dom } h$ ,  $x^+ = T_{\lambda}(x)$ ,  $(\lambda > 0)$ .



Algorithm NoLips for  $\inf \{f(x) + g(x) : x \in C \equiv \overline{\operatorname{dom}} h\}$ 

Main Algorithmic Operator- [Reduces to classical prox-grad, when h quadratic]

$$\mathsf{T}_\lambda(\mathsf{x}) := \operatorname{argmin} \left\{ \mathsf{f}(\mathsf{u}) + \mathsf{g}(\mathsf{x}) + \langle \nabla \mathsf{g}(\mathsf{x}), \mathsf{u} - \mathsf{x} \rangle + \frac{1}{\lambda} \mathsf{D}_\mathsf{h}(\mathsf{u},\mathsf{x}) : \mathsf{u} \in \overline{\mathsf{dom}} \, \mathsf{h} \right\} \ (\lambda > \mathbf{0}).$$

**NoLips Main Iteration:**  $x \in \text{int dom } h$ ,  $x^+ = T_{\lambda}(x)$ ,  $(\lambda > 0)$ .

#### Algorithm NoLips – in More Details

- 0. Input. Choose a convex function h such that there exists L > 0 with Lh g convex on int dom h.
- 1. Initialization. Start with any  $x^0 \in \text{int dom } h$ .
- 2. Recursion. For each  $k \ge 1$  with  $\lambda_k > 0$ , generate  $\{x^k\}_{k \in \mathbb{N}} \in \text{int dom } h$  via

$$x^{k} = T_{\lambda_{k}}(x^{k-1}) = \underset{x}{\operatorname{argmin}} \left\{ f(x) + \left\langle \nabla g(x^{k-1}), x - x^{k-1} \right\rangle + \frac{1}{\lambda_{k}} D_{h}(x, x^{k-1}) \right\}$$

Main Issues / Questions for NoLips

- Well posedness and Computation of  $T_{\lambda}(\cdot)$ ?
- What is the complexity of NoLips?
- Does NoLips converge to an optimal solution?
- In particular: Can we identify the <u>most aggressive</u> step-size in terms of problem's data?

## NoLips is Well Defined

We assume *h* is a Legendre function [Rockafellar 70].

▶ *h* is strictly convex and differentiable on int dom  $h \neq \emptyset$  and

dom  $\partial h$  = int dom h with  $\partial h(x) = \{\nabla h(x)\}, \forall x \in \text{int dom } h$ .

▶  $\|\nabla h(x^k)\| \to \infty$  whenever  $\{x^k\} \subset \text{int dom } h, x^k \to x \in \mathsf{Bdy}(\mathsf{dom } h.)$ 

With *h* Legengre:  $\nabla h$  is a *bijection* from int dom  $h \rightarrow \text{int dom } h^*$  and  $(\nabla h)^{-1} = \nabla h^*$ .

# NoLips is Well Defined

We assume *h* is a Legendre function [Rockafellar 70].

▶ *h* is strictly convex and differentiable on int dom  $h \neq \emptyset$  and

dom  $\partial h$  = int dom h with  $\partial h(x) = \{\nabla h(x)\}, \forall x \in \text{int dom } h$ .

▶  $\|\nabla h(x^k)\| \to \infty$  whenever  $\{x^k\} \subset \text{int dom } h, x^k \to x \in \mathsf{Bdy}(\mathsf{dom } h.)$ 

With *h* Legengre:  $\nabla h$  is a *bijection* from int dom  $h \rightarrow$  int dom  $h^*$  and  $(\nabla h)^{-1} = \nabla h^*$ .

#### Note:

- ► Legendre functions "abound" for defining useful *D<sub>h</sub>*. (All previous examples and more..).
- Crucial for deriving meaningful convergence results.

Equipped with the above, one can prove (see technical details in our paper.)

## Lemma (Well posedness of the method)

The proximal gradient map  $T_{\lambda} \neq \emptyset$ , is single-valued and maps int dom h in int dom h.



# NoLips – Decomposition of $T_{\lambda}(\cdot)$ into Elementary Steps

 $T_{\lambda}$  shares the same structural decomposition as the usual proximal gradient. It splits into *"elementary" steps* useful for computational purposes.



# NoLips – Decomposition of $T_{\lambda}(\cdot)$ into Elementary Steps

 $T_{\lambda}$  shares the same structural decomposition as the usual proximal gradient. It splits into *"elementary" steps* useful for computational purposes.

#### Define Bregman gradient map

$$p_{\lambda}(x) := \operatorname{argmin} \left\{ \lambda \langle 
abla g(x), u 
angle + D_h(u, x) : u \in \mathbb{R}^d 
ight\} \equiv 
abla h^*(
abla h(x) - \lambda 
abla g(x))$$

Clearly reduces to the usual explicit gradient step when  $h = \frac{1}{2} \| \cdot \|^2$ .

#### Define the proximal Bregman map

$$\operatorname{prox}_{\lambda f}^{h}(y) := \operatorname{argmin} \left\{ \lambda f(u) + D_{h}(u, y) : u \in \mathbb{R}^{d} \right\}, \ y \in \operatorname{int} \operatorname{dom} h$$

# NoLips – Decomposition of $T_{\lambda}(\cdot)$ into Elementary Steps

 $T_{\lambda}$  shares the same structural decomposition as the usual proximal gradient. It splits into *"elementary" steps* useful for computational purposes.

#### ⊕ Define Bregman gradient map

$$p_{\lambda}(x) := \operatorname{argmin} \left\{ \lambda \langle 
abla g(x), u 
angle + D_h(u, x) : u \in \mathbb{R}^d 
ight\} \equiv 
abla h^* (
abla h(x) - \lambda 
abla g(x))$$

Clearly reduces to the usual explicit gradient step when  $h = \frac{1}{2} \| \cdot \|^2$ .

#### **⊕** Define the proximal Bregman map

$$\operatorname{prox}_{\lambda f}^{h}(y) := \operatorname{argmin} \left\{ \lambda f(u) + D_{h}(u, y) : u \in \mathbb{R}^{d} \right\}, \ y \in \operatorname{int} \operatorname{dom} h$$

One can show NoLips  $\equiv$  Composition of these two Bregman maps:

**NoLips Main Iteration:**  $x \in \text{int dom } h$ ,  $x^+ = \text{prox}_{\lambda f}^h \circ p_{\lambda}(x)$  ( $\lambda > 0$ )

For Specific and Useful Examples, see the paper.

# The Key Estimation Inequality for Analyzing NoLips

Lemma (Descent inequality for NoLips) Let  $\lambda > 0$ . For all x in int dom h, let  $x^+ := T_{\lambda}(x)$ . Then,

$$\lambda\left(\Phi(x^+)-\Phi(u)
ight)\leq D_h(u,x)-D_h(u,x^+)-(1-\lambda L)D_h(x^+,x),\ orall u\in dom\,h.$$



## The Key Estimation Inequality for Analyzing NoLips

Lemma (Descent inequality for NoLips) Let  $\lambda > 0$ . For all x in int dom h, let  $x^+ := T_{\lambda}(x)$ . Then,

 $\lambda\left(\Phi(x^+)-\Phi(u)\right) \leq D_h(u,x)-D_h(u,x^+)-(1-\lambda L)D_h(x^+,x), \ \forall u \in \textit{dom} \ h.$ 

Proof simply combines the NoLips Descent Lemma with known old results:

[Lemma 3.1 and Lemma 3.2 – Chen and T. (1993)].

1. (The three points identity) For any  $x, y \in int(dom h)$  and  $u \in dom h$ :

$$D_h(u,y) - D_h(u,x) - D_h(x,y) = \langle \nabla h(y) - \nabla h(x), x - u \rangle.$$

2. (Bregman Based Proximal Inequality) Given  $z \in int \text{ dom } h$ , define

$$u^+:= \operatorname{argmin} \{ arphi(u) + t^{-1} D_h(u,z): \ u \in X \}; \ arphi ext{ convex}, \quad t > 0.$$

Then, for any  $u \in \text{dom } h$ ,

$$t(\varphi(u^+)-\varphi(u)) \leq D_h(u,z) - D_h(u,u^+) - D_h(u^+,z).$$

Complexity for NoLips: O(1/k)

## Theorem (NoLips: Complexity)

(i) (Global estimate in function values) Let  $\{x^k\}_{k\in\mathbb{N}}$  be the sequence generated by NoLips with  $\lambda \in (0, 1/L]$ . Then

$$\Phi(x^k) - \Phi(u) \leq \frac{LD_h(u, x^0)}{k} \qquad \forall u \in dom h.$$

 (ii) (Complexity for h with closed domain) Assume in addition, that dom h = dom h and that (P) has at least a solution. Then for any solution x of (P),

$$\Phi(x^k) - \min_C \Phi \leq \frac{LD_h(\bar{x}, x^0)}{k}$$

**Notes**  $\diamond$  The entropies of Boltzmann-Shannon, Fermi-Dirac and Hellinger are non trivial examples for which the assumption ( $\overline{\text{dom } h} = \text{dom } h$ ) holds.

 $\diamond$  When  $h(x) = \frac{1}{2} ||x||^2$ ,  $g \in C_L^{1,1}$ , and we thus recover the classical sublinear global rate of the usual proximal gradient method.



Marc Teboulle (Tel Aviv University)

A New Look at First Order Methods Lifting the Lipschitz Gradient Continuity Restriction

#### Yes! by introducing an interesting notion of symmetry for $D_h$ .

Bregman distances are in general not symmetric, except when h is the energy.



#### Yes! by introducing an interesting notion of symmetry for $D_h$ .

Bregman distances are in general not symmetric, except when h is the energy.

Definition (Symmetry coefficient-Measures Lack of Symmetry) Let  $h : \mathbb{R}^d \to (-\infty, \infty]$  be a Legendre function. Its symmetry coefficient is defined by

$$\alpha(h) := \inf \left\{ \frac{D_h(x, y)}{D_h(y, x)} : (x, y) \in \operatorname{int} \operatorname{dom} h \times \operatorname{int} \operatorname{dom} h, \ x \neq y \right\}$$

#### Yes! by introducing an interesting notion of symmetry for $D_h$ .

Bregman distances are in general not symmetric, except when h is the energy.

Definition (Symmetry coefficient-Measures Lack of Symmetry) Let  $h : \mathbb{R}^d \to (-\infty, \infty]$  be a Legendre function. Its symmetry coefficient is defined by

$$\alpha(h) := \inf \left\{ \frac{D_h(x,y)}{D_h(y,x)} : (x,y) \in \operatorname{int} \operatorname{dom} h \times \operatorname{int} \operatorname{dom} h, \ x \neq y \right\}$$

#### **Properties of the Symmetry Coefficient** $\alpha(h)$ :

- $\alpha(h) \in [0,1]$ . The closer is  $\alpha(h)$  to 1 the more symmetric  $D_h$  is.
- $\alpha(h) = 1$  Perfect symmetry! when h is the energy.
- $\alpha(h) = 0$  Total lack of symmetry, e.g.,  $h(x) = x \log x$  and  $h(x) = -\log x$ .
- $\alpha(h) > 0$  Some symmetry.., e.g.,  $h(x) = x^4$ ,  $\alpha(h) = 2 \sqrt{3}$ .

The symmetry coefficient allows to determine the best step size of NoLips for pointwise convergence of the generated sequence  $\{x^k\}$ .

# The Step-Size Choice $\lambda$ for Global Convergence of NoLips

efining Step Size in Terms of Problem's Data
$$0 < \lambda \leq rac{(1 + lpha(h)) - \delta}{L}$$
 for some  $\delta \in (0, 1 + lpha(h)),$ 

- ▶  $[0,1] \ni \alpha(h)$  is the symmetry coefficient of h.
- L > 0 is the constant in condition (LC) Lh g convex.
  - When h(·) := || · ||<sup>2</sup>/2, then α(h) = 1, L is usual Lipchitz constant for ∇g and above reduces to

$$0 < \lambda \le \frac{2-\delta}{L}$$

recovers the classical step size allowed for pointwise convergence of the classical proximal gradient method [Combettes-Wajs 05].

 With the above step-size choice, we can establish global convergence of the sequence {x<sup>k</sup>} generated by NoLips.

D

# Pointwise Convergence for NoLips

Theorem (NoLips: Point convergence - With  $\lambda \in (0, L^{-1}(1 + \alpha(h)))$ 

Assume that the solution set  $S^*$  of  $(\mathcal{P})$  is nonsempty. Then, the following holds.

- (i) (Subsequential convergence) If S<sup>\*</sup> is compact, any limit point of {x<sup>k</sup>}<sub>k∈ℕ</sub> is a solution to (P).
- (ii) (Global convergence) Assume that dom h = dom h and that (H) is satisfied. Then the sequence  $\{x^k\}_{k \in \mathbb{N}}$  converges to some solution  $x^*$  of  $(\mathcal{P})$ .

**Note** Nontrivial examples: Boltzmann-Shannon, Fermi-Dirac and Hellinger entropies satisfy the set of assumptions in **H** and  $\overline{\text{dom } h} = \text{dom } h$ .

Additional assumption on  $D_h$  is to ensure separation properties of  $D_h$  at the boundary.

#### Assumption H:

- (i) For every  $x \in \text{dom } h$  and  $\beta \in \mathbb{R}$ , the level set  $\{y \in \text{int dom } h : D_h(x, y) \le \beta\}$  is bounded.
- (ii) If  $\{x^k\}_{k\in\mathbb{N}}$  converges to some x in dom h then  $D_h(x, x^k) \to 0$ .
- (iii) Reciprocally, if x is in dom h and if  $\{x^k\}_{k\in\mathbb{N}}$  is such that  $D_h(x, x^k) \to 0$ , then  $x^k \to x$ .

Applications - A Prototype: Linear Inverse Problems with Poisson Noise **A very large class of problems arising in Statistical and Image Sciences areas:** inverse problems where data measurements are collected by counting discrete events (e.g., photons, electrons) contaminated by noise described by a Poisson process.

Huge amount of literature: astronomy, nuclear medicine (PET), electronic microscropy, statistical estimation (EM), image deconvolution, denoising speckle (multiplicative) noise, ect....

**Problem:** Given a matrix  $A \in \mathbb{R}^{m \times n}_+$  and  $b \in \mathbb{R}^m_{++}$  the goal is to reconstruct the signal/image  $x \in \mathbb{R}^n_+$  from the noisy measurements b such that  $Ax \simeq b$ .

Applications - A Prototype: Linear Inverse Problems with Poisson Noise **A very large class of problems arising in Statistical and Image Sciences areas:** inverse problems where data measurements are collected by counting discrete events (e.g., photons, electrons) contaminated by noise described by a Poisson process.

Huge amount of literature: astronomy, nuclear medicine (PET), electronic microscropy, statistical estimation (EM), image deconvolution, denoising speckle (multiplicative) noise, ect....

**Problem:** Given a matrix  $A \in \mathbb{R}^{m \times n}_+$  and  $b \in \mathbb{R}^m_{++}$  the goal is to reconstruct the signal/image  $x \in \mathbb{R}^n_+$  from the noisy measurements b such that  $Ax \simeq b$ .

A natural proximity measure in  $\mathbb{R}^{n}_{+}$  - (Kullback-Liebler Divergence):

$$\mathcal{D}(b,Ax) := \sum_{i=1}^m \{b_i \log \frac{b_i}{(Ax)_i} + (Ax)_i - b_i\}.$$

which (up to some const.) is the negative Poisson log-likelihood function.

- ▶ The optimization problem: ( $\mathbb{E}$ ) minimize  $\{\mu f(x) + g(x) : x \in \mathbb{R}^n_+\}$
- $g(x) \equiv D(d, Ax)$ , f a regularizer smooth or nonsmooth,  $\mu > 0$
- $x \to \mathcal{D}(b, Ax)$  convex, but does not admit a globally Lipschitz continuous gradient.

## NoLips in Action : New Simple Schemes for Many Problems

The optimization problem will be of the form:

 $(\mathbb{E}) \qquad \min_{x} \{ \mu f(x) + \mathcal{D}_{\phi}(b, Ax) \} \text{ or } \qquad \min_{x} \{ \mu f(x) + \mathcal{D}_{\phi}(Ax, b) \}$ 

where  $g(x) := D_{\phi}(b, Ax)$  for some convex  $\phi$ , and f(x) some convex regularizer.

# NoLips in Action : New Simple Schemes for Many Problems

The optimization problem will be of the form:

(E)  $\min_{\mathbf{v}} \{ \mu f(\mathbf{x}) + \mathcal{D}_{\phi}(\mathbf{b}, A\mathbf{x}) \}$  or  $\min_{\mathbf{v}} \{ \mu f(\mathbf{x}) + \mathcal{D}_{\phi}(A\mathbf{x}, \mathbf{b}) \}$ 

where  $g(x) := D_{\phi}(b, Ax)$  for some convex  $\phi$ , and f(x) some convex regularizer.

Applying NoLips requires:

- 1. To pick an adequate h, so that Lh g convex; L in terms of problem's data.
- 2. In turns, this determines the step-size  $\lambda$  defined through  $(L, \alpha(h))$ .
- 3. Compute  $p_{\lambda}(\cdot)$  and  $\operatorname{prox}_{\lambda f}^{h}(\cdot)$ ) Bregman gradient and proximal steps.

Our convergence/complexity results hold and produce new simple algorithms:

Simple schemes via explicit map  $M_j(\cdot)$ 

$$x > 0,$$
  $x_j^+ = M_j(b, A, x; \mu, \lambda) \cdot x_j,$   $j = 1, \ldots, n.$ 

 $\checkmark$ 

#### Two Simple Algorithms for Poisson Linear Inverse Problems

Given  $g(x) := D_{\phi}(b, Ax)$  (  $\phi(u) = u \log u$ ), to apply NoLips:

- We take  $h(x) = -\sum_{j=1}^{n} \log x_j$ , dom  $h = \mathbb{R}^n_{++}$ .
- We need to find L > 0 such that Lh g is convex in  $\mathbb{R}^{n}_{++}$ .

#### Two Simple Algorithms for Poisson Linear Inverse Problems

Given  $g(x) := D_{\phi}(b, Ax)$  (  $\phi(u) = u \log u$ ), to apply NoLips:

- We take  $h(x) = -\sum_{j=1}^{n} \log x_j$ , dom  $h = \mathbb{R}^n_{++}$ .
- We need to find L > 0 such that Lh g is convex in  $\mathbb{R}^{n}_{++}$ .

**Lemma.** With (g, h) above, Lh - g is convex on  $\mathbb{R}^{n}_{++}$  for any  $L \geq ||b||_1 := \sum_{i=1}^{m} b_i$ .

#### Two Simple Algorithms for Poisson Linear Inverse Problems

Given  $g(x) := D_{\phi}(b, Ax)$  (  $\phi(u) = u \log u$ ), to apply NoLips:

- We take  $h(x) = -\sum_{j=1}^{n} \log x_j$ , dom  $h = \mathbb{R}^n_{++}$ .
- We need to find L > 0 such that Lh g is convex in  $\mathbb{R}^{n}_{++}$ .

**Lemma.** With (g, h) above, Lh - g is convex on  $\mathbb{R}^{n}_{++}$  for any  $L \ge \|b\|_{1} := \sum_{i=1}^{m} b_{i}$ .

Thus, we can take  $\lambda = L^{-1} = \|b\|_1^{-1}$ , and applying NoLips with  $x \in \mathbb{R}_{++}^n$  reads:

$$x^+ = \operatorname{argmin} \left\{ \mu f(u) + \langle 
abla g(x), u 
angle + \|b\|_1 \sum_{j=1}^n \left( rac{u_j}{x_j} - \log rac{u_j}{x_j} - 1 
ight) : u > 0 
ight\}.$$

The above yields closed form algorithms for Poisson reconstruction problems with two typical regularizers.



#### Example 1 – Sparse Poisson Linear Inverse Problem

**Sparse regularization.** Let  $f(x) := \mu ||x||_1$ , known to promote sparsity. Define,

$$c_j(x):=\sum_{i=1}^m b_irac{a_{ij}}{\langle a_i,x
angle}, \; r_j:=\sum_i a_{ij}>0.$$

**NoLips for Sparse Poisson Linear Inverse Problems** 

$$x_j > 0, \ x_j^+ = \frac{\|b\|_1 x_j}{\|b\|_1 + (\mu x_j + x_j(r_j - c_j(x)))}, \ j = 1, \dots n$$



#### Example 1 – Sparse Poisson Linear Inverse Problem

**Sparse regularization.** Let  $f(x) := \mu ||x||_1$ , known to promote sparsity. Define,

$$c_j(x):=\sum_{i=1}^m b_irac{a_{ij}}{\langle a_i,x
angle}, \; r_j:=\sum_i a_{ij}>0.$$

**NoLips for Sparse Poisson Linear Inverse Problems** 

$$x_j > 0, \ x_j^+ = \frac{\|b\|_1 x_j}{\|b\|_1 + (\mu x_j + x_j(r_j - c_j(x)))}, \ j = 1, \dots n$$

**Special Case:**  $\mu = 0$ , ( $\mathbb{E}$ ) is the Poisson Maximum Likelihood Estimation.

NoLips yields in that case: A New Scheme for Poisson MLE

$$x_j > 0, \ x_j^+ = \frac{\|b\|_1 x_j}{\|b\|_1 + x_j(r_j - c_j(x))}, \ j = 1, \dots n.$$

\$

## Example 2 - Thikhonov - Poisson Linear Inverse Problems

**Tikhonov regularization.** Let  $f(x) := \mu ||x||^2/2$ . Recall that this term is used as a penalty in order to promote solutions of Ax = b with *small Euclidean norms*.



#### Example 2 - Thikhonov - Poisson Linear Inverse Problems

**Tikhonov regularization.** Let  $f(x) := \mu ||x||^2/2$ . Recall that this term is used as a penalty in order to promote solutions of Ax = b with *small Euclidean norms*.

Using previous notation, NoLips yields a

" A Poisson-Thikonov method" : Set  $\lambda = \|b\|_1^{-1}$  and start with  $x \in \mathbb{R}^n_{++}$ 

$$x_j^+=rac{\sqrt{
ho_j^2(x)+4\mu\lambda x_j^2-
ho_j(x)}}{2\mu\lambda x_j},\ j=1,\ldots,n.$$

where

$$\rho_j(x) := 1 + \lambda x_j \left( r_j - c_j(x) \right), \ j = 1, \ldots, n.$$

As just mentioned, many other interesting methods can be considered

- By choosing different kernels for  $\phi$ , or
- By reversing the order of the arguments in the proximity measure (which is not symmetric!..hence defining different problems, see the paper.)

# Conclusion and More Details/Results on NoLips

Proposed framework offers a new paragdim for FOM

- **•** Breaks the longstanding question asking for L-smooth gradient.
- ▶ Proven Complexity and Pointwise Convergence as Classical case.
- ▶ Allows to derive new FOM without Lipschitz gradient.

Details and More Results: Bauschke H., Bolte J., and Teboulle M.

"A Descent Lemma beyond Lipshitz Gradient Continuity: First Order Methods Revisited and Applications". *Mathematics of Operations Research*, (2017), 330–348.

Available Online http://dx.doi.org/10.1287/moor.2016.0817



# Conclusion and More Details/Results on NoLips

Proposed framework offers a new paragdim for FOM

- **•** Breaks the longstanding question asking for L-smooth gradient.
- ▶ Proven Complexity and Pointwise Convergence as Classical case.
- ▶ Allows to derive new FOM without Lipschitz gradient.

Details and More Results: Bauschke H., Bolte J., and Teboulle M.

"A Descent Lemma beyond Lipshitz Gradient Continuity: First Order Methods Revisited and Applications". *Mathematics of Operations Research*, (2017), 330–348.

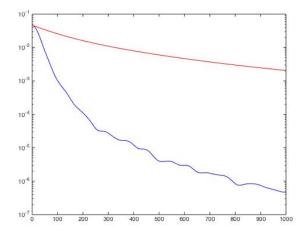
Available Online http://dx.doi.org/10.1287/moor.2016.0817



#### THANK YOU FOR LISTENING!



NoLips (Red) Versus a FAST NoLips (Blue)...





Marc Teboulle (Tel Aviv University)

A New Look at First Order Methods Lifting the Lipschitz Gradient Continuity Restriction