

BREGMAN PROXIMAL ALGORITHMS FOR COMPOSITE AND FINITE-SUM NONCONVEX MINIMIZATION PROBLEMS

Themelis, Andreas
Kyushu University

Latafat, Puya
KU Leuven

Ahookhosh, Masoud
University of Antwerp

Patrinos, Panagiotis
KU Leuven

<https://hdl.handle.net/2324/6790349>

出版情報：2021-07-22
バージョン：
権利関係：



KYUSHU UNIVERSITY

BREGMAN PROXIMAL ALGORITHMS

FOR COMPOSITE AND FINITE-SUM NONCONVEX MINIMIZATION PROBLEMS



Andreas Themelis

Kyushu University

andreas.themelis@ees.kyushu-u.ac.jp

Puya Latafat

KU Leuven



Masoud Ahookhosh

University of Antwerp
Universiteit
Antwerpen



Panagiotis Patrinos

KU Leuven



First-Order Methods with Convergence Analysis @ SIAM OP21

July 22, 2021

Outline

Proximal algorithm(s)

Proximal point algorithm

Simplest analysis

Bregman proximal algorithm(s)

Why Bregman

Linesearch B-FBS: the BELLA algorithm

Newton-type?

Linesearch?

Block-coordinate B-FBS: Bregman Finito/MISO

Conclusions



Proximal algorithm(s)

Proximal point algorithm

Problem

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \varphi(x) \quad (\text{P})$$

A1. $\varphi : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ proper, lsc

A2. $\inf \varphi > -\infty$

Proximal mapping

$$\text{prox}_{\gamma\varphi}(x) = \arg \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + \frac{1}{2\gamma} \|w - x\|^2 \right\} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n \quad (\gamma > 0)$$

Method

Proximal point algorithm (PPA)

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

$\overline{\mathbb{R}} := \mathbb{R} \cup \{\infty\}$ is the extended-real line

lsc: lower semicontinuous

Proximal algorithm(s)

Proximal point algorithm: simplest convergence analysis

Tool

Moreau envelope

$$\varphi^\gamma(x) = \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + \frac{1}{2\gamma} \|w - x\|^2 \right\} : \mathbb{R}^n \rightarrow \mathbb{R}$$

P1. $\varphi^\gamma \leq \varphi$

P3. φ^γ continuous

P2. $\varphi^\gamma(x) = \varphi(\bar{x}) + \frac{1}{2\gamma} \|\bar{x} - x\|^2$

P4. $\frac{x - \bar{x}}{\gamma} \in \hat{\partial}\varphi(\bar{x}) \quad (\bar{x} \in \text{prox}_{\gamma\varphi}(x))$

Convergence of PPA

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

$$\varphi(x^{k+1})$$

$\hat{\partial}\varphi(z) = \left\{ v \mid \liminf_{z \neq w \rightarrow z} \frac{\varphi(w) - \varphi(z) - \langle v, w - z \rangle}{\|w - z\|} \geq 0 \right\}$ is the (Fréchet) subdifferential of φ at z

Proximal algorithm(s)

Proximal point algorithm: simplest convergence analysis

Tool

Moreau envelope

$$\varphi^\gamma(x) = \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + \frac{1}{2\gamma} \|w - x\|^2 \right\} : \mathbb{R}^n \rightarrow \mathbb{R}$$

P1. $\varphi^\gamma \leq \varphi$

P3. φ^γ continuous

P2. $\varphi^\gamma(x) = \varphi(\bar{x}) + \frac{1}{2\gamma} \|\bar{x} - x\|^2$

P4. $\frac{x - \bar{x}}{\gamma} \in \hat{\partial}\varphi(\bar{x}) \quad (\bar{x} \in \text{prox}_{\gamma\varphi}(x))$

Convergence of PPA

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

$$\varphi^\gamma(x^{k+1}) \stackrel{\text{P1}}{\leq} \varphi(x^{k+1})$$

Proximal algorithm(s)

Proximal point algorithm: simplest convergence analysis

Tool

Moreau envelope

$$\varphi^\gamma(x) = \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + \frac{1}{2\gamma} \|w - x\|^2 \right\} : \mathbb{R}^n \rightarrow \mathbb{R}$$

P1. $\varphi^\gamma \leq \varphi$

P3. φ^γ continuous

P2. $\varphi^\gamma(x) = \varphi(\bar{x}) + \frac{1}{2\gamma} \|\bar{x} - x\|^2$

P4. $\frac{x - \bar{x}}{\gamma} \in \hat{\partial}\varphi(\bar{x}) \quad (\bar{x} \in \text{prox}_{\gamma\varphi}(x))$

Convergence of PPA

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

$$\varphi^\gamma(x^{k+1}) \stackrel{\text{P1}}{\leq} \varphi(x^{k+1}) \stackrel{\text{P2}}{=} \varphi^\gamma(x^k) - \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \quad (\text{SD})$$

Proximal algorithm(s)

Proximal point algorithm: simplest convergence analysis

Tool

Moreau envelope

$$\varphi^\gamma(x) = \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + \frac{1}{2\gamma} \|w - x\|^2 \right\} : \mathbb{R}^n \rightarrow \mathbb{R}$$

P1. $\varphi^\gamma \leq \varphi$

P3. φ^γ continuous

P2. $\varphi^\gamma(x) = \varphi(\bar{x}) + \frac{1}{2\gamma} \|\bar{x} - x\|^2$

P4. $\frac{x - \bar{x}}{\gamma} \in \hat{\partial}\varphi(\bar{x}) \quad (\bar{x} \in \text{prox}_{\gamma\varphi}(x))$

Convergence of PPA

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

$$\varphi^\gamma(x^{k+1}) \stackrel{\text{P1}}{\leq} \varphi(x^{k+1}) \stackrel{\text{P2}}{=} \varphi^\gamma(x^k) - \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \quad (\text{SD})$$

► (SD) $\Rightarrow \varphi^\gamma(x^k) \& \varphi(x^k) \searrow \varphi_\star$

Proximal algorithm(s)

Proximal point algorithm: simplest convergence analysis

Tool

Moreau envelope

$$\varphi^\gamma(x) = \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + \frac{1}{2\gamma} \|w - x\|^2 \right\} : \mathbb{R}^n \rightarrow \mathbb{R}$$

P1. $\varphi^\gamma \leq \varphi$

P3. φ^γ continuous

P2. $\varphi^\gamma(x) = \varphi(\bar{x}) + \frac{1}{2\gamma} \|\bar{x} - x\|^2$

P4. $\frac{x - \bar{x}}{\gamma} \in \hat{\partial}\varphi(\bar{x}) \quad (\bar{x} \in \text{prox}_{\gamma\varphi}(x))$

Convergence of PPA

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

$$\varphi^\gamma(x^{k+1}) \stackrel{\text{P1}}{\leq} \varphi(x^{k+1}) \stackrel{\text{P2}}{=} \varphi^\gamma(x^k) - \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \quad (\text{SD})$$

► (SD) $\Rightarrow \varphi^\gamma(x^k) \& \varphi(x^k) \searrow \varphi_\star$

► if $(x^k)_{k \in K} \rightarrow x^\star$, then $\varphi(x^\star) \stackrel{lsc}{\leq} \varphi_\star = \varphi^\gamma(x^\star) \stackrel{\text{P1}}{\leq} \varphi(x^\star)$

Proximal algorithm(s)

Proximal point algorithm: simplest convergence analysis

Tool

Moreau envelope

$$\varphi^\gamma(x) = \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + \frac{1}{2\gamma} \|w - x\|^2 \right\} : \mathbb{R}^n \rightarrow \mathbb{R}$$

P1. $\varphi^\gamma \leq \varphi$

P3. φ^γ continuous

P2. $\varphi^\gamma(x) = \varphi(\bar{x}) + \frac{1}{2\gamma} \|\bar{x} - x\|^2$

P4. $\frac{x - \bar{x}}{\gamma} \in \hat{\partial}\varphi(\bar{x}) \quad (\bar{x} \in \text{prox}_{\gamma\varphi}(x))$

Convergence of PPA

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

$$\varphi^\gamma(x^{k+1}) \stackrel{\text{P1}}{\leq} \varphi(x^{k+1}) \stackrel{\text{P2}}{=} \varphi^\gamma(x^k) - \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \quad (\text{SD})$$

► (SD) $\Rightarrow \varphi^\gamma(x^k) \& \varphi(x^k) \searrow \varphi_\star = \varphi(x^\star)$

► if $(x^k)_{k \in K} \rightarrow x^\star$, then $\varphi(x^\star) \stackrel{lsc}{\leq} \varphi_\star = \varphi^\gamma(x^\star) \stackrel{\text{P1}}{\leq} \varphi(x^\star)$

Proximal algorithm(s)

Proximal point algorithm: simplest convergence analysis

Tool

Moreau envelope

$$\varphi^\gamma(x) = \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + \frac{1}{2\gamma} \|w - x\|^2 \right\} : \mathbb{R}^n \rightarrow \mathbb{R}$$

P1. $\varphi^\gamma \leq \varphi$

P3. φ^γ continuous

P2. $\varphi^\gamma(x) = \varphi(\bar{x}) + \frac{1}{2\gamma} \|\bar{x} - x\|^2$

P4. $\frac{\bar{x} - x}{\gamma} \in \hat{\partial}\varphi(\bar{x}) \quad (\bar{x} \in \text{prox}_{\gamma\varphi}(x))$

Convergence of PPA

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

$$\varphi^\gamma(x^{k+1}) \stackrel{\text{P1}}{\leq} \varphi(x^{k+1}) \stackrel{\text{P2}}{=} \varphi^\gamma(x^k) - \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \quad (\text{SD})$$

- ▶ (SD) $\Rightarrow \varphi^\gamma(x^k) \& \varphi(x^k) \searrow \varphi_\star = \varphi(x^\star)$
- ▶ if $(x^k)_{k \in K} \rightarrow x^\star$, then $\varphi(x^\star) \stackrel{lsc}{\leq} \varphi_\star = \varphi^\gamma(x^\star) \stackrel{\text{P3}}{\leq} \varphi^\gamma(x^k) \stackrel{\text{P1}}{\leq} \varphi(x^k)$
- ▶ (SD) $\Rightarrow \sum_{k \in \mathbb{N}} \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \leq \varphi^\gamma(x^0) - \inf \varphi (< \infty)$

Proximal algorithm(s)

Proximal point algorithm: simplest convergence analysis

Tool

Moreau envelope

$$\varphi^\gamma(x) = \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + \frac{1}{2\gamma} \|w - x\|^2 \right\} : \mathbb{R}^n \rightarrow \mathbb{R}$$

P1. $\varphi^\gamma \leq \varphi$

P3. φ^γ continuous

P2. $\varphi^\gamma(x) = \varphi(\bar{x}) + \frac{1}{2\gamma} \|\bar{x} - x\|^2$

P4. $\frac{x - \bar{x}}{\gamma} \in \hat{\partial}\varphi(\bar{x}) \quad (\bar{x} \in \text{prox}_{\gamma\varphi}(x))$

Convergence of PPA

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

$$\varphi^\gamma(x^{k+1}) \stackrel{\text{P1}}{\leq} \varphi(x^{k+1}) \stackrel{\text{P2}}{=} \varphi^\gamma(x^k) - \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \quad (\text{SD})$$

- ▶ (SD) $\Rightarrow \varphi^\gamma(x^k) \& \varphi(x^k) \searrow \varphi_\star = \varphi(x^\star)$
- ▶ if $(x^k)_{k \in K} \rightarrow x^\star$, then $\varphi(x^\star) \stackrel{lsc}{\leq} \varphi_\star = \varphi^\gamma(x^\star) \leq \varphi(x^\star)$
- ▶ (SD) $\Rightarrow \sum_{k \in \mathbb{N}} \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \leq \varphi^\gamma(x^0) - \inf \varphi (< \infty)$
$$\Rightarrow 0 \leftarrow \frac{x^{k-1} - x^k}{\gamma}$$

Proximal algorithm(s)

Proximal point algorithm: simplest convergence analysis

Tool

Moreau envelope

$$\varphi^\gamma(x) = \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + \frac{1}{2\gamma} \|w - x\|^2 \right\} : \mathbb{R}^n \rightarrow \mathbb{R}$$

P1. $\varphi^\gamma \leq \varphi$

P3. φ^γ continuous

P2. $\varphi^\gamma(x) = \varphi(\bar{x}) + \frac{1}{2\gamma} \|\bar{x} - x\|^2$

P4. $\frac{x - \bar{x}}{\gamma} \in \hat{\partial}\varphi(\bar{x}) \quad (\bar{x} \in \text{prox}_{\gamma\varphi}(x))$

Convergence of PPA

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

$$\varphi^\gamma(x^{k+1}) \stackrel{\text{P1}}{\leq} \varphi(x^{k+1}) \stackrel{\text{P2}}{=} \varphi^\gamma(x^k) - \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \quad (\text{SD})$$

- ▶ (SD) $\Rightarrow \varphi^\gamma(x^k) \& \varphi(x^k) \searrow \varphi_\star = \varphi(x^\star)$
- ▶ if $(x^k)_{k \in K} \rightarrow x^\star$, then $\varphi(x^\star) \stackrel{lsc}{\leq} \varphi_\star = \varphi^\gamma(x^\star) \leq \varphi(x^\star)$
- ▶ (SD) $\Rightarrow \sum_{k \in \mathbb{N}} \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \leq \varphi^\gamma(x^0) - \inf \varphi (< \infty)$
 $\Rightarrow 0 \leftarrow \frac{x^{k+1} - x^k}{\gamma} \stackrel{\text{P4}}{\in} \hat{\partial}\varphi(x^k)$

Proximal algorithm(s)

Proximal point algorithm: simplest convergence analysis

Tool

Moreau envelope

$$\varphi^\gamma(x) = \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + \frac{1}{2\gamma} \|w - x\|^2 \right\} : \mathbb{R}^n \rightarrow \mathbb{R}$$

P1. $\varphi^\gamma \leq \varphi$

P3. φ^γ continuous

P2. $\varphi^\gamma(x) = \varphi(\bar{x}) + \frac{1}{2\gamma} \|\bar{x} - x\|^2$

P4. $\frac{\bar{x} - x}{\gamma} \in \hat{\partial}\varphi(\bar{x}) \quad (\bar{x} \in \text{prox}_{\gamma\varphi}(x))$

Convergence of PPA

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

$$\varphi^\gamma(x^{k+1}) \stackrel{\text{P1}}{\leq} \varphi(x^{k+1}) \stackrel{\text{P2}}{=} \varphi^\gamma(x^k) - \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \quad (\text{SD})$$

► (SD) $\Rightarrow \varphi^\gamma(x^k)$ & $\varphi(x^k) \searrow \varphi_\star = \varphi(x^\star)$

► if $(x^k)_{k \in K} \rightarrow x^\star$, then $\varphi(x^\star) \stackrel{lsc}{\leq} \varphi_\star = \varphi^\gamma(x^\star) \leq \varphi(x^\star)$

► (SD) $\Rightarrow \sum_{k \in \mathbb{N}} \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \leq \varphi^\gamma(x^0) - \inf \varphi (< \infty)$

$$\Rightarrow 0 \leftarrow \frac{x^{k-1} - x^k}{\gamma} \stackrel{\text{P4}}{\in} \hat{\partial}\varphi(x^k), \quad \varphi(x^k) \rightarrow \varphi(x^\star)$$

Proximal algorithm(s)

Proximal point algorithm: simplest convergence analysis

Tool

Moreau envelope

$$\varphi^\gamma(x) = \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + \frac{1}{2\gamma} \|w - x\|^2 \right\} : \mathbb{R}^n \rightarrow \mathbb{R}$$

P1. $\varphi^\gamma \leq \varphi$

P3. φ^γ continuous

P2. $\varphi^\gamma(x) = \varphi(\bar{x}) + \frac{1}{2\gamma} \|\bar{x} - x\|^2$

P4. $\frac{\bar{x} - x}{\gamma} \in \hat{\partial}\varphi(\bar{x}) \quad (\bar{x} \in \text{prox}_{\gamma\varphi}(x))$

Convergence of PPA

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

$$\varphi^\gamma(x^{k+1}) \stackrel{\text{P1}}{\leq} \varphi(x^{k+1}) \stackrel{\text{P2}}{=} \varphi^\gamma(x^k) - \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \quad (\text{SD})$$

- ▶ (SD) $\Rightarrow \varphi^\gamma(x^k) \& \varphi(x^k) \searrow \varphi_\star = \varphi(x^\star)$
- ▶ if $(x^k)_{k \in K} \rightarrow x^\star$, then $\varphi(x^\star) \stackrel{lsc}{\leq} \varphi_\star = \varphi^\gamma(x^\star) \leq \varphi(x^\star)$
- ▶ (SD) $\Rightarrow \sum_{k \in \mathbb{N}} \frac{1}{2\gamma} \|x^{k+1} - x^k\|^2 \leq \varphi^\gamma(x^0) - \inf \varphi (< \infty)$
 $\Rightarrow 0 \leftarrow \frac{x^{k-1} - x^k}{\gamma} \stackrel{\text{P4}}{\in} \hat{\partial}\varphi(x^k), \varphi(x^k) \rightarrow \varphi(x^\star) \Rightarrow 0 \in \partial\varphi(x^\star)$

$\partial\varphi(z) := \{v \mid \exists (z^k, v^k) \rightarrow (z, v) \text{ with } v^k \in \hat{\partial}\varphi(z^k) \& \varphi(z^k) \rightarrow \varphi(z)\}$

Proximal algorithm(s)

Proximal point algorithm: simplest convergence analysis

PPA

φ proper, lsc, lower bounded

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

ω set of limit points of $(x^k)_{k \in \mathbb{N}}$:

1. $\varphi(x^k), \varphi^\gamma(x^k) \searrow \varphi_\star$
2. $\varphi \equiv \varphi^\gamma \equiv \varphi_\star$ on ω
3. $0 \in \partial\varphi(x^\star)$ for any $x^\star \in \omega$
4. φ coercive $\Rightarrow (x^k)_{k \in \mathbb{N}}$ bounded, and $\omega \neq \emptyset$ compact & connected

φ is coercive (AKA level bounded) if $\liminf_{\|x\| \rightarrow \infty} \varphi(x) = \infty$

Proximal algorithm(s)

Proximal point algorithm: simplest convergence analysis

PPA

φ proper, lsc, lower bounded

$$x^{k+1} \in \text{prox}_{\gamma\varphi}(x^k)$$

ω set of limit points of $(x^k)_{k \in \mathbb{N}}$:

1. $\varphi(x^k), \varphi^\gamma(x^k) \searrow \varphi_\star$
2. $\varphi \equiv \varphi^\gamma \equiv \varphi_\star$ on ω
3. $0 \in \partial\varphi(x^\star)$ for any $x^\star \in \omega$
4. φ coercive $\Rightarrow (x^k)_{k \in \mathbb{N}}$ bounded, and $\omega \neq \emptyset$ compact & connected

PPA $_\lambda$ (λ -relaxed PPA)

$$\begin{cases} \bar{x}^k \in \text{prox}_{\gamma\varphi}(x^k) \\ x^{k+1} = (1 - \lambda)x^k + \lambda\bar{x}^k \end{cases} \stackrel{\text{same arguments}^*}{\Rightarrow} \varphi^\gamma(x^{k+1}) \leq \varphi^\gamma(x^k) - \frac{\lambda(2-\lambda)}{2\gamma} \|\bar{x}^k - x^k\|^2$$

All the claims of PPA hold for PPA $_\lambda$ for any $\lambda \in (0, 2)$

$$\star \varphi^\gamma(x^+) \leq \varphi(\bar{x}) + \frac{1}{2\gamma} \|\bar{x} - x^+\|^2 = \varphi^\gamma(x) - \frac{1}{2\gamma} \|\bar{x} - x\|^2 + \frac{1}{2\gamma} \|\bar{x} - x^+\|^2 = \varphi^\gamma(x) - \frac{\lambda(2-\lambda)}{2\gamma} \|\bar{x} - x\|^2$$

Outline

Proximal algorithm(s)

 Proximal point algorithm

 Simplest analysis

Bregman proximal algorithm(s)

 Why Bregman



Linesearch B-FBS: the BELLA algorithm

 Newton-type?

 Linesearch?

Block-coordinate B-FBS: Bregman Finito/MISO

Conclusions

Bregman proximal algorithm(s)

Why? ~ Structured problems

Proximal Point

- 😊 extremely simple analysis
- 😢 not “practical”: $\text{prox}_{\gamma\varphi}$ usually hard to compute

Proximal gradient method (AKA Forward-backward splitting)

In many applications, $\varphi = \text{“smooth } f” + \text{“easy-to-prox } g”$

$$\begin{aligned} x^{k+1} &\in \text{prox}_{\gamma g}\left(x^k - \gamma \nabla f(x^k)\right) \\ &= \arg \min_{w \in \mathbb{R}^n} \left\{ f(x^k) + \langle \nabla f(x^k), w - x^k \rangle + g(w) + \frac{1}{2\gamma} \|w - x^k\|^2 \right\} \end{aligned}$$

smooth: differentiable with ∇f Lipschitz continuous (Lipschitz constant := L_f)

Bregman proximal algorithm(s)

Why? ~ Structured problems

Proximal Point

- ⊕ extremely simple analysis
- ⊖ not “practical”: $\text{prox}_{\gamma\varphi}$ usually hard to compute

Proximal gradient method (AKA Forward-backward splitting)

In many applications, $\varphi = \text{“smooth } f\text{”} + \text{“easy-to-prox } g\text{”}$

$$\begin{aligned} x^{k+1} &\in \text{prox}_{\gamma g}\left(x^k - \gamma \nabla f(x^k)\right) \\ &= \arg \min_{w \in \mathbb{R}^n} \left\{ f(x^k) + \langle \nabla f(x^k), w - x^k \rangle + g(w) + \frac{1}{2\gamma} \|w - x^k\|^2 \right\} \\ &= \arg \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + D_{\hat{h}}(w, x^k) \right\} \end{aligned}$$

$\hat{h} := \frac{1}{2\gamma} \|\cdot\|^2 - f$ strongly convex when $\gamma < 1/L_f$

$$D_h(z, x) := h(z) - h(x) - \langle \nabla h(x), z - x \rangle$$

for a convex $h \in C^1(\mathbb{R}^n)$ is the **Bregman distance** induced by h

Bregman proximal algorithm(s)

Why? ~ Structured problems

$$\varphi = "L_f\text{-smooth } f" + \text{"easy-to-prox } g"$$

$$\begin{aligned} x^{k+1} &\in \text{prox}_{\gamma g}\left(x^k - \gamma \nabla f(x^k)\right) \\ &= \arg \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + D_{\hat{h}}(w, x^k) \right\} \\ &:= \text{prox}_{\varphi}^{\hat{h}}(x^k) \end{aligned}$$

FBS = Bregman PPA

with **dfg** $\hat{h} = \frac{1}{2\gamma} \|\cdot\|^2 - f$



dfg: distance-generating function (should comply with *essential smoothness, Legendreness, essential strict convexity*)

Bregman proximal algorithm(s)

Why? ~ Structured problems

$$\varphi = "L_f\text{-smooth } f" + \text{"easy-to-prox } g"$$

$$\begin{aligned} x^{k+1} &\in \text{prox}_{\gamma g}\left(x^k - \gamma \nabla f(x^k)\right) \\ &= \arg \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + D_{\hat{h}}(w, x^k) \right\} \\ &:= \text{prox}_{\varphi}^{\hat{h}}(x^k) \end{aligned}$$

FBS = Bregman PPA

with dfg $\hat{h} = \frac{1}{2\gamma} \|\cdot\|^2 - f$



✓ Same results of **PPA** $\gamma < 1/L_f$

$$\begin{cases} \frac{1}{2\gamma} \|\cdot - \cdot\|^2 & \rightarrow D_{\hat{h}} \\ \varphi^\gamma & \rightarrow \varphi^{\hat{h}} := \min_w \left\{ \varphi(w) + D_{\hat{h}}(w, \cdot) \right\} \end{cases}$$

✓ For **FBS**: take $0 < \lambda < \frac{2}{1+\gamma L_f}$

$$\varphi^{\hat{h}}(x^+) \leq \varphi^{\hat{h}}(x) - \frac{2-\lambda(1+\gamma L_f)}{\lambda(1-\gamma L_f)} D_{\hat{h}}(x^+, x)$$

✓ ($\varphi^{\hat{h}}$ is the **forward-backward envelope**)

P. Patrinos and A. Bemporad, *Proximal Newton methods for convex composite optimization*, In: 52th IEEE CDC, 2013

AT, L. Stella and P. Patrinos, *Forward-Backward Envelope for the Sum of Two Nonconvex Functions: Further Properties and Nonmonotone Linesearch Algorithms*, SIAM J Opt 28(3):2274-2303, 2018

Bregman proximal algorithm(s)

Why? ~ ...

With a suitable h

- ✓ **Uni-simplify** analysis

≈ M. Teboulle, *A simplified view of first order methods for optimization*, MathProg 170(1):67–96, Springer, 2018



Bregman proximal algorithm(s)

Why? ~ ...

With a suitable h

✓ **Uni-simplify** analysis

≈ M. Teboulle, *A simplified view of first order methods for optimization*, MathProg 170(1):67–96, Springer, 2018

- ▶ Can enforce $x^k \in \text{dom } h$: blend proximal and barrier methods

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \varphi(x) \quad \text{s.t.} \quad x \in C$$

where C is the closure of $\text{dom } h$

a lot!
(some caution if φ nonconvex)

Bregman proximal algorithm(s)

Why? ~ ...

With a suitable h

- ✓ **Uni-simplify** analysis

≈ M. Teboulle, *A simplified view of first order methods for optimization*, MathProg 170(1):67–96, Springer, 2018

- ▶ Can enforce $x^k \in \text{dom } h$: blend proximal and barrier methods

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \varphi(x) \quad \text{s.t.} \quad x \in C$$

where C is the closure of $\text{dom } h$

a lot!
(~~some~~ caution if φ nonconvex)

- ✓ Extend applicability
 - ✓ make prox tractable
 - ✓ relax assumptions

f is L_f -**relative smooth** wrt h if $L_f h \pm f$ are convex

Bregman proximal algorithm(s)

Why? ~ ...

Bregman FBS

f L_f -smooth relative to a dgf $\textcolor{brown}{h}$

$$\begin{aligned}x^{k+1} &\in \arg \min_{w \in \mathbb{R}^n} \left\{ f(x^k) + \langle \nabla f(x^k), w - x^k \rangle + g(w) + \frac{1}{\gamma} D_{\textcolor{brown}{h}}(w, x^k) \right\} \\&= \arg \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + D_{\hat{h}}(w, x^k) \right\} \\&= \text{prox}_{\varphi}^{\hat{h}}(x^k) \quad \text{where } \hat{h} := \frac{1}{\gamma} h - f\end{aligned}$$



Bregman proximal algorithm(s)

Why? ~ ...

Bregman FBS

f L_f -smooth relative to a dgf h

$$\begin{aligned}x^{k+1} &\in \arg \min_{w \in \mathbb{R}^n} \left\{ f(x^k) + \langle \nabla f(x^k), w - x^k \rangle + g(w) + \frac{1}{\gamma} D_h(w, x^k) \right\} \\&= \arg \min_{w \in \mathbb{R}^n} \left\{ \varphi(w) + D_{\hat{h}}(w, x^k) \right\} \\&= \text{prox}_{\varphi}^{\hat{h}}(x^k) \quad \text{where } \hat{h} := \frac{1}{\gamma} h - f\end{aligned}$$

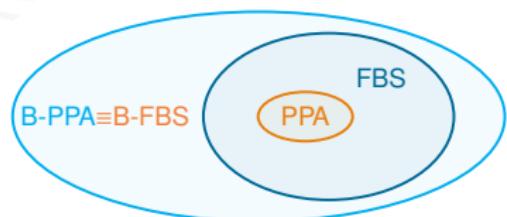
B-BFS \equiv B-PPA

f L_f -smooth relative to dgf h

$$\Rightarrow \hat{h} := \frac{1}{\gamma} h - f \text{ is a dgf}$$

(with same properties)

- ▶ $\text{dom } \hat{h} = \text{dom } h$
- ▶ h (loc.) str.cvx/smooth $\Rightarrow \hat{h}$ (loc.) str.cvx/smooth
- ▶ $D_h(x^k, y^k) \rightarrow 0 \Rightarrow D_{\hat{h}}(x^k, y^k) \rightarrow 0$
- ▶ ...



Outline

Proximal algorithm(s)

Proximal point algorithm

Simplest analysis

Bregman proximal algorithm(s)

Why Bregman

Linesearch B-FBS: the BELLA algorithm

Newton-type?

Linesearch?



Block-coordinate B-FBS: Bregman Finito/MISO

Conclusions

Linesearch B-FBS: the BELLA algorithm

Problem & objective

Problem

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \varphi(x) \quad (\mathbf{P})$$

A1. $\varphi : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ proper, lsc

A2. $\inf \varphi > -\infty$

Toolbox

Bregman proximal mapping

$$\text{prox}_{\varphi}^h(x) = \arg \min_{w \in \mathbb{R}^n} \{\varphi(w) + D_h(w, x)\} : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$$

A3. h dgf with $\text{dom } h = \mathbb{R}^n$

Goal

- Globalize Newton-type methods $x^+ = x + \tau d$
- by using only the B-PPA oracle

Linesearch B-FBS: the BELLA algorithm

Newton type?

$$x^+ = x + \tau d$$

But...

:(φ^h usually nonsmooth

:(or, anyway, $\nabla \varphi^h = \nabla^2 h [\text{id} - \text{prox}_\varphi^h]$

- ▶ even if applicable, Armijo-type linesearch would require $\nabla^2 h$
- ▶ (this was the original approach of minFBE)

Directions

Necessary optimality condition: $0 \in R_\varphi^h(x^\star) := [\text{id} - \text{prox}_\varphi^h](x^\star)$

⇒ Get d via **quasi-Newton** on R_φ^h

Requires (only)

- ✓ prox_φ^h
- ✓ direct linear algebra (scalar products)

Linesearch B-FBS: the BELLA algorithm

Linesearch?

$$x^+ = x + \tau d$$

Still...

- :(φ^h usually **nonsmooth**)
- :(or, anyway, $\nabla\varphi^h = \nabla^2 h [\text{id} - \text{prox}_\varphi^h]$)
- :(is d of “*descent*”?)

And yet...

- :) φ^h is **continuous**
- :) $\varphi^h(\bar{x}) \leq \varphi(\bar{x}) = \varphi^h(x) - D_h(\bar{x}, x)$ $(\bar{x} \in \text{prox}_\varphi^h(x))$

Linesearch B-FBS: the BELLA algorithm

Linesearch?

$$\cancel{x^+ = x + \tau d}$$

Still...

- ⌚ φ^h usually **nonsmooth**
- ⌚ or, anyway, $\nabla\varphi^h = \nabla^2 h [\text{id} - \text{prox}_\varphi^h]$
- ⌚ is d of “*descent*”?

And yet...

- 😊 φ^h is **continuous**
- 😊 $\varphi^h(\bar{x}) \leq \varphi(\bar{x}) = \varphi^h(x) - D_h(\bar{x}, x)$ $(\bar{x} \in \text{prox}_\varphi^h(x))$

Solution (works for **any** d !)

$$x^+ = (1 - \tau)\bar{x} + \tau(x + d)$$

reducing τ until

$$\varphi^h(x^+) \leq \varphi^h(x) - \sigma D_h(\bar{x}, x) \quad (\sigma \in (0, 1) \text{ fixed})$$

Linesearch B-FBS: the BELLA algorithm

The BELLA algorithm

Bregman EnveLope Linesearch Algorithm

1. $\bar{x} \in \text{prox}_\varphi^h(x)$
2. Choose $d \in \mathbb{R}^n$ and set $\tau = 1$
3. **While** $\varphi^h((1 - \tau)\bar{x} + \tau(x + d)) \geq \varphi^h(x) - \sigma D_h(\bar{x}, x)$ **do**
 $\tau \leftarrow \tau/2$
4. $x^+ = (1 - \tau)\bar{x} + \tau(x + d)$

Convergence

$$\varphi^h(x^{k+1}) \leq \varphi^h(x^k) - \sigma D_h(\bar{x}^k, x^k) \quad (\text{SD})$$

Same convergence properties of **(B)-PPA!**

Extras...

(under assumptions)

- ▶ global/R-linear with **KL**
- ▶ If directions d are “good”, $\tau = 1$ eventually always accepted

Outline

Proximal algorithm(s)

 Proximal point algorithm

 Simplest analysis

Bregman proximal algorithm(s)

 Why Bregman

Linesearch B-FBS: the BELLA algorithm

 Newton-type?

 Linesearch?



Block-coordinate B-FBS: Bregman Finito/MISO

Conclusions

Block coordinate B-FBS

Regularized finite-sum

Problem $\underset{x \in \mathbb{R}^n}{\text{minimize}} \varphi(x) \equiv \frac{1}{N} \sum_{i=1}^N f_i(x) + g(x)$ (P)

- A1. $g : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ proper, lsc
- A2. f_i L_{f_i} -smooth relative to a dgf h_i
- A3. $\inf \varphi > -\infty$

Consensus formulation

$$\underset{x \in \mathbb{R}^{nN}}{\text{minimize}} \Phi(\mathbf{x}) \equiv \underbrace{\frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x}_i)}_{F(\mathbf{x})} + \underbrace{\frac{1}{N} \sum_{i=1}^N g(\mathbf{x}_i) + \delta_{\Delta}(\mathbf{x})}_{G(\mathbf{x})} \quad (\text{P}_{\Delta})$$

$\Delta := \left\{ \mathbf{x} \in \mathbb{R}^{nN} \mid x_1 = x_2 = \dots = x_N \right\}$ is the **consensus set**

$$\delta_E(x) = \begin{cases} 0 & \text{if } x \in E, \\ \infty & \text{otherwise} \end{cases}$$

is the indicator function of the set $E \subseteq \mathbb{R}^n$

Block coordinate B-FBS

Bregman Finito/MISO

Block-coordinate B-PPA (or B-FBS)

with dgf $\hat{H} := \hat{h}_1 \times \dots \times \hat{h}_N$ $\hat{h}_i := \frac{1}{\gamma_i} h_i - f_i$, $\gamma_i \in (0, {}^N/L_{f_i})$

1. $\mathbf{u} \in \text{prox}_{\Phi}^{\hat{H}}(\mathbf{x})$ % B-PPA = B-FBS

2. select/sample $i \in \{1, \dots, N\}$

3. $\mathbf{x}_j^+ = \begin{cases} u_i & \text{if } j = i, \\ \mathbf{x}_j & \text{otherwise} \end{cases}$

$f \times g(x, y) := (f(x), f(y))$ is the separable stacking of functions f and g

Block coordinate B-FBS

Bregman Finito/MISO

Block-coordinate B-PPA (or B-FBS)

with dgf $\hat{H} := \hat{h}_1 \times \dots \times \hat{h}_N$ $\hat{h}_i := \frac{1}{\gamma_i} h_i - f_i$, $\gamma_i \in (0, N/L_{f_i})$

1. $\mathbf{u} \in \text{prox}_{\Phi}^{\hat{H}}(\mathbf{x})$ % B-PPA = B-FBS

2. select/sample $i \in \{1, \dots, N\}$

3. $\mathbf{x}_j^+ = \begin{cases} u_i & \text{if } j = i, \\ \mathbf{x}_j & \text{otherwise} \end{cases}$

Not wasteful! it is Bregman MISO/Finito⁺

1. $\mathbf{z} \in \arg \min_w \left\{ g(w) + \sum_i \frac{1}{\gamma_i} h_i(w) - \langle \tilde{s}, w \rangle \right\}$

✓ nonsmooth/nonconvex g

✓ independent stepsizes γ_i

✓ relative-smooth f_i

2. select/sample $i \in \{1, \dots, N\}$

3. update gradients table

$$s_j^+ = \begin{cases} \frac{1}{\gamma_i} \nabla h_i(\mathbf{z}) - \frac{1}{N} \nabla f_i(\mathbf{z}) & \text{if } j = i, \\ s_j & \text{otherwise} \end{cases}$$

4. update aggregated B-FBS step $\tilde{s}^+ = \tilde{s} + (s_i^+ - s_i)$

Block coordinate B-FBS

Convergence results

For ANY sampling

- ✓ Decrease on $\Phi^{\hat{H}}$
- ✓ φ coercive \Rightarrow bounded sequence

Randomized/Shuffled/Cyclic

- ✓ Sufficient decrease on $\Phi^{\hat{H}}$
 - ✓ Convergence results of **(B)-PPA** apply!
 - ✓ (including global/ R -linear with **KL**)

Extras...

- ✓ Q -linear rates with φ strongly convex (g or f_i can even be nonconvex!)

Analysis agnostic to “our” decomposition of φ !

Outline

Proximal algorithm(s)

 Proximal point algorithm

 Simplest analysis

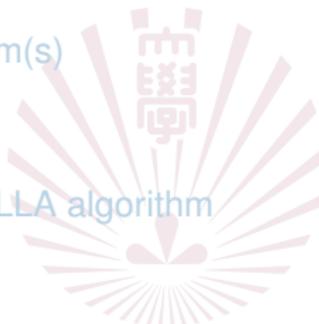
Bregman proximal algorithm(s)

 Why Bregman

Linesearch B-FBS: the BELLA algorithm

 Newton-type?

 Linesearch?



Block-coordinate B-FBS: Bregman Finito/MISO

Conclusions

Conclusions

In this talk...

What we knew

- ▶ Proximal point: elegant, simple, and powerful
- ▶ Bregman divergence: the key enhancement

What we knew but maybe didn't care

- ▶ **B-PPA = B-FBS**
- ▶ Separating tasks: $\underset{\text{analysis}}{\text{minimize } \varphi}$ **VS** $\underset{\text{algorithm}}{\text{minimize } f + g \dots \text{ s.t. } \dots}$

What maybe we didn't know

- ▶ Continuity + sufficient decrease \rightarrow globalize fast local methods
😊 the **BELLA** algorithm
- ▶ Block-coordinate proximal point \rightarrow incremental proximal algorithms
😊 **Bregman MISO/Finito⁺**

M. Ahookhosh, AT and P. Patrinos, *A Bregman Forward-Backward Linesearch Algorithm for Nonconvex Composite Optimization: Superlinear Convergence to Nonisolated Local Minima*, SIAM J Opt 31(1):653-685, 2021

P. Latafat, AT, M. Ahookhosh and P. Patrinos, *Bregman Finito/MISO for nonconvex regularized finite sum minimization without Lipschitz gradient continuity*, arXiv:2102.10312, 2021