Splitting Envelopes Accelerated Second-order Proximal Methods

Panos Patrinos (joint work with Lorenzo Stella, Alberto Bemporad)

September 8, 2014



Outline

- ✓ forward-backward envelope (FBE)
- ✓ forward-backward Newton method (FBN)
- ✓ dual FBE and Augmented Lagrangian
- ✓ alternating minimization Newton method (AMNM)
- ✓ Douglas Rachford envelope (DRE)
- ✓ accelerated Douglas Rachford splitting (ADRS)

based on

- P. Patrinos and A. Bemporad. Proximal Newton methods for convex composite optimization. In Proc. 52nd IEEE Conference on Decision and Control (CDC), pages 2358-2363, Florence, Italy, 2013.
- P. Patrinos, L. Stella, and A. Bemporad, Forward-backward truncated Newton methods for convex composite optimization. submitted, arXiv:1402.6655, 2014.
- P. Patrinos, L. Stella, and A. Bemporad. Douglas-Rachford splitting: complexity estimates and accelerated variants. In Proc. 53rd IEEE Conference on Decision and Control (CDC), Los Angeles, CA, arXiv:1407.6723, 2014.
- 4. L. Stella, P. Patrinos, and A. Bemporad, Alternating minimization Newton method for separable convex optimization, 2014 (submitted).

fixed point implementation for MPC

 A. Guiggiani, P. Patrinos, and A. Bemporad. Fixed-point implementation of a proximal Newton method for embedded model predictive control. In 19th IFAC, South Africa, 2014.

Convex composite optimization

minimize F(x) = f(x) + g(x)

 $\checkmark f: \mathbb{R}^n \to \mathbb{R}$ convex, twice continuously differentiable with

 $\|\nabla f(x) - \nabla f(y)\| \le L_f \|x - y\|, \text{ for all } x, y \in \mathbb{R}^n$

 $\checkmark~g:{\rm I\!R}^n\to \overline{\rm I\!R}$ convex, nonsmooth with inexpensive proximal mapping

$$\operatorname{prox}_{\gamma g}(x) = \operatorname*{arg\,min}_{z \in \mathbb{R}^n} \left\{ g(z) + \frac{1}{2\gamma} \|z - x\|^2 \right\}$$

- ✓ many problem classes: QPs, cone programs, sparse least-squares, rank minimization, total variation minimization,...
- ✓ applications: control, system identification, signal processing, image analysis, machine learning,...

Proximal mappings

$$\operatorname{prox}_{\gamma g}(x) = \operatorname*{arg\,min}_{z \in \mathbb{R}^n} \left\{ g(z) + \frac{1}{2\gamma} \|z - x\|^2 \right\} \qquad \gamma > 0$$

 \checkmark resolvent of maximal monotone operator ∂g

$$\operatorname{prox}_{\gamma g}(x) = (I + \gamma \partial g)^{-1}(x)$$

- ✓ single-valued and (firmly) **nonexpansive**
- explicitly computable for many functions (see Parikh, Boyd '14, Combettes, Pesquet '10)
- \checkmark reduces to **projection** when g is indicator of convex set

$$\operatorname{prox}_{\gamma\delta_C}(x) = \Pi_C(x)$$

✓ $z = \text{prox}_{\gamma g}(x)$ is implicit a subgradient step $(0 \in \partial g(z) + \gamma^{-1}(z-x))$

$$z = x - \gamma v$$
 $v \in \partial g(z)$

Proximal Minimization Algorithm

minimize g(x), $g: \mathbb{R}^n \to \overline{\mathbb{R}}$ closed proper convex

given $x^0 \in {\rm I\!R}^n$, repeat

$$x^{k+1} = \operatorname{prox}_{\gamma g}(x^k) \qquad \gamma > 0$$

✓ fixed point iteration for optimality conditions

$$0 \in \partial g(x^{\star}) \iff x^{\star} \in (I + \gamma \partial g)(x^{\star}) \iff x^{\star} = \operatorname{prox}_{\gamma g}(x^{\star})$$

- ✓ special case of proximal point algorithm (Martinet '70, Rockafellar '76)
- $\checkmark\,$ converges under very general conditions
- ✓ mostly conceptual algorithm

Moreau envelope

Moreau envelope of closed proper convex $g: \mathbb{R}^n \to \overline{\mathbb{R}}$

$$g^{\gamma}(x) = \inf_{z \in \mathbb{R}^n} \left\{ g(z) + \frac{1}{2\gamma} \|z - x\|^2 \right\}, \qquad \gamma > 0$$

✓ g^{γ} is real-valued, convex, differentiable with $1/\gamma$ -Lipschitz gradient

$$\nabla g^{\gamma}(x) = (1/\gamma)(x - \operatorname{prox}_{\gamma g}(x))$$

 \checkmark minimizing nonsmooth g is equivalent to minimizing smooth g^γ

 \checkmark proximal minimization algorithm = gradient method for g^{γ}

$$x^{k+1} = x^k - \gamma \nabla g^{\gamma}(x^k)$$

 $\checkmark\,$ can use any method of unconstrained smooth minimization for g^γ

Forward-Backward Splitting (FBS)

minimize
$$F(x) = f(x) + g(x)$$

 \checkmark optimality condition: $x_{\star} \in \mathbb{R}^n$ is optimal if and only if

$$x_{\star} = \operatorname{prox}_{\gamma g}(x_{\star} - \gamma \nabla f(x_{\star})), \qquad \gamma > 0$$

forward-backward splitting (aka proximal gradient)

$$x^{k+1} = \operatorname{prox}_{\gamma g}(x^k - \gamma \nabla f(x^k)), \qquad \gamma \in (0, 2/L_f)$$

FBS is a fixed point iteration

- ✓ g = 0: gradient method, $g = \delta_C$: gradient projection, f = 0: prox min
- ✓ accelerated versions (Nesterov)

Forward-Backward Envelope

$$x - \operatorname{prox}_{\gamma g}(x - \gamma \nabla f(x)) = 0$$

$$\checkmark \text{ use } \operatorname{prox}_{\gamma g}(y) = y - \gamma \nabla g^{\gamma}(y) \text{ for } y = x - \gamma \nabla f(x)$$
$$\gamma \nabla f(x) + \gamma \nabla g^{\gamma}(x - \gamma \nabla f(x)) = 0$$

✓ multiply with $\gamma^{-1}(I - \gamma \nabla^2 f(x))$ (positive definite for $\gamma \in (0, 1/L_f)$)

✓ gradient of the Forward Backward Envelope (FBE)

$$F_{\gamma}^{\mathrm{FB}}(x) = f(x) - \frac{\gamma}{2} \|\nabla f(x)\|_{2}^{2} + g^{\gamma}(x - \gamma \nabla f(x))$$

 \checkmark alternative expression for FBE

$$F_{\gamma}^{\mathrm{FB}}(x) = \inf_{z \in \mathbb{R}^n} \{\underbrace{f(x) + \langle \nabla f(x), z - x \rangle}_{\text{linearize } f \text{ around } x} + g(z) + \frac{1}{2\gamma} \|z - x\|^2 \}$$

Properties of FBE

- ✓ stationary points of $F_{\gamma}^{\rm FB} =$ minimizers of F
- ✓ reformulates original nonsmooth problem into a smooth one

$$\underset{x \in \mathbb{R}^n}{\operatorname{minimize}} \ F_{\gamma}^{\operatorname{FB}}(x) \qquad \text{equivalent to} \qquad \underset{x \in \mathbb{R}^n}{\operatorname{minimize}} \ F(x)$$

 \checkmark F_{γ}^{FB} is real-valued, continuously differentiable

$$\nabla F_{\gamma}^{\text{FB}}(x) = \gamma^{-1} (I - \gamma \nabla^2 f(x)) (x - \text{prox}_{\gamma g}(x - \gamma \nabla f(x)))$$

✓ FBS is a variable metric gradient method for FBE

$$x^{k+1} = x^k - \gamma D_k^{-1} \nabla F_{\gamma}^{\text{FB}}(x^k)$$





Forward-Backward Newton Method (FBN)

Input: $x^0 \in \mathbb{R}^n$, $\gamma \in (0, 1/L_f)$, $\sigma \in (0, 1/2)$ for k = 0, 1, 2, ... do

Newton direction

Choose $H^k \in \hat{\partial}^2 F^{\text{FB}}_{\gamma}(x^k)$. Compute d^k by solving (approximately)

$$H^k d = -\nabla F_{\gamma}^{\rm FB}(x^k)$$

Line search

Compute stepsize by backtracking

end

$$F_{\gamma}^{\rm FB}(x^k+\tau_k d^k) \leq F_{\gamma}^{\rm FB}(x^k) + \sigma \tau_k \langle \nabla F_{\gamma}^{\rm FB}(x^k), d^k \rangle$$

Jpdate: $x^{k+1} = x^k + \tau_k d^k$

Linear Newton approximation

 $\mathsf{FBE} \text{ is } C^1 \text{ but not } C^2$

$$Hd=-
abla F_{\gamma}^{\mathrm{FB}}(x)$$
 where

$$\nabla F_{\gamma}^{\text{FB}}(x) = \gamma^{-1} (I - \gamma \nabla^2 f(x)) (x - \text{prox}_{\gamma g}(x - \gamma \nabla f(x)))$$

and $\hat{\partial}^2 F_{\gamma}(x)$ is an approximate generalized Hessian

$$H = \gamma^{-1} (I - \gamma \nabla^2 f(x)) (I - \mathbf{P} (I - \gamma \nabla^2 f(x))) \in \hat{\partial}^2 F_{\gamma}(x)$$

where $P \in \underbrace{\partial_C(\operatorname{prox}_{\gamma g})}_{\operatorname{Clarke's generalized}} (x - \gamma \nabla f(x))$

✓ preserves all favorable properties of the Hessian for C^2 functions ✓ "Gauss-Newton" generalized Hessian: we omit 3rd order terms

Generalized Jacobians of proximal mappings

✓ $\partial_C \operatorname{prox}_{\gamma g}(x)$ is the following set of matrices (Clarke, 1983)

 $\operatorname{conv} \left\{ \begin{array}{l} \text{limits of (ordinary) Jacobians for every sequence that converges} \\ \text{to } x, \text{ consisting of points where } \operatorname{prox}_{\gamma g} \text{ is differentiable} \end{array} \right\}$

▶ $\operatorname{prox}_{\gamma g}(x)$ simple to compute $\Longrightarrow P \in \partial_C(\operatorname{prox}_{\gamma g})(x)$ for free

▶ g (block) separable $\implies P \in \partial_C(\operatorname{prox}_{\gamma g})(x)$ (block) diagonal

example– ℓ_1 **norm** more examples in Patrinos, Stella, Bemporad (2014)

$$\checkmark g(x) = \|x\|_1 \qquad \operatorname{prox}_{\gamma f}(x)_i = \begin{cases} x_i + \gamma, & \text{if } x_i \leq -\gamma, \\ 0, & \text{if } -\gamma \leq x_i \leq \gamma \\ x_i - \gamma, & \text{if } x_i \geq \gamma \end{cases}$$

✓ $P \in \partial_C(\mathrm{prox}_{\gamma g})(x)$ are diagonal matrices with

$$P_{ii} = \begin{cases} 1, & \text{if } i \in \{i \mid |x_i| > \gamma\}, \\ \in [0, 1], & \text{if } i \in \{i \mid |x_i| = \gamma\}, \\ 0, & \text{if } i \in \{i \mid |x_i| < \gamma\}. \end{cases}$$

Convergence of FBN

 \checkmark every limit point of $\{x^k\}$ converges to $\mathop{\rm arg\,min}_{x\in {\rm I\!R}^n} F(x)$

✓ all $H \in \hat{\partial}^2 F_{\gamma}(x_{\star})$ nonsingular \implies **Q-quadratic asymptotic rate**

extension: FBN II

- ✓ apply FB step after a Newton step
- ✓ same asymptotic rate + global complexity estimates

▶ non-strongly convex f: sublinear rate for $F(x^k) - F(x_\star)$

► strongly convex
$$f$$
: linear rate for $\begin{cases} F(x^k) - F(x_\star) \\ \|x^k - x_\star\|^2 \end{cases}$

FBN-CG

large problems

conjugate gradient (CG) on regularized Newton system until

$$\|\underbrace{(H^k + \delta_k I)d^k + \nabla F_{\gamma}^{\mathrm{FB}}(x^k)}_{\text{residual}}\| \le \eta_k \|\nabla F_{\gamma}^{\mathrm{FB}}(x^k)\|$$

with $\eta_{\pmb{k}} = O(\|\nabla F_{\gamma}^{\mathrm{FB}}(x^k)\|), \ \delta_{\pmb{k}} = O(\|\nabla F_{\gamma}^{\mathrm{FB}}(x^k)\|)$

properties

✓ no need to form $\nabla^2 f(x)$ and H^k – only matvec products

✓ same convergence properties

Box-constrained convex programs

 $\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & \ell \leq x \leq u \end{array}$

Newton direction solves

minimize
$$\frac{1}{2}\langle d, \nabla^2 f(x^k)d \rangle + \langle \nabla f(x^k), d \rangle$$

subject to $d_i = \ell_i - x_i^k, i \in \beta_1, \quad d_i = u_i - x_i^k, i \in \beta_2$

where

$$\begin{split} \beta_1 &= \{i \mid x_i^k - \gamma \nabla_i f(x^k) \leq \ell_i\} & \text{estimate of } x_i^\star = \ell_i \\ \beta_2 &= \{i \mid x_i^k - \gamma \nabla_i f(x^k) \geq u_i\} & \text{estimate of } x_i^\star = u_i \end{split}$$

Newton system becomes

$$Q_{\delta\delta}d_{\delta} = -(\nabla_{\delta}f(x^{k}) + \nabla_{\delta\beta}f(x^{k})d_{\beta}), \quad (\beta = \beta_{1} \cup \beta_{2}, \ \delta \setminus \beta = [n])$$

Example



FBN: much less sensitive to bad conditioning

Sparse least-squares

minimize
$$\frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1$$

✓ Newton system becomes

$$d_{\beta} = -x_{\beta}$$
$$A_{\cdot\delta}^{\top}A_{\cdot\delta}d_{\delta} = -[A_{\cdot\delta}^{\top}(A_{\cdot\delta}x_{\delta} - b) + \lambda\operatorname{sign}(x_{\delta} - \gamma\nabla_{\delta}f(x))]$$

 \checkmark δ is an estimate of the nonzero components of x_{\star}

$$\delta = \{i \mid |x_i - \gamma \nabla_i f(x)| > \lambda \gamma \}$$

 \checkmark close to solution δ small

FBN Methods are robust

✓ 548 datasets taken from

http://wwwopt.mathematik.tu-darmstadt.de/spear/

- ✓ compared against AFBS, YALL1 (ADDM-based), SpaRSA (Barzilai-Borwein), I1-Is (interior point)
- ✓ performance plot: point (x, y) indicates algorithm is at most x times slower in a fraction of y problems



Sparse Logistic Regression



Augmented Lagrangians and Moreau envelopes

 $\begin{array}{ll} \text{minimize} & f(x) & f: \mathbb{R}^n \to \overline{\mathbb{R}} \text{ convex (can be nonsmooth)}, \\ \text{subject to} & Ax = b & A \in \mathbb{R}^{p \times n} \end{array}$

Augmented Lagrangian

$$L_{\gamma}(x,y) = f(x) + \langle y, Ax - b \rangle + \frac{\gamma}{2} ||Ax - b||^2$$

Augmented Lagrangian method (ALM or Method of Multipliers)

$$\begin{aligned} x^k &= \inf_{x \in \mathbb{R}^n} L_{\gamma}(x, y^k) & \text{Hestenes (1969), Powell (1969)} \\ y^{k+1} &= y^k + \gamma(Ax^k - b) \end{aligned}$$

ALM = proximal minimization for dual Rockafellar (1973,1976) = gradient method for Moreau envelope of dual

Separable convex problems

 $\begin{array}{ll} \text{minimize} & f(x) + g(z), & f,g \text{ convex (can be nonsmooth)} \\ \text{subject to} & Ax + Bz = b & A \in \mathbb{R}^{p \times n}, B \in \mathbb{R}^{p \times m} \end{array}$

 \checkmark f is strongly convex with convexity parameter μ_f

- $\checkmark~f,~g$ are nice, e.g. separable. coupling introduced through constraints
- \checkmark ALM not amenable to decomposition: x and z updates are coupled

Alternating Minimization Method (AMM): FBS applied to the dual

$$\begin{aligned} x^{k+1} &= \underset{x \in \mathbb{R}^n}{\arg \min} \ L_0(x, z^k, y^k) \\ z^{k+1} &= \underset{z \in \mathbb{R}^m}{\arg \min} \ L_\gamma(x^{k+1}, z^k, y^k) \qquad \gamma \in (0, 2\mu_f / \|A\|^2) \\ y^{k+1} &= y^k + \gamma(Ax^{k+1} + Bz^{k+1} - b) \end{aligned}$$

Augmented Lagrangian

 $L_{\gamma}(x,z,\lambda) = f(x) + g(z) + \langle y, Ax + Bz - b \rangle + \frac{\gamma}{2} \|Ax + Bz - b\|^2$

Dual FBE

FBE for dual problem is augmented Lagrangian!

$$F_{\gamma}^{\mathrm{FB}}(\boldsymbol{y}) = L_{\gamma}(\boldsymbol{x}(\boldsymbol{y}), \boldsymbol{z}(\boldsymbol{y}), \boldsymbol{y})$$

where x(y), z(y) are the **AMM updates**

$$\begin{aligned} x(y) &= \underset{x \in \mathbb{R}^n}{\arg\min} \ f(x) + \langle y, Ax \rangle \\ z(y) &= \underset{z \in \mathbb{R}^m}{\arg\min} \ g(z) + \langle y, Bz \rangle + \frac{\gamma}{2} \|Ax(y) + Bz - b\|^2 \end{aligned}$$

Connection between AMM and FBE

dual problem is equivalent to

$$\underset{y \in \mathbb{R}^p}{\text{maximize }} F_{\gamma}^{\text{FB}}(y) = L_{\gamma}(x(y), z(y), y) \qquad \gamma \in \left(0, \mu_f / \|A\|^2\right)$$

$$\checkmark f \in C^2(\mathbb{R}^n) \implies F_{\gamma}^{\mathrm{FB}} \in C^1(\mathbb{R}^p)$$
$$\nabla F_{\gamma}^{\mathrm{FB}}(y) = \overbrace{\left(I - \gamma A(\nabla^2 f(x(y)))^{-1} A^{\top}\right)}^{D(y)} (Ax(y) + Bz(y) - b)$$

AMM as a variable metric gradient method on dual FBE

$$y^{k+1} = y^k + \gamma D(y^k)^{-1} \nabla F_{\gamma}^{\mathrm{FB}}(\boldsymbol{y^k})$$

✓ AMNM: FBN to dual

Strictly convex QPs

minimize $\frac{1}{2}\langle x, Qx \rangle + \langle q, x \rangle$ $\operatorname{cond}(Q) = 10^4$ subject to $\ell \leq Ax \leq u$ $A \in \mathbb{R}^{2000 \times 1000}$



Projection onto Convex Sets

minimize
$$\frac{1}{2} ||x - p||^2 + \sum_{i=1}^m \delta_{C_i}(z_i)$$

subject to $x = z_i, \quad i = 1, \dots, M$

✓ δ_C is the **indicator** of *C* ✓ x^* is the **projection** of *p* onto $C_1 \cap \ldots \cap C_m$



Random problem with m=100 random hyperplanes in ${\rm I\!R}^{120}$



Distributed MPC

$$\begin{array}{ll} \text{minimize} & \sum_{i=1}^{M} \sum_{t=1}^{N-1} \left(\|\xi_i(t)\|_{Q_i}^2 + \|u_i(t)\|_{R_i}^2 \right) + \|\xi_i(N)\|_{P_i}^2 \\ \text{subject to} & \xi_i(0) = \xi_i^0, & i \in \mathbb{N}_{[1,M]} \\ & \xi_i(t+1) = \sum_{j \in \mathcal{N}_i} \Phi_{ij}\xi_j(t) + \Gamma_{ij}u_j(t), \ t \in \mathbb{N}_{[0,N-1]}, i \in \mathbb{N}_{[1,M]} \\ & (\xi_i(t), u_i(t)) \in Y_i, & t \in \mathbb{N}_{[0,N-1]}, i \in \mathbb{N}_{[1,M]} \\ & \xi_i(N) \in Z_i, & i \in \mathbb{N}_{[1,M]} \end{array}$$

✓ solve an optimal control problem over a network of agents
 ✓ local constraint sets Y_i, Z_i are simple, coupling in the dynamics
 ✓ can handle complicated local and coupled constraints as well

DMPC simulations

 $M=100 \ {\rm subsystems} \\ {\rm 23600 \ vars, \ 20600 \ cons}$

average over 20 instances 8 states, 3 inputs, N = 10



	FAMM	AMNM	
M	local	local	global
5	29.5	2.1	2.1
10	47.0	2.2	2.1
20	65.7	2.4	2.4
50	104.8	3.3	3.3
100	139.2	3.9	3.8
200	159.3	4.4	4.3

communication rounds (in thousands)

Douglas-Rachford Splitting

minimize F(x) = f(x) + g(x)

✓ optimality condition: x_{\star} is optimal if and only if $x_{\star} = \text{prox}_{\gamma f}(\tilde{x})$, where \tilde{x} solves

$$\operatorname{prox}_{\gamma g}(2\operatorname{prox}_{\gamma f}(x) - x) - \operatorname{prox}_{\gamma f}(x) = 0$$

🗸 DRS

$$y^{k} = \operatorname{prox}_{\gamma f}(x^{k})$$
$$z^{k} = \operatorname{prox}_{\gamma g}(2y^{k} - x^{k})$$
$$x^{k+1} = x^{k} + \lambda_{k}(z^{k} - y^{k})$$

✓ $\gamma > 0$ and $\lambda_k \in [0, 2]$ with $\sum_{k \in \mathbb{N}} \lambda_k (2 - \lambda_k) = +\infty$ ✓ DRS is a relaxed fixed point iteration

ADMM

 $\begin{array}{ll} \text{minimize} & f(x) + g(z), & f,g \text{ convex (can be nonsmooth)} \\ \text{subject to} & Ax + Bz = b & A \in \mathbb{R}^{p \times n}, B \in \mathbb{R}^{p \times m} \end{array}$

Alternating Direction Method of Multipliers

$$\begin{aligned} x^{k+1} &= \underset{x \in \mathbb{R}^n}{\arg \min} \ L_{\gamma}(x, z^k, y^k) \\ z^{k+1} &= \underset{z \in \mathbb{R}^m}{\arg \min} \ L_{\gamma}(x^{k+1}, z, y^k) \\ y^{k+1} &= y^k + \gamma (Ax^{k+1} + Bz^{k+1} - b) \end{aligned}$$

DRS applied to the dual (Eckstein, Bertsekas, 1992)

Douglas Rachford Envelope

assume f convex, $C^2 \implies$ Moreau envelope f^γ is C^2

$$\|\nabla f(x) - \nabla f(y)\| \le L_f \|x - y\|$$
 for all $x, y \in \mathbb{R}^n$

 \checkmark optimality condition

$$\begin{aligned} &\operatorname{prox}_{\gamma f}(x) - \operatorname{prox}_{\gamma g}(2\operatorname{prox}_{\gamma f}(x) - x) = 0 \\ \checkmark \quad &\operatorname{use} \, \nabla h^{\gamma}(x) = \gamma^{-1}(x - \operatorname{prox}_{\gamma h}(x)) \\ &\nabla f^{\gamma}(x) + \nabla g^{\gamma}(x - 2\gamma \nabla f^{\gamma}(x)) = 0 \end{aligned}$$

✓ multiply by $(I - 2\gamma \nabla^2 f^{\gamma}(x))$, $\gamma \in (0, 1/L_f)$ and "integrate" Douglas Rachford Envelope (DRE)

$$F_{\gamma}^{\mathrm{DR}}(x) = f^{\gamma}(x) - \gamma \|\nabla f^{\gamma}(x)\|^{2} + g^{\gamma}(x - 2\gamma \nabla f^{\gamma}(x))$$

Properties of DRE

✓ $\gamma \in (0, 1/L_f)$: miniziming DRE = minimizing nonsmooth F

 $\inf F(x) = \inf F_{\gamma}^{\text{DR}}(x)$ $\arg\min F(x) = \operatorname{prox}_{\gamma f}(\arg\min F_{\gamma}^{\text{DR}}(x))$

 $\checkmark \ f \in C^2({\rm I\!R}^n) \implies {\rm DRE} \text{ is } C^1 \text{ on } {\rm I\!R}^n$

✓ f quadratic \implies DRE is convex with $(1/\gamma)$ -Lipschitz gradient

Connection between DRE and FBE

 \checkmark partial linearization of F around x

$$\ell_F(z;x) = f(x) + \langle \nabla f(x), z - x \rangle + g(z)$$

🗸 FBE is

$$F_{\gamma}^{\mathrm{FB}}(x) = \min_{z \in \mathbb{R}^n} \left\{ \ell_F(z; x) + \frac{1}{2\gamma} \|z - x\|^2 \right\}$$

✓ DRE is [use $\nabla f^{\gamma}(x) = \gamma^{-1}(x - \operatorname{prox}_{\gamma f}(x)) = \nabla f(\operatorname{prox}_{\gamma f}(x))$]

$$F_{\gamma}^{\mathrm{DR}}(x) = \min_{z \in \mathbb{R}^n} \left\{ \ell_F(z; \operatorname{prox}_{\gamma f}(x)) + \frac{1}{2\gamma} \|z - \operatorname{prox}_{\gamma f}(x)\|^2 \right\}$$

✓ DRE is equal to FBE evaluated at $prox_{\gamma f}(x)$

$$F_{\gamma}^{\mathrm{DR}}(x) = F_{\gamma}^{\mathrm{FB}}(\mathrm{prox}_{\gamma f}(x))$$

Connection between DRS and FBS

✓ FBS (with relaxation)

$$x^{k+1} = x^k + \lambda_k \left(\operatorname{prox}_{\gamma g} \left(x^k - \gamma \nabla f(x^k) \right) - x^k \right)$$

🗸 DRS

$$y^{k} = \operatorname{prox}_{\gamma f}(x^{k})$$
$$x^{k+1} = x^{k} + \lambda_{k} \left(\operatorname{prox}_{\gamma g}(2y^{k} - x^{k}) - y^{k} \right)$$

✓
$$f \in C^1(\mathbb{R}^n)$$

 $y^k = \operatorname{prox}_{\gamma f}(x^k) = x^k - \gamma \nabla f(\operatorname{prox}_{\gamma f}(x^k))$

DRS becomes

$$y^{k} = \operatorname{prox}_{\gamma f}(x^{k})$$
$$x^{k+1} = x^{k} + \lambda_{k} \left(\operatorname{prox}_{\gamma g} \left(y^{k} - \gamma \nabla f(y^{k}) \right) - y^{k} \right)$$

✓ **DRS** iteration = **FBS** iteration at "shifted" point $y^k = prox_{\gamma f}(x^k)$

DRS as a variable metric method

DRS

$$x^{k+1} = x^k + \lambda_k \left(\operatorname{prox}_{\gamma g}(2\operatorname{prox}_{\gamma f}(x^k) - x^k) - \operatorname{prox}_{\gamma f}(x^k) \right)$$

gradient of DRE

$$\nabla F_{\gamma}^{\mathrm{DR}}(x) = \left(I - 2\gamma \nabla^2 f^{\gamma}(x^k)\right) \left(\operatorname{prox}_{\gamma f}(x^k) - \operatorname{prox}_{\gamma g}(2\operatorname{prox}_{\gamma f}(x^k) - x^k)\right)$$

DRS: variable metric method applied to DRE

$$x^{k+1} = x^k - \lambda_k D^k \nabla F_{\gamma}^{\mathrm{DR}}(x^k)$$

where $D^k = (I-2\gamma \nabla^2 f^\gamma(x^k))^{-1}$

- \checkmark relaxation parameter λ_k of DRS \implies stepsize for gradient method
- \checkmark can use backtracking for selecting λ_k

DRS – complexity estimates

✓ assume f quadratic \implies DRE is convex for $\gamma \in (0, 1/L_f)$

DRS is preconditioned gradient method under change of variables

$$x = Sw \qquad \qquad S = D^{1/2}$$

✓ convergence rate of DRS $\lambda_k = \lambda = (1 - \gamma L_f)/(1 + \gamma L_f)$

$$F(z^{k+1}) - F_{\star} \le \frac{1}{(2\gamma\lambda)k} \|x^0 - \tilde{x}\|^2$$

 \checkmark optimal prox-parameter γ for DRE

$$\gamma_{\star} = \frac{\sqrt{2} - 1}{L_f}$$

 \checkmark linear convergence rate if F is strongly convex

Accelerated DRS

- ✓ DRE is **convex** with $(1/\gamma)$ -Lipschitz gradient for f **quadratic** and $\gamma \in (0, 1/L_f)$
- ✓ Nesterov's FGM applied to (preconditioned) DRE: $x^0 = x^{-1} \in {\rm I\!R}^n$

$$u^{k} = x^{k} + \beta_{k}(x^{k} - x^{k-1})$$
$$y^{k} = \operatorname{prox}_{\gamma f}(u^{k})$$
$$z^{k} = \operatorname{prox}_{\gamma g}(2y^{k} - u^{k})$$
$$x^{k+1} = u^{k} + \lambda(z^{k} - y^{k})$$

✓ $\lambda = (1 - \gamma L_f)/(1 + \gamma L_f)$. can choose $\beta_k = \frac{k-1}{k+2}$ ✓ convergence rate is $O(1/k^2)$

$$F(z^k) - F_\star \le \frac{4}{\gamma\lambda(k+2)^2} \|x^0 - \tilde{x}\|^2.$$

 \checkmark linear convergence for f or g strongly convex

Sparse least-squares

minimize
$$\frac{1}{2} ||Ax - b||_2^2 + \lambda ||x||_1$$
,



Take home message I

- \checkmark proximal minimization = gradient method on Moreau envelope (70's)
- \checkmark **FBS**= (variable metric) gradient method on**FBE**(this talk)
- \checkmark **DRS** = (variable metric) gradient method on **DRE** (this talk)

Take home message II

\checkmark ALM = proximal minimization for dual

Moreau envelope of dual =
$$\underset{x \in \mathbb{R}^n}{\operatorname{arg min}} L_{\gamma}(x, y)$$

(Rockafellar, 1973)

 \checkmark AMM = FBS for dual

FBE of dual
$$= L_{\gamma}(x(y), z(y), y)$$

where

$$x(y) = \operatorname*{arg\,min}_{x \in \mathbb{R}^n} L_0(x, z, y) \qquad z(y) = \operatorname*{arg\,min}_{z \in \mathbb{R}^n} L_\gamma(x(y), z, y)$$

to conclude

- \checkmark interpretation of operator splitting algorithms as gradient methods
- ✓ splitting envelopes can lead to new exciting algorithms
- $\checkmark\,$ examples: FBN, AMNM and ADRS