A VU-Point of View of Nonsmooth Optimization

Claudia Sagastizábal IMECC-UNICAMP BRAZIL, adjunct researcher sagastiz@unicamp.br

ENSTA, January 14th and 15th, 2019

(supported by ENSTA, PGMO, and Fondation X)

Program

1. Yesterday morning: Introduction to nonsmooth convex optimization

- 2. Yesterday afternoon: Models and the proximal point algorithm
- 3. Today morning: Bundle methods and the Moreau-Yosida regularization
- 4. Today afternoon: Beyond first order: VU-decomposition methods

Bibliography and one comment

References in the conference paper

https://eta.impa.br/dl/064.pdf



Proc. Int. Cong. of Math. – 2018 Rio de Janeiro, Vol. 3 (3785–3806)

A UU-POINT OF VIEW OF NONSMOOTH OPTIMIZATION

Claudia Sagastizábal



Bibliography and one comment

References in the conference paper

https://eta.impa.br/dl/064.pdf



Proc. Int. Cong. of Math. – 2018 Rio de Janeiro, Vol. 3 (3785–3806)

A UU-POINT OF VIEW OF NONSMOOTH OPTIMIZATION

Claudia Sagastizábal



Also the forthcoming book chapter

Beyond first order: \mathcal{VU} -decomposition methods

Shuai Liu, Claudia Sagastizábal

Bibliography and one comment

References in the conference paper

https://eta.impa.br/dl/064.pdf



Proc. Int. Cong. of Math. – 2018 Rio de Janeiro, Vol. 3 (3785–3806)

A UU-POINT OF VIEW OF NONSMOOTH OPTIMIZATION

Claudia Sagastizábal



Also the forthcoming book chapter

Beyond first order: \mathcal{VU} -decomposition methods

Shuai Liu, Claudia Sagastizábal

Beware that requiring the knowledge of the Lipschitz constant of a function is not far from require to know the full subdifferential

$$x^{k+1}= p^f_{t_k}(x^k)=x^k-t_kG^{k+1}$$
 with $G^{k+1}\in \partial_{arepsilon_{k+1}}f(x^k)$

$$x^{k+1}=p^f_{t_k}(x^k)=x^k-t_kG^{k+1}$$
 with $G^{k+1}\in\partial_{arepsilon_{k+1}}f(x^k)$

(DR) is satisfied at all iterations (no hats needed)

$$x^{k+1}= p^f_{t_k}(x^k)=x^k-t_kG^{k+1}$$
 with $G^{k+1}\in \partial_{arepsilon_{k+1}}f(x^k)$

(DR) is satisfied at all iterations (no hats needed) may not be implementable

► Bundle method as inexact PPA $x^{k+1} = p_{t_k}^{\mathbf{M}_k}(\hat{x}^k) = \hat{x}^k - t_k G^k$ with $G^k \in \partial_{\varepsilon_k} f(\hat{x}^k)$

$$x^{k+1}= p^f_{t_k}(x^k)=x^k-t_kG^{k+1}$$
 with $G^{k+1}\in \partial_{arepsilon_{k+1}}f(x^k)$

(DR) is satisfied at all iterations (no hats needed) may not be implementable

► Bundle method as inexact PPA $x^{k+1} = p_{t_k}^{\mathbf{M}_k}(\hat{x}^k) = \hat{x}^k - t_k G^k$ with $G^k \in \partial_{\varepsilon_k} f(\hat{x}^k)$ $\hat{x}^{k+1} = \hat{x}^k$ whenever (DR) is satisfied

Bundle methods

- **0** Choose x^1 , $t_1 > 0$, and set $\hat{x}^1 = x^1$, k = 1.
- 1 Given \hat{x}^k , \mathbf{M}_k and t_k , compute $x^{k+1} = \arg\min \mathbf{M}_k(x) + \frac{1}{2t_k} ||x - \hat{x}^k||^2$ and

$$\delta_{k+1} = f(\hat{x}^k) - \mathbf{M}_k(x^{k+1}) = \varepsilon^k + t_k \|G^k\|^2$$

- **2** Call the oracle at x^{k+1} . If $\delta_{k+1} \leq tol$ **STOP**
- 3 (Descent Rule)

$$f(x^{k+1}) \le f(\hat{x}^k) - m\delta_k? \quad \begin{cases} \text{ yes } \\ \text{ no } \end{cases} \frac{\mathsf{SS:} \hat{x}^{k+1} = x^{k+1}}{\mathsf{NS:} \hat{x}^{k+1} = \hat{x}^k}$$

4 Choose a new model : $f(\cdot) \ge M_{k+1}(\cdot)$ and $M_{k+1}(\cdot) \ge \max\left(M_k(x^{k+1}) + G^{k\top}(\cdot - x^{k+1}), f^{k+1} + g^{k+1\top}(\cdot - x^{k+1})\right)$ Choose a bounded below stepsize t_{k+1} that, if NS, is nondecreasing

5 Set
$$k = k + 1$$
, loop to 1.

Relatives in the family of bundle methods

- **Proximal BM**, parameter t_k
- **•** Trust-region BM, parameter Δ_k
- **Level BM**, parameter ℓ_k

Relatives in the family of bundle methods

- **Proximal BM**, parameter t_k
- **•** Trust-region BM, parameter Δ_k
- ► Level BM, parameter ℓ_k

Using (DR) and the theory in [CL93], we showed convergence for the proximal family

Relatives in the family of bundle methods

- **Proximal BM**, parameter t_k
- **•** Trust-region BM, parameter Δ_k
- ► Level BM, parameter ℓ_k

Using (DR) and the theory in [CL93], we showed convergence for the proximal family

What about speed?

Springer-Verlag 1999

Digital Object Identifier (DOI) 10.1007/s101079900088 Stephen M. Robinson

Linear convergence of epsilon-subgradient descent methods for a class of convex functions*

Springer-Verlag 1999

Digital Object Identifier (DOI) 10.1007/s101079900088 Stephen M. Robinson

Linear convergence of epsilon-subgradient descent methods for a class of convex functions*

Consider a method defining iterates

 $x^{k+1} = x^k - t_k g^k$ for $g^k \in \partial_{\mathcal{E}} f(x^k)$ and such that

- (DR) holds with $\delta_k = \varepsilon + t_k \|g^k\|^2$
- $t_k \in [t_{\min}, t_{\max}]$ with $0 < t_{\min} \le t_{\max} < \infty$

Digital Object Identifier (DOI) 10.1007/s101079900088 Stephen M. Robinson

Linear convergence of epsilon-subgradient descent methods for a class of convex functions*

Consider a method defining iterates

 $x^{k+1} = x^k - t_k g^k$ for $g^k \in \partial_{\mathcal{E}} f(x^k)$ and such that

• (DR) holds with $\delta_k = \varepsilon + t_k \|g^k\|^2$

• $t_k \in [t_{\min}, t_{\max}]$ with $0 < t_{\min} \le t_{\max} < \infty$

If f is not too "flat" around its set of minimizers X,

 $f(x) \ge \inf f + cdist(x, \overline{X})$ for any $x \in \overline{X} + \eta \mathbb{B}$,

Digital Object Identifier (DOI) 10.1007/s101079900088 Stephen M. Robinson

Linear convergence of epsilon-subgradient descent methods for a class of convex functions*

Consider a method defining iterates

 $x^{k+1} = x^k - t_k g^k$ for $g^k \in \partial_{\mathcal{E}} f(x^k)$ and such that

• (DR) holds with $\delta_k = \varepsilon + t_k \|g^k\|^2$

• $t_k \in [t_{\min}, t_{\max}]$ with $0 < t_{\min} \le t_{\max} < \infty$

If f is not too "flat" around its set of minimizers \overline{X} ,

 $f(x) \ge \inf f + cdist(x, \overline{X})$ for any $x \in \overline{X} + \eta \mathbb{B}$,

the rate of convergence if *R*-linear

Springer-Verlag 1999

Digital Object Identifier (DOI) 10.1007/s101079900088 Stephen M. Robinson

Linear convergence of epsilon-subgradient descent methods for a class of convex functions^{*}

Consider a method defining iterates

 $x^{k+1} = x^k - t_k g^k$ for $g^k \in \partial_{\mathcal{E}} f(x^k)$ and such that

• (DR) holds with $\delta_k = \varepsilon + t_k ||q^k||^2$

▶ $t_k \in [t_{\min}, t_{\max}]$ with $0 < t_{\min} \le t_{\max} < \infty$

If f is not too "flat" around its set of minimizers X.

 $f(x) \geq \inf f + cdist(x, \overline{X})$ for any $x \in \overline{X} + \eta \mathbb{B}$,

the rate of convergence if *R*-linear

note: condition \iff "inverse growth condition" Kurdyka-Lojasiewicz inequality

Jean-Jacques Moreau and its envelope

The Moreau-envelope of f is a $C^{1,1}$ -smoothing of f

$$F_t(x) := \min\left\{f(y) + \frac{1}{2t}||y - x||^2\right\}$$

Jean-Jacques Moreau and its envelope

The Moreau-envelope of f is a $C^{1,1}$ -smoothing of f

$$F_t(x) := \min\left\{f(y) + \frac{1}{2t} \|y - x\|^2\right\}$$



•the unique minimizer is the proximal point mapping $p_t^f(x)$

•the envelope's gradient is $\nabla F_t(x) = \frac{1}{t} \left(x - p_t^f(x) \right)$



min f(x) is equivalent to $x^* = p_t(x^*)$



min
$$f(x)$$
 is equivalent to $x^* = p_t(x^*)$

Picard's iteration
$$x^{k+1} = p_t(x^k)$$



 \mathbf{x}^{k}

min
$$f(x)$$
 is equivalent to $x^* = p_t(x^*)$

Picard's iteration

$$= p_t(x^k)$$

$$= x^k - \frac{t}{t} \left(x^k - p_t(x^k) \right)$$



min
$$f(x)$$
 is equivalent to $x^* = p_t(x^*)$

Picard's iteration X

$$k^{k+1} = p_t(x^k)$$

$$= x^k - \frac{t}{t} \left(x^k - p_t(x^k) \right)$$

$$= x^k - t \nabla F_t(x^k)$$





a gradient method to minimize Moreau's envelope





a gradient method to minimize Moreau's envelope

note: the "stepsize" t is NOT computed by a linesearch

(curve-search on the metric)

Replace the proximal parameter by a *matrix*:

Instead of $\frac{1}{t}$

Replace the proximal parameter by a *matrix*:

Instead of $\frac{1}{t}$ use a positive definite matrix M :

$$F_M(x) := \min\left\{f(y) + \frac{1}{2} \|y - x\|_M^2\right\}$$

Replace the proximal parameter by a *matrix*:

Instead of $\frac{1}{t}$ use a positive definite matrix M:

$$F_M(x) := \min\left\{f(y) + \frac{1}{2}||y - x||_M^2\right\}$$

- There is a unique minimizer $p_M(x)$
- The proximal point operator is Lipschitzian
- The relation $\nabla F_M(x) = M(x p_M(x))$ holds
- $\blacktriangleright \min f \iff \min F_M \iff x = p_M(x)$

Replace the proximal parameter by a *matrix*:

Instead of $\frac{1}{t}$ use a positive definite matrix M:

$$F_M(x) := \min\left\{f(y) + \frac{1}{2}||y - x||_M^2\right\}$$

- There is a unique minimizer $p_M(x)$
- The proximal point operator is Lipschitzian
- The relation $\nabla F_M(x) = M(x p_M(x))$ holds
- $\blacktriangleright \min f \iff \min F_M \iff x = p_M(x)$

Interest: now we perceive better the role of M

The role of the metric

Picard's iteration
$$x^{k+1} = p_M(x^k)$$

= $x^k - M^{-1}M(x^k - p_t(x^k))$
= $x^k - M^{-1}\nabla F_M(x^k)$

a gradient method to minimize Moreau's envelope, preconditioned by the matrix $M^{-1} \approx \nabla^{-2} F_M(x^k)$

The role of the metric

Picard's iteration
$$x^{k+1} = p_M(x^k)$$

= $x^k - M^{-1}M(x^k - p_t(x^k))$
= $x^k - M^{-1}\nabla F_M(x^k)$

a gradient method to minimize Moreau's envelope, preconditioned by the matrix $M^{-1} \approx \nabla^{-2} F_M(x^k)$ We know $F_M \in C^{1,1}$... when does the Hessian exist? Existe-t-il une géneralisation adéquate de la notion de Hessien?...Cette question est la plus passionnante qui se pose actuellement, et une réponse satisfaisante marquerait probablement pour longtemps une étape décisive dans les recherches fondamentales en programmation mathématique.

"Does an adequate generalization for the notion of a Hessian exist?... This is today's most interesting question, to which a satisfactory answer would probably start a new, longlasting and decisive era for basic research in Mathematical Programming."

Chapter 5 in:

THÈSE DOCTORAT D'ÉTAT ES SCIENCES MATHÉMATIQUES L'UNIVERSITÉ PARIS IX Claude LEMARÉCHAL pour obtenir LE GRADE DE DOCTEUR ES SCIENCES Sujet de la thèse EXTENSIONS DIVERSES DES MÉTHODES DE GRADIENT ET APPLICATIONS Soutenue le 4 décembre 1980 devant la Commission composée de : LP ALIDIN AUSLENDER

December 4th, 1980

How structural nonsmoothness has been exploited

Improving speed only possible if algorithm incorporates **structure** information

1995-2000: \mathcal{U} -Lagrangian (Lemaréchal, Oustry, Sagastiz)

1999: *VU*-decomposition (Mifflin, Sagastiz)

2002: \mathcal{M} -manifolds, partly smooth functions (Lewis, Hare)

2003: Composite objective functions (Shapiro)

Later on, special minimization of composite objective functions revisited: Lewis & Wright, Nesterov, Planiden, Hare & Sagastiz, Liu, Sagastiz & Solodov

Illustrative examples

The half-and-half function in \mathbb{R}^2 $f(x_1, x_2) = |x_1| + bx_2^2$


Proximal point: calculus rules

► separable sum: $f(x,y) = g(x) + h(y) \Longrightarrow$ $p_t^f(x) = \left(p_t^g(x), p_t^h(y)\right)$

- ► scalar factor ($\alpha \neq 0$) and translation ($v \neq 0$): $f(x) = g(\alpha x + v) \Longrightarrow$ $p_t^f(x) = \frac{1}{\alpha} \left(p_t^{\alpha^2 g}(\alpha x + v) - v \right)$
- "perspective" ($\alpha > 0$): $f(x) = \alpha g(\frac{1}{\alpha}x) \Longrightarrow p_t^f(x) = \alpha p_t^{g/\alpha}(\frac{x}{\alpha})$

Proximal point: special functions

$$p_t^f(x) = \left(p_t^{f_1}(x_1), p_t^{f_2}(x_2)\right)$$

$$p_t^f(x) = \left(p_t^{f_1}(x_1), p_t^{f_2}(x_2)\right)$$

$$p_t^{f_1}(x_1) = x_1 - proj_{[-t,t]}(x_1)$$

$$p_t^{f}(x) = \left(p_t^{f_1}(x_1), p_t^{f_2}(x_2)\right)$$

$$p_t^{f_1}(x_1) = x_1 - proj_{[-t,t]}(x_1)$$

$$p_t^{f_2}(x_2) = \frac{1}{1+bt}x_2$$

$$p_t^{f}(x) = \left(p_t^{f_1}(x_1), p_t^{f_2}(x_2)\right)$$

$$p_t^{f_1}(x_1) = x_1 - proj_{[-t,t]}(x_1)$$

$$p_t^{f_2}(x_2) = \frac{1}{1+bt} x_2$$

$$F_t(x) = ?$$
 and $\nabla F_t(x) = ?$

$$p_t^{f}(x) = \left(p_t^{f_1}(x_1), p_t^{f_2}(x_2)\right)$$

$$p_t^{f_1}(x_1) = x_1 - proj_{[-t,t]}(x_1)$$

$$p_t^{f_2}(x_2) = \frac{1}{1+bt}x_2$$

$$F_t(x) = ?$$
 and $\nabla F_t(x) = ?$
What about a Hessian?

- Let *M* be the current matrix
- ▶ for example, let $u = \Delta x$, and $v = \Delta g$

- ► Let *M* be the current matrix
- for example, let $u = \Delta x$, and $v = \Delta g$
- The quasi-Newton equation
- $M_+u = v$ should hold for the update.

- Let M be the current matrix
- for example, let $u = \Delta x$, and $v = \Delta g$
- The quasi-Newton equation
- $M_+u = v$ should hold for the update.
- If $M = \mu I$, a scalar multiple of the identity satisfying the qN equation may be imposible, as it writes down

$$\mu_{k+1}u = v$$

with *u* and *v* vectors. Instead, μ_{k+1} solves

$$\min_{\mu}\frac{1}{2}\|v/\mu-u\|^2$$



- Let M be the current matrix
- ▶ for example, let $u = \Delta x$, and $v = \Delta g \leftarrow$ options!
- The quasi-Newton equation
- $M_+u = v$ should hold for the update.
- If $M = \mu I$, a scalar multiple of the identity satisfying the qN equation may be imposible, as it writes down

$$\mu_{k+1}u=v$$

with *u* and *v* vectors. Instead, μ_{k+1} solves

$$\min_{\mu}\frac{1}{2}\|v/\mu-u\|^2$$



Reversal quasi-Newton update

- Given u and Δg , invert the prox
- compute v accordingly

Reversal quasi-Newton update

- Given u and Δg , invert the prox
- compute v accordingly



Reversal quasi-Newton update

- Given u and Δg , invert the prox
- compute v accordingly

 $\mu=?$ As in Lemaréchal C., Sagastizábal C. (1994) An approach to variable metric bundle methods. In: Henry

J., Yvon JP. (eds) System Modelling and Optimization, LNCIS vol 197. Springer, Berlin, Heidelberg

$$\frac{1}{\mu_{n+1}} = \frac{1}{\mu_n} + \frac{\langle v, \Delta x \rangle}{|v|^2}$$

whenever $\langle v, u \rangle > 0$.

Theorem 10. If ∇f is locally Lipschitzian, then $\mu_n \to 0$. Make the following additional assumptions: f has a (unique) minimal point \bar{x} and a quadratic growth condition holds: for some $\alpha > 0$,

$$f(x) \ge f(\bar{x}) + \alpha |x - \bar{x}|^2.$$

Then $f(x_n)$ tends to $f(\bar{x})$ q-superlinearly.

drop M in F_M

For Newton's-type methods to apply, we need F to have an invertible Hessian at \bar{x}

- For Newton's-type methods to apply, we need F to have an invertible Hessian at \bar{x}
- A convex function φ has a *generalized Hessian* $H\varphi(x_0)$ at x_0 if the gradient $\nabla \varphi(x_0)$ exists and there exists a symmetric positive semidefinite operator $H\varphi(x_0)$ such that

$$\partial \varphi(x_0+d) \subset
abla \varphi(x_0) + H \varphi(x_0) d + \|d\| \mathbb{B}$$

If f at a generalized Hessian at p(x₀), then the Hessian of F exists at x₀ and

$$\nabla^2 F(x_0) = M - M[Hf(p(x_0)) + M]^{-1}M$$

drop M in F_M

For Newton's-type methods to apply, we need F to have an invertible Hessian at \bar{x}

For Newton's-type methods to apply, we need F to have an invertible Hessian at \bar{x}

We know that

f convex and lower sci \implies F is $C^{1,1}$ everywhere

For Newton's-type methods to apply, we need F to have an invertible Hessian at \bar{x}

We know that

f convex and lower sci \implies F is $C^{1,1}$ everywhere

- **BUT:** *F* is C^2 everywhere \Longrightarrow
 - f is C^2 everywhere, too!

For Newton's-type methods to apply, we need F to have an invertible Hessian at \bar{x}

We know that

f convex and lower sci \implies F is $C^{1,1}$ everywhere

BUT: *F* is C^2 everywhere \Longrightarrow

f is C^2 everywhere, too!

► The Hessian of *F* exists at x_0 if and only if *p* has a Jacobian: $\nabla^2 F(x_0) = M(I - p'(x_0))$.

For Newton's-type methods to apply, we need F to have an invertible Hessian at \bar{x}

We know that

f convex and lower sci \implies F is $C^{1,1}$ everywhere

- **BUT:** *F* is C^2 everywhere \Longrightarrow
 - f is C^2 everywhere, too!
- ► The Hessian of *F* exists at x_0 if and only if *p* has a Jacobian: $\nabla^2 F(x_0) = M(I p'(x_0))$.
- When the Jacobian exists, its image lies in a very specific subspace

For Newton's-type methods to apply, we need F to have an invertible Hessian at \bar{x}

We know that

f convex and lower sci \implies F is $C^{1,1}$ everywhere

- **BUT:** F is C^2 everywhere \Longrightarrow
 - f is C^2 everywhere, too!
- ► The Hessian of *F* exists at x_0 if and only if *p* has a Jacobian: $\nabla^2 F(x_0) = M(I p'(x_0))$.
- When the Jacobian exists, its image lies in a very specific subspace
- By considering trajectories not far from that subspace we can define a 2nd-order object of *f*, even if *f* is not differentiable

Views of some functions The half-and-half $f(x_1, x_2) = |x_1| + bx_2^2$



Views of some functions

A max-variation $f(x_1, x_2) = \max(|x_1|, bx_2^2)$



The key is in decomposing the space

$\blacktriangleright \mathbb{R}^n = \mathcal{V} \oplus \mathcal{U}$

- On the \mathcal{V} -subspace the function is sharp
- On the U-subspace the function looks smooth



The subspace definition depends on *x*

The key is in decomposing the space

- $\blacktriangleright \mathbb{R}^n = \mathcal{V} \oplus \mathcal{U} + \bar{x}$
- On the \mathcal{V} -subspace the function is sharp
- On the U-subspace the function looks smooth



The subspace definition depends on \bar{x}

The key is in decomposing the space

- $\blacktriangleright \mathbb{R}^n = \mathcal{V} \oplus \mathcal{U} + \bar{x}$
- On the \mathcal{V} -subspace the function is sharp
- On the U-subspace the function looks smooth



The subspace definition depends on \bar{x}

Nonsmoothness appears in a structured manner

The *n*-dimensional half-and-half function For $x \in \mathbb{R}^n$, given matrices *A* with nontrivial kernel, $B \succ 0$, $f(x) = \max(|x^TAx|, x^TBx)$

has a unique minimizer at $\bar{x} = 0$.

The *n*-dimensional half-and-half function For $x \in \mathbb{R}^n$, given matrices *A* with nontrivial kernel, $B \succ 0$,

 $f(x) = \max(|x^{\top}Ax|, x^{\top}Bx)$ has a unique minimizer at $\bar{x} = 0$. On $\mathcal{N}(A)$ the

function is not differentiable, and the first term vanishes: $f|_{\mathcal{N}(A)}$ looks smooth

The *n*-dimensional half-and-half function For $x \in \mathbb{R}^n$, given matrices *A* with nontrivial kernel, $B \succ 0$,

 $f(x) = \max(|x^{\top}Ax|, x^{\top}Bx)$ has a unique minimizer at $\bar{x} = 0$. On $\mathcal{N}(A)$ the function is not differentiable, and the first term vanishes: $f|_{\mathcal{N}(A)}$ looks smooth



$\mathcal{V}\mathcal{U}s$ of the function

 $f(x) = \max\left(\left|x^{\top}Ax\right|, x^{\top}Bx\right)$












\mathcal{VU} s of the function $f(x) = \max(|x^{T}Ax|, x^{T}Bx)$



\mathcal{VU} s of the function $f(x) = \max(|x^{T}Ax|, x^{T}Bx)$



$\mathcal{VU}_{\text{How can we track}}^{\text{the function } f(x)} = \max(|x^{\top}Ax|, x^{\top}Bx)$ the valley of nonsmoothness,

where *f* looks "nice"?



\mathcal{VU} s of the function $f(x) = \max(|x^{T}Ax|, x^{T}Bx)$



\mathcal{VU} -theory & primal-dual tracks

U -Lagrangian:

$$\begin{array}{rcl} L_U(u,\bar{g}) &:=& \inf_{v \in V} \{f(\bar{x}+u \oplus v) - \langle \bar{g}, v \rangle \} \\ &=& f(\bar{x}+u \oplus v(u)) - \langle \bar{g}, v(u) \rangle \\ \to L_U(0,\bar{g}) = f(\bar{x}), &\to L_U(u,\bar{g}) \in C^1(U) \end{array}$$

 \rightarrow minimizer v = v(u) generates trajectories

smooth tangent to U



 $orall ar{g}\in ri\,\partial f(ar{x})$ same v(u) and $L_U\in C^2$ i.e., has a U-Hessian

Dual tracks



- proximal point Newton method

Approximating primal-dual tracks Fundamental theoretical result:

Proximal Points are on the primal track If

 $arg = 0 \in ri\partial f(\bar{x})$, then $\exists u(x) :$

$$p(x) := \operatorname{argmin}\left\{f(y) + \frac{1}{2}\mu|y-x|^2\right\} = \chi(u(x))$$

for all $x \approx \bar{x}$ with $\mu = \mu(x) : \mu |x - \bar{x}| \to 0$ as $x \to \bar{x}$

 $\Rightarrow \text{ use a bundle subroutine} \\ \text{to approximate the prox} \\ \text{and estimate the pair} \\ (\chi(u), \gamma(u)) \end{aligned}$



By-product: local V U -decomposition, VU

Newton-like corrector-predictor V U algorithm



- **Corrector step:** Solve $(\chi - \text{and } \gamma \text{-} \text{QP})$'s

 \Rightarrow new *p*, *s*, VU, and determine *H* (*U*-Hessian)

- **Predictor step:** Solve $H\Delta u = -U^{\top}s \Rightarrow x^{+} = p + U\Delta u$

Convergence properties

- If infinite bundle steps, the inner sequence converges to a minimizer of *f*
- Otherwise, either the outer sequence $\{p\}$ is finite with s = 0 and last p minimizes f
- or, $\{f(p)\}$ is infinite and decreasing \Rightarrow
 - either f unbounded below,
 - or $s \to 0$ and any $acc(\{p\})$ minimizes fwhen $\{\sigma/\mu\} \to 0$.

If the U -Hessian at \bar{x} is positive definite, $0 \in ri \partial f(\bar{x})$, and

- $-rac{\sigma}{\mu^2}=O(|s^-|^2),$ bounded $\{H^{-1}\},$ $\mathbb{U}
 ightarrow U,$
- Dennis-Moré-like condition for $\{H\}$
- -s approximates γ superlinearly

 \Rightarrow superlinear convergence of $\{p\}$ to \bar{x}

VU-Algorithm:

As $u \rightarrow 0$:



VU-Algorithm:

As $u \rightarrow 0$: \mathcal{V} -step: The **fast** track $\chi(u) \rightarrow x^*$ ▶ \mathcal{U} -step: The \mathcal{U} -gradient, $\nabla L_{\mathcal{U}}(u; g_{\mathcal{V}}^*) \rightarrow 0$ \overline{x}

VU-Algorithm:







$\mathcal{V}\mathcal{U}\text{-}Algorithm$ is globally and superlinearly convergent



$\mathcal{V}\mathcal{U}\text{-}Algorithm$ is globally and superlinearly convergent

Comparison with BFGS method



Concluding comments

On-going & Future work

- O: Derivative-free variant
- O: *ɛ*-V U variant
- F: Application to 2-stage stochastic programming problems