

Convex V-models for nonconvex minimization

Claudia Sagastizábal

`mailto:sagastiz@impa.br, http://www.impa.br/~sagastiz`

Optimization Day, Grenoble, Nov. 5, 2014

Joint work with R. mifflin

With thanks to:

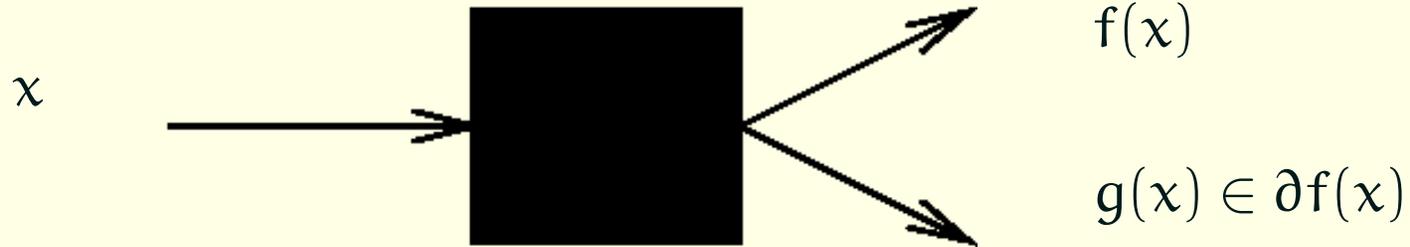
**AFOSR Grant FA9550-08-1-0370, NSF Grant DMS 0707205,
and CNPq & Faperj from Brazil**

NSO Algorithms

For a nonsmooth function, solving

$$\min f(x)$$

with a black-box method



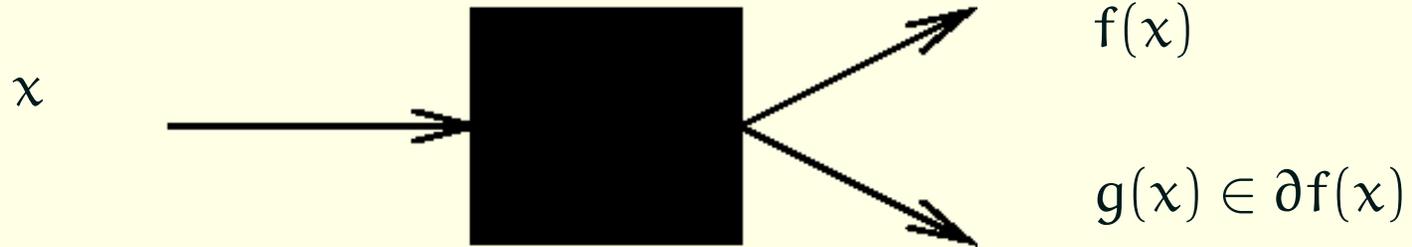
is doomed to slow convergence speed.

NSO Algorithms

For a nonsmooth function, solving

$$\min f(x)$$

with a black-box method



is doomed to slow convergence speed.

Better performance possible by exploiting structure

How does structure appear?

- Implicitly

 - U-Lagrangian

 - VU-decomposition

 - partly smooth functions

- Explicitly

 - as a sum

 - as a composition

How does structure appear?

– Implicitly

U-Lagrangian

VU-decomposition

partly smooth functions

How does structure appear?

– Implicitly

U-Lagrangian

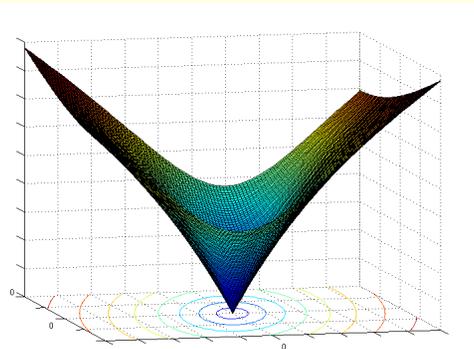
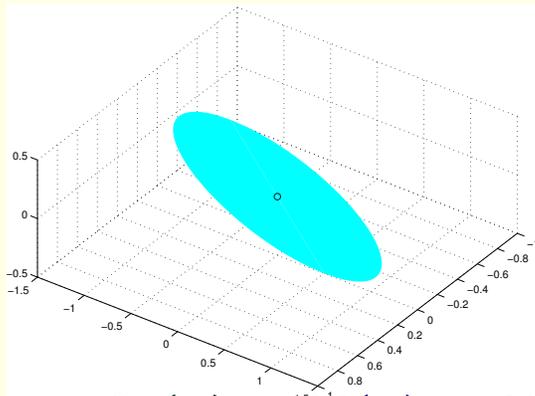
VU-decomposition

partly smooth functions

$$f(x_1, x_2, x_3) = \frac{1}{2}x_1^2 + \frac{1}{2} \ln \left(1 + \sqrt{(x_1^2 - 2x_2)^2 + (x_3 - x_2)^2} \right)$$

$x^* = (0, 0, 0)$ is a stationary point (minimizer)

zero subgradient $\in \partial f(x^*)$

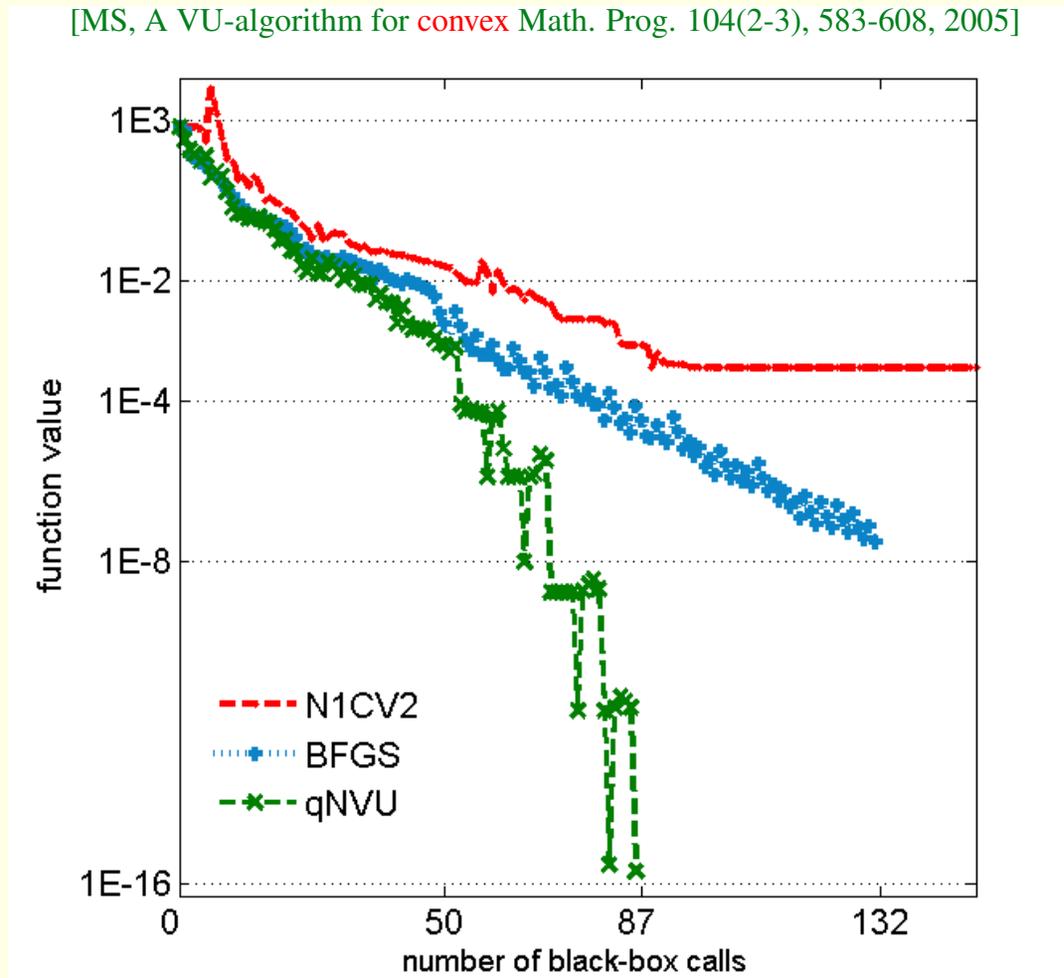


$$\bar{g} \in \partial f(\bar{x}), \quad \mathcal{V}(\bar{x}) := \text{lin}(\partial f(\bar{x}) - \bar{g}) \quad \text{and} \quad \mathcal{U}(\bar{x}) := \mathcal{V}(\bar{x})^\perp$$

CONVEX CASE

Lewis and Overton 8-variable half-and-half function

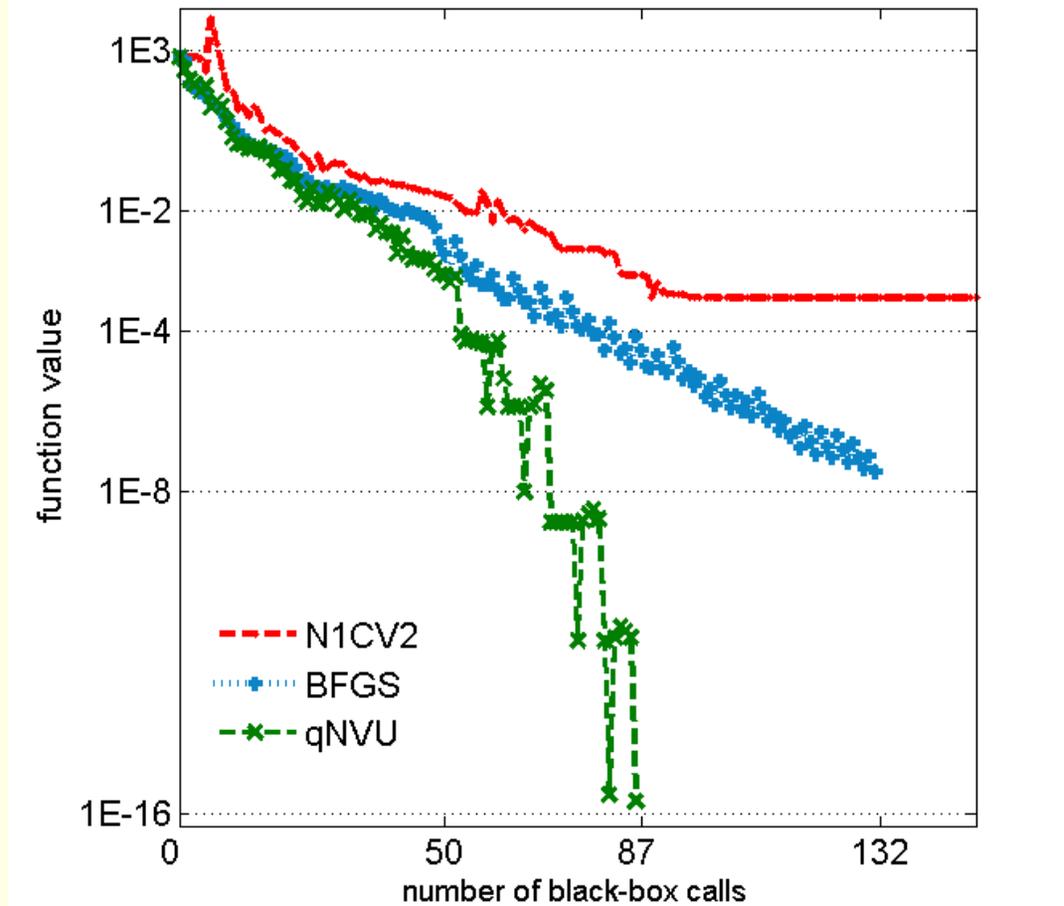
[MS, A VU-algorithm for **convex** Math. Prog. 104(2-3), 583-608, 2005]



Sublinear, linear, and superlinear convergence - convex case

Lewis and Overton 8-variable half-and-half function

[MS, A VU-algorithm for **convex** Math. Prog. 104(2-3), 583-608, 2005]



Sublinear, linear, and superlinear convergence -

What can be done for nonconvex case?

How does structure appear?

– Implicitly

U-Lagrangian

VU-decomposition

partly smooth functions

digging tools

– Explicitly

as a sum

as a composition

blackboxes

How does structure appear?

– Implicitly

U-Lagrangian

VU-decomposition

partly smooth functions

digging tools

– Explicitly

as a sum

as a composition

≠ black boxes

Psyche Opening the Golden Box, 1903 -J.W. Waterhouse

Explicit Structure: Opening the Black Box



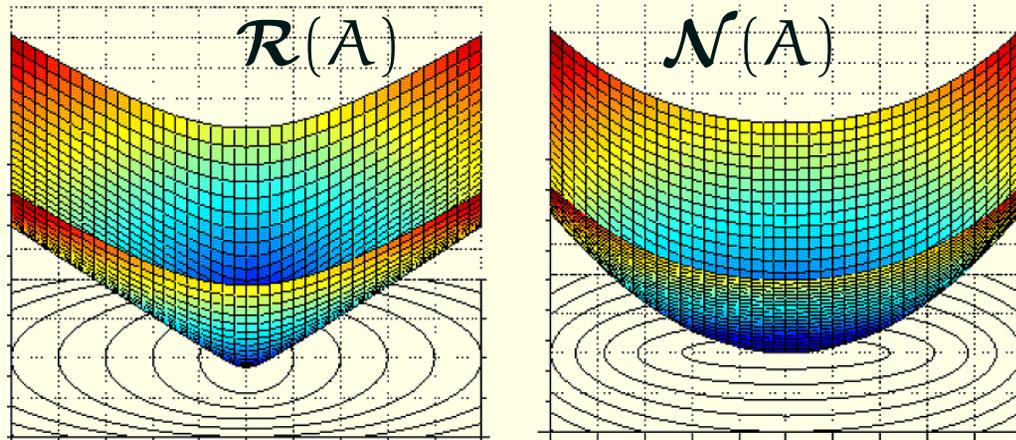
The half-and-half function

For $x \in \mathbb{R}^n$, given matrices $A \succeq 0$, $B \succ 0$,

$$f(x) = \sqrt{\langle x, Ax \rangle} + \langle x, Bx \rangle$$

has a unique minimizer at 0.

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

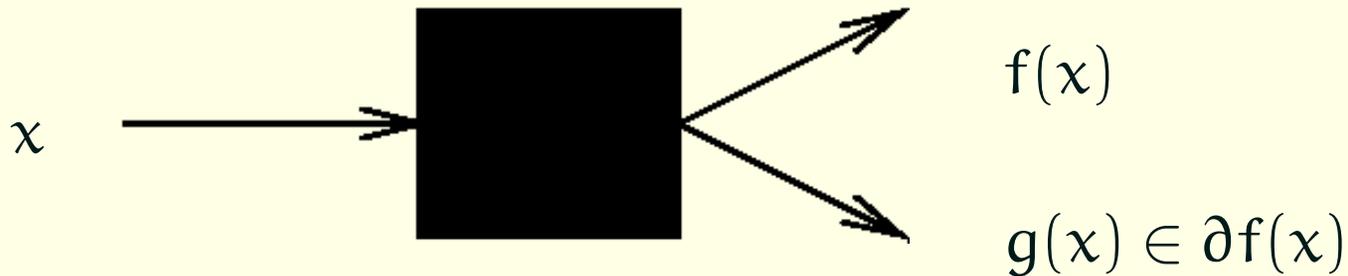


The half-and-half function has several interesting structures

If no structure at all

$$f(x) = \sqrt{\langle x, Ax \rangle} + \langle x, Bx \rangle$$

This defines the **black box**:

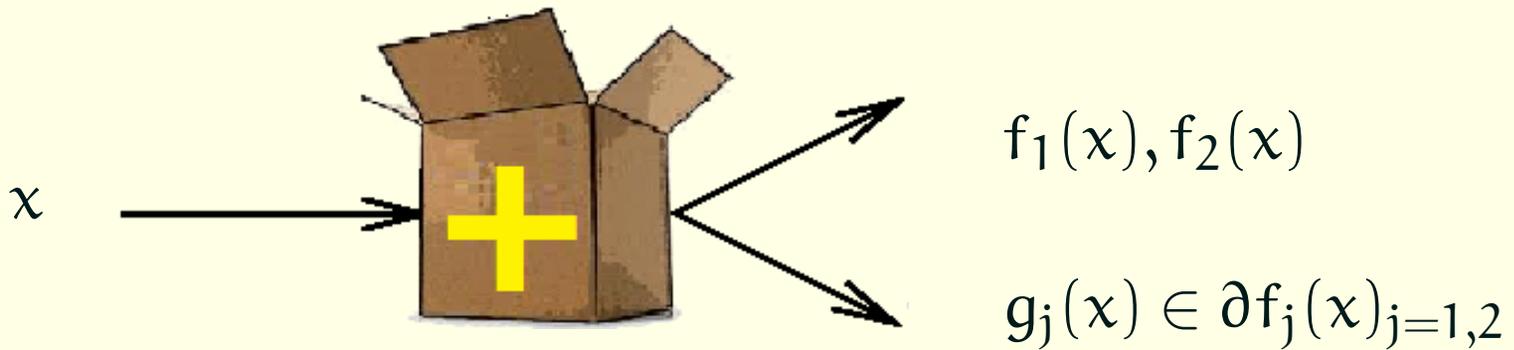


The half-and-half function has several interesting structures

Sum structure

$$f(x) = f_1(x) + f_2(x) \text{ with } \begin{cases} f_1(x) = \sqrt{\langle x, Ax \rangle} \\ f_2(x) = \langle x, Bx \rangle \quad (\text{smooth } \Psi) \end{cases}$$

This defines a **sum black box**:



The half-and-half function has several interesting structures

Composite structure

$$f(x) = (h \circ c)(x) \text{ with } \begin{cases} c(x) = (x, \langle x, Bx \rangle) \in \mathbb{R}^{n+1} \\ h(C) = \sqrt{\langle C_{1:n}, AC_{1:n} \rangle} + C_{n+1} \end{cases}$$

for C smooth and h positively homogeneous

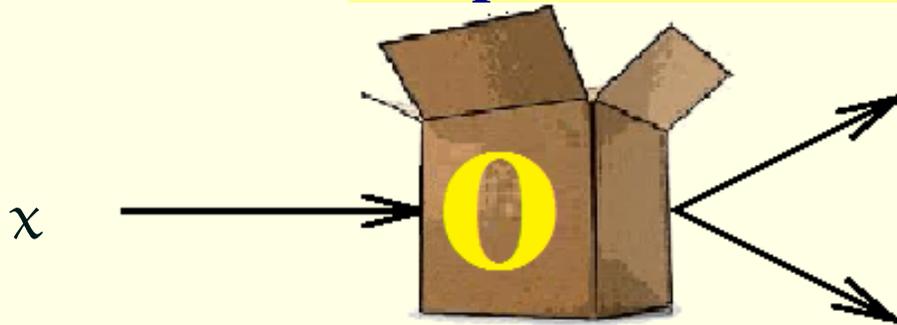
The half-and-half function has several interesting structures

Composite structure

$$f(x) = (h \circ c)(x) \text{ with } \begin{cases} c(x) = (x, \langle x, Bx \rangle) \in \mathbb{R}^{n+1} \\ h(C) = \sqrt{\langle C_{1:n}, AC_{1:n} \rangle} + C_{n+1} \end{cases}$$

for C smooth and h positively homogeneous

This defines a **composite black box**:



$C := c(x)$ and $h(C)$

Jacobian $c'(x)$ and

$G(C) \in \partial h(C)$

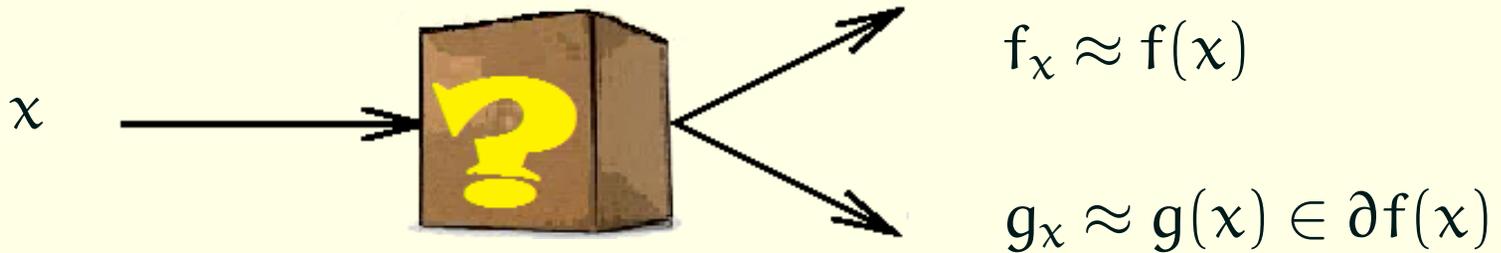
The half-and-half function has several interesting structures

Missing information structure

Suppose not all of A/B is known/accessible,

so that only **estimates** are available for f

This defines a **noisy black box**:



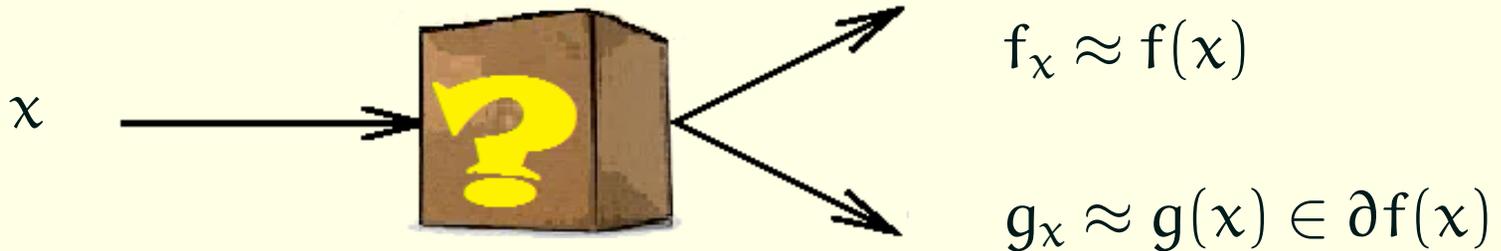
The half-and-half function has several interesting structures

Missing information structure

Suppose not all of A/B is known/accessible,

so that only **estimates** are available for f

This defines a **noisy black box**:



Not considered in this talk

(Convex: OISagLem, NonConvex: HareSagSol)

The Bottom Line

- The same function can be seen in different manners**
- Depending on how the oracle outputs its f/g information, different structure is revealed:**
 - A sum of two nonsmooth functions**
 - A sum of a nonsmooth function and a smooth one**
 - A composition of a nonsmooth function with a smooth mapping**

The Bottom Line

- The same function can be seen in different manners
- Depending on how the oracle outputs its f/g information, different structure is revealed:
 - A sum of two nonsmooth functions
 - A sum of a nonsmooth function and a smooth one
 - A composition of a nonsmooth function with a smooth mapping
- What is the impact of structure in a bundle algorithm?

f/g

The Bottom Line

- The same function can be seen in different manners
- Depending on how the oracle outputs its f/g information, different structure is revealed:
 - A sum of two nonsmooth functions
 - A sum of a nonsmooth function and a smooth one
 - A composition of a nonsmooth function with a smooth mapping
- What is the impact of structure in a bundle algorithm?
- To which extent does the convergence analysis rely on the specific f/g output?

The Bottom Line

- The same function can be seen in different manners
- Depending on how the oracle outputs its f/g information, different structure is revealed:
 - A sum of two nonsmooth functions
 - A sum of a nonsmooth function and a smooth one
 - A composition of a nonsmooth function with a smooth mapping
- What is the impact of structure in a bundle algorithm?
- To which extent does the convergence analysis rely on the specific f/g output?
- Is the composite structure worth looking at?

Is the composite structure worth looking at?

- **Max-functions:**

$$f(\mathbf{x}) = \max\{f_1(\mathbf{x}), \dots, f_m(\mathbf{x})\} \text{ so } \begin{cases} c_i(\mathbf{x}) = f_i(\mathbf{x}), i = 1 : m \\ h(\mathbf{c}) = \max(c_1, \dots, c_m) \end{cases}$$

- **Sparse least squares (compressed sensing):**

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{R}\mathbf{x} - \mathbf{b}\|_2^2 + \tau \|\mathbf{x}\|_1 \text{ so } \begin{cases} c_1(\mathbf{x}) = \frac{1}{2} \|\mathbf{R}\mathbf{x} - \mathbf{b}\|_2^2 \\ c_{i+1}(\mathbf{x}) = x_i, i = 1 : m \\ h(\mathbf{c}) = c_1 + \tau \|\mathbf{c}_{2:n+1}\|_1 \end{cases}$$

- **Sum of norms** $f(\mathbf{x}) = \sum_{i=1}^m |f_i(\mathbf{x})|$ so $\begin{cases} c_i(\mathbf{x}) = f_i(\mathbf{x}), i = 1 : m \\ h(\mathbf{c}) = \sum_{i=1}^m |c_i| \end{cases}$

- **Sum of leading eigenvalues, ℓ_1 -penalization of NLPs, etc**

Is the composite structure worth looking at?

- **Max-functions:**

$$f(x) = \max\{f_1(x), \dots, f_m(x)\} \text{ so } \begin{cases} c_i(x) = f_i(x), i = 1 : m \\ h(c) = \max(c_1, \dots, c_m) \end{cases}$$

- **Sparse least squares (compressed sensing):**

$$f(x) = \frac{1}{2} \|Rx - b\|_2^2 + \tau \|x\|_1 \text{ so } \begin{cases} c_1(x) = \frac{1}{2} \|Rx - b\|_2^2 \\ c_{i+1}(x) = x_i, i = 1 : m \\ h(c) = c_1 + \tau \|c_{2:n+1}\|_1 \end{cases}$$

- **Sum of norms** $f(x) = \sum_{i=1}^m |f_i(x)|$ so $\begin{cases} c_i(x) = f_i(x), i = 1 : m \\ h(c) = \sum_{i=1}^m |c_i| \end{cases}$

- **Sum of leading eigenvalues, ℓ_1 -penalization of NLPs, etc**

YES!

NOTE: f may fail to be convex

Impact of Structure Knowledge

Black-box information defines **pieces** that put together create a **model \mathcal{M}** of the function f . The model is used to define iterates not too far away from a “good” past iterate, \hat{x} .

At iteration i ,

$$y^{i+1} \text{ minimizes } \mathcal{M}(y) + \frac{1}{2}\mu|y - \hat{x}|^2$$

Impact of Structure Knowledge

Black-box information defines **pieces** that put together create a **model M** of the function f . The model is used to define iterates not too far away from a “good” past iterate, \hat{x} .

At iteration i ,

$$y^{i+1} \text{ minimizes } M(y) + \frac{1}{2}\mu|y - \hat{x}|^2$$

“pieces” chosen to make minimization simple (convex QP)

for example, “piece”=linearization:

$$y^i \longrightarrow \blacksquare \begin{cases} f^i = f(y^i) \\ g^i = g(y^i) \end{cases} \implies f^i + \langle g^i, y - y^i \rangle$$

Impact of Structure Knowledge

Black-box information defines **pieces** that put together create a **model M** of the function f . The model is used to define iterates not too far away from a “good” past iterate, \hat{x} .

At iteration i ,

$$y^{i+1} \text{ minimizes } M(y) + \frac{1}{2}\mu|y - \hat{x}|^2$$

“pieces” chosen to make minimization simple (convex QP)

for example, “piece”=linearization:

$$y^i \rightarrow \blacksquare \begin{cases} f^i = f(y^i) \\ g^i = g(y^i) \end{cases} \implies M(y) = \max_i \left\{ f^i + \langle g^i, y - y^i \rangle \right\}$$

$$y^{i+1} = \arg \min_x M(y) + \frac{1}{2} \mu |y - \hat{x}|^2$$

for example, “piece”=linearization:

$$y^i \rightarrow \blacksquare \begin{cases} f^i = f(y^i) \\ g^i = g(y^i) \end{cases} \implies M(y) = \max_i \left\{ f^i + \langle g^i, y - y^i \rangle \right\}$$

Some jargon

\hat{x} is a serious point

$\bigcup_i (y^i, f^i, g^i)$ is the bundle B

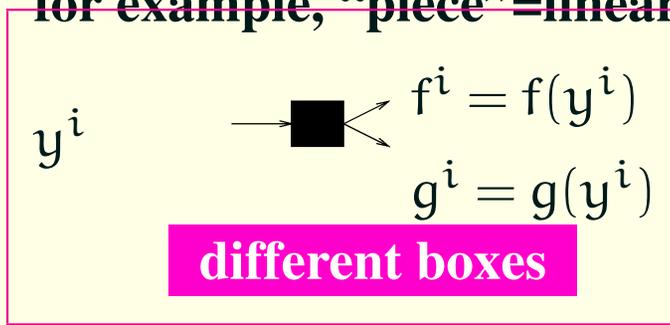
If y^{i+1} gives sufficient decrease for f (in terms of a progress measure),

it becomes the next \hat{x}

Otherwise, it is declared a null point

$$y^{i+1} = \arg \min_x M(y) + \frac{1}{2} \mu |y - \hat{x}|^2$$

for example, “piece”=linearization:



$$\implies M(y) = \max_i \left\{ f^i + \langle g^i, y - y^i \rangle \right\}$$

different models

Some jargon

\hat{x} is a serious point

$\bigcup_i (y^i, f^i, g^i)$ is the bundle B

If y^{i+1} gives sufficient decrease for f (in terms of a progress measure),
it becomes the next \hat{x}

Otherwise, it is declared a null point

Back to our questions

OK! Is the composite structure worth looking at?

OK! What is the impact of structure in a bundle algorithm?

- To which extent does the convergence analysis rely on the specific f/g output?

Back to our questions

OK! Is the composite structure worth looking at?

OK! What is the impact of structure in a bundle algorithm?

- To which extent does the convergence analysis

rely on the specific f/g output?

rely on the specific model?

Back to our questions

OK! Is the composite structure worth looking at?

OK! What is the impact of structure in a bundle algorithm?

- To which extent does the convergence analysis

rely on the specific f/g output?

rely on the specific model?

What about convexity?

The setting

$\min_{\mathbf{y} \in \mathbb{R}^n} f(\mathbf{y})$; f locally Lipschitz,
only one (Clarke) gradient $g(\mathbf{y})$,
computed by a black box at each \mathbf{y} .

Black-box information used to build a model function

$M(\mathbf{d}) \approx f(\hat{\mathbf{x}} + \mathbf{d})$ near a center $\hat{\mathbf{x}}$

To which extent does the convergence analysis
rely on the specific model?

Some Model Functions

($y = \hat{x} + d$)

$$f(y) \quad \Rightarrow \quad M(d) = \max_{i \in B} \left\{ f^i + \left\langle g^i, \hat{x} + d - y^i \right\rangle \right\}$$

f_1

f_1

h

Some Model Functions

$$(y = \hat{x} + d)$$

$$\begin{aligned} f(y) \quad \Rightarrow \quad M(d) &= \max_{i \in B} \left\{ f^i + \langle g^i, \hat{x} + d - y^i \rangle \right\} \\ &= f(\hat{x}) + \max_{i \in B} \left\{ -e^i + \langle g^i, d \rangle \right\} \quad \text{(blind CP)} \end{aligned}$$

f_1

f_1

h

Some Model Functions

$$(y = \hat{x} + d)$$

$$\begin{aligned} f(y) \quad \Rightarrow \quad M(d) &= \max_{i \in B} \left\{ f^i + \langle g^i, \hat{x} + d - y^i \rangle \right\} \\ &= f(\hat{x}) + \max_{i \in B} \left\{ -e^i + \langle g^i, d \rangle \right\} \quad \text{(blind CP)} \end{aligned}$$

$$\begin{aligned} f_1(y) + \Psi(y) \quad \Rightarrow \quad M(d) &= \text{blind CP for } f_1 \\ &\quad + \Psi(\hat{x}) + \langle \nabla \Psi(\hat{x}), d \rangle \quad \text{(structured CP)} \end{aligned}$$

f_1

h

Some Model Functions

$$(y = \hat{x} + d)$$

$$\begin{aligned} f(y) \quad \Rightarrow \quad M(d) &= \max_{i \in B} \left\{ f^i + \langle g^i, \hat{x} + d - y^i \rangle \right\} \\ &= f(\hat{x}) + \max_{i \in B} \left\{ -e^i + \langle g^i, d \rangle \right\} \quad \text{(blind CP)} \end{aligned}$$

$$\begin{aligned} f_1(y) + \Psi(y) \quad \Rightarrow \quad M(d) &= \text{blind CP for } f_1 \\ &\quad + \Psi(\hat{x}) + \langle \nabla \Psi(\hat{x}), d \rangle \quad \text{(structured CP)} \\ \Rightarrow \quad M(d) &= \text{blind CP for } f_1 + \Psi(\hat{x}) + \langle \nabla \Psi(\hat{x}), d \rangle \\ &\quad + \frac{1}{2} \langle d, \nabla^2 \Psi(\hat{x}) d \rangle \quad \text{(structured nonpolyhedral)} \end{aligned}$$

h

Some Model Functions

$$(y = \hat{x} + d)$$

$$\begin{aligned} f(y) \quad \Rightarrow \quad M(d) &= \max_{i \in B} \left\{ f^i + \langle g^i, \hat{x} + d - y^i \rangle \right\} \\ &= f(\hat{x}) + \max_{i \in B} \left\{ -e^i + \langle g^i, d \rangle \right\} \quad (\mathbf{blind\ CP}) \end{aligned}$$

$$\begin{aligned} f_1(y) + \Psi(y) \quad \Rightarrow \quad M(d) &= \mathbf{blind\ CP\ for\ } f_1 \\ &\quad + \Psi(\hat{x}) + \langle \nabla \Psi(\hat{x}), d \rangle \quad (\mathbf{structured\ CP}) \\ \Rightarrow \quad M(d) &= \mathbf{blind\ CP\ for\ } f_1 + \Psi(\hat{x}) + \langle \nabla \Psi(\hat{x}), d \rangle \\ &\quad + \frac{1}{2} \langle d, \nabla^2 \Psi(\hat{x}) d \rangle \quad (\mathbf{structured\ nonpolyhedral}) \end{aligned}$$

$$(h \circ c)(y) \quad \Rightarrow \quad M(d) = h \circ (c(\hat{x}) + c'(\hat{x})d) \quad (\mathbf{composite\ nonpolyhedral})$$

h

Some Model Functions

$$(y = \hat{x} + d)$$

$$\begin{aligned} f(y) \quad \Rightarrow \quad M(d) &= \max_{i \in B} \left\{ f^i + \left\langle g^i, \hat{x} + d - y^i \right\rangle \right\} \\ &= f(\hat{x}) + \max_{i \in B} \left\{ -e^i + \left\langle g^i, d \right\rangle \right\} \quad \text{(blind CP)} \end{aligned}$$

$$\begin{aligned} f_1(y) + \Psi(y) \quad \Rightarrow \quad M(d) &= \text{blind CP for } f_1 \\ &\quad + \Psi(\hat{x}) + \left\langle \nabla \Psi(\hat{x}), d \right\rangle \quad \text{(structured CP)} \\ \Rightarrow \quad M(d) &= \text{blind CP for } f_1 + \Psi(\hat{x}) + \left\langle \nabla \Psi(\hat{x}), d \right\rangle \\ &\quad + \frac{1}{2} \left\langle d, \nabla^2 \Psi(\hat{x}) d \right\rangle \quad \text{(structured nonpolyhedral)} \end{aligned}$$

$$\begin{aligned} (h \circ c)(y) \quad \Rightarrow \quad M(d) &= h \circ (c(\hat{x}) + c'(\hat{x})d) \quad \text{(composite nonpolyhedral)} \\ M(d) &= \text{blind CP for } h \circ (c(\hat{x}) + c'(\hat{x})d) \\ &\quad \text{(composite polyhedral)} \end{aligned}$$

Some Model Functions

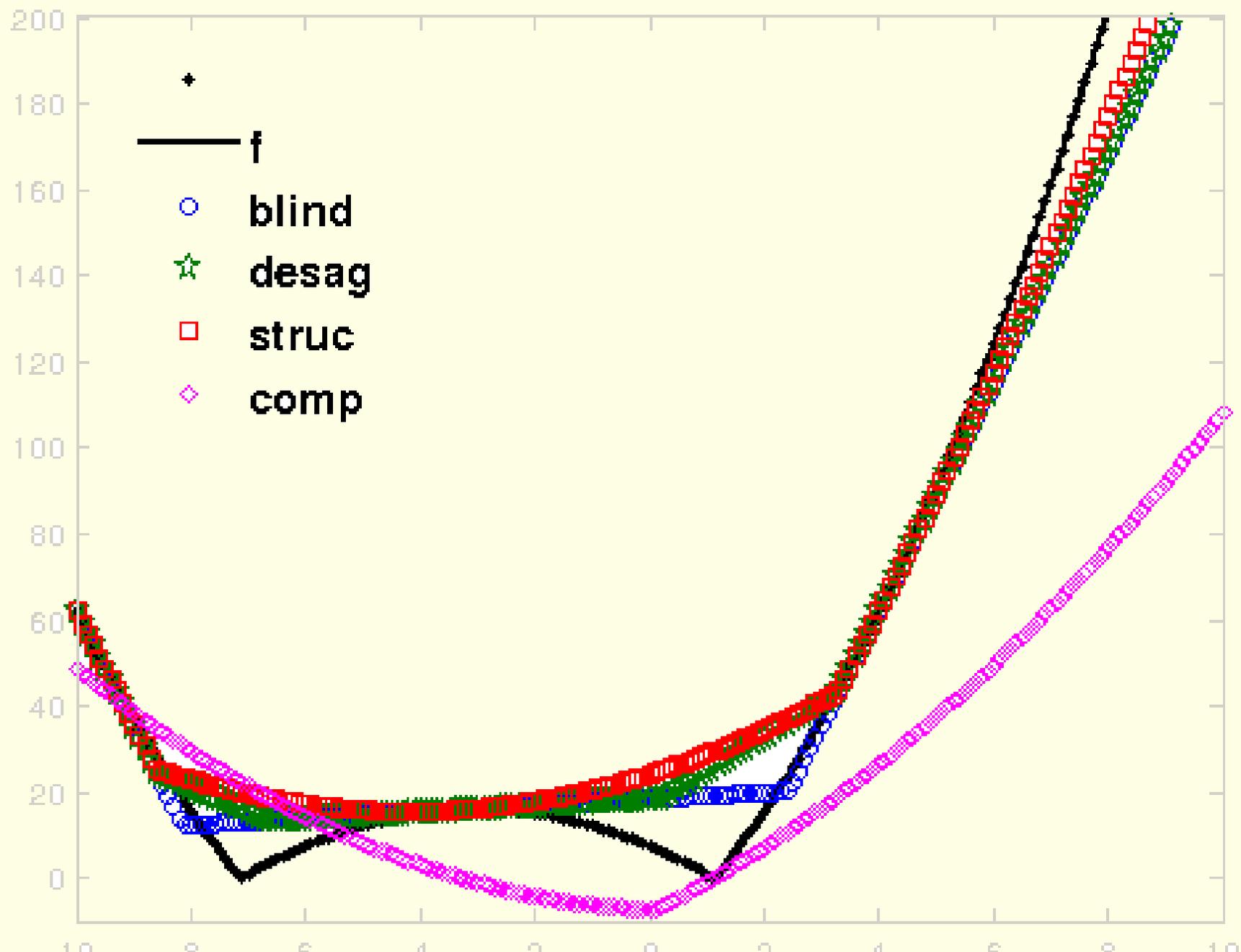
$$(y = \hat{x} + d)$$

$$\begin{aligned} f(y) \quad \Rightarrow \quad M(d) &= \max_{i \in B} \left\{ f^i + \left\langle g^i, \hat{x} + d - y^i \right\rangle \right\} \\ &= f(\hat{x}) + \max_{i \in B} \left\{ -e^i + \left\langle g^i, d \right\rangle \right\} \quad \text{(blind CP)} \end{aligned}$$

$$\begin{aligned} f_1(y) + \Psi(y) \quad \Rightarrow \quad M(d) &= \text{blind CP for } f_1 \\ &\quad + \Psi(\hat{x}) + \left\langle \nabla \Psi(\hat{x}), d \right\rangle \quad \text{(structured CP)} \\ \Rightarrow \quad M(d) &= \text{blind CP for } f_1 + \Psi(\hat{x}) + \left\langle \nabla \Psi(\hat{x}), d \right\rangle \\ &\quad + \frac{1}{2} \left\langle d, \nabla^2 \Psi(\hat{x}) d \right\rangle \quad \text{(structured nonpolyhedral)} \end{aligned}$$

$$\begin{aligned} (h \circ c)(y) \quad \Rightarrow \quad M(d) &= h \circ (c(\hat{x}) + c'(\hat{x})d) \quad \text{(composite nonpolyhedral)} \\ M(d) &= \text{blind CP for } h \circ (c(\hat{x}) + c'(\hat{x})d) \\ &\quad \text{(composite polyhedral)} \end{aligned}$$

or we can add second order terms for c, etc



Some Model Functions

$$(y = \hat{x} + d)$$

$$\begin{aligned} f(y) \quad \Rightarrow \quad M(d) &= \max_{i \in B} \left\{ f^i + \langle g^i, \hat{x} + d - y^i \rangle \right\} \\ &= f(\hat{x}) + \max_{i \in B} \left\{ -e^i + \langle g^i, d \rangle \right\} \quad (\text{blind CP}) \end{aligned}$$

$$\begin{aligned} f_1(y) + \Psi(y) \quad \Rightarrow \quad M(d) &= \text{blind CP for } f_1 \\ &\quad + \Psi(\hat{x}) + \langle \nabla \Psi(\hat{x}), d \rangle \quad (\text{structured CP}) \\ \Rightarrow \quad M(d) &= \text{blind CP for } f_1 + \Psi(\hat{x}) + \langle \nabla \Psi(\hat{x}), d \rangle \\ &\quad + \frac{1}{2} \langle d, \nabla^2 \Psi(\hat{x}) d \rangle \quad (\text{structured nonpolyhedral}) \end{aligned}$$

$$\begin{aligned} (h \circ c)(y) \quad \Rightarrow \quad M(d) &= h \circ (c(\hat{x}) + c'(\hat{x})d) \quad (\text{composite nonpolyhedral}) \\ M(d) &= \text{blind CP for } h \circ (c(\hat{x}) + c'(\hat{x})d) \\ &\quad (\text{composite polyhedral}) \end{aligned}$$

or we can add second order terms for c, etc

What are the least conditions we can set on M ??

Model and f at the center: lower type ($M(0) \leq f(\hat{x})$)

$$f(\hat{x}) \quad \Rightarrow \quad M(0) \quad = \quad \max_{i \in B} \left\{ f^i + \left\langle g^i, \hat{x} - y^i \right\rangle \right\}$$

$$= f(\hat{x}) + \max_{i \in B} \left\{ -e^i + \left\langle g^i, 0 \right\rangle \right\} \quad (\text{blind CP})$$

$$f_1(\hat{x}) + \Psi(\hat{x}) \quad \Rightarrow \quad M(0) \quad = \quad \text{blind CP for } f_1$$

$$+ \Psi(\hat{x}) + \left\langle \nabla \Psi(\hat{x}), 0 \right\rangle \quad (\text{structured CP})$$

$$\Rightarrow \quad M(0) \quad = \quad \text{blind CP for } f_1 + \Psi(\hat{x})$$

$$+ \frac{1}{2} \left\langle 0, \nabla^2 \Psi(\hat{x}) 0 \right\rangle \quad (\text{structured nonpolyhedral})$$

$$(h \circ c)(\hat{x}) \quad \Rightarrow \quad M(0) \quad = \quad h \circ (c(\hat{x}) + c'(\hat{x})0) \quad (\text{composite nonpolyhedral})$$

$$M(0) \quad = \quad \text{blind CP for } h \circ (c(\hat{x}) + c'(\hat{x})0)$$

(composite polyhedral)

or we can add second order terms for c , etc

Proximal Bundle Method iteration elements

For a prox-parameter $\mu > 0$ and a

model function $M(d) \approx f(\hat{x} + d)$ **near a center** \hat{x}

search direction $d(\hat{x}) := \arg \min M(\cdot) + \frac{1}{2}\mu|\cdot|^2,$

Proximal Bundle Method iteration elements

For a prox-parameter $\mu > 0$ and a
model function $M(d) \approx f(\hat{x} + d)$ near a **center** \hat{x}
search direction $d(\hat{x}) := \arg \min M(\cdot) + \frac{1}{2}\mu|\cdot|^2$,
model Gradient

$$G(\hat{x}) := -\mu d(\hat{x}) \in \partial M(d(\hat{x}))$$

Affine support for the model

$$A(d) := M(d(\hat{x})) + \langle G(\hat{x}), d - d(\hat{x}) \rangle$$

Error $E(\hat{x}) := M(0) - A(0)$

Proximal Bundle Method iteration elements

For a prox-parameter $\mu > 0$ and a
model function $M(d) \approx f(\hat{x} + d)$ near a center \hat{x}
search direction $d(\hat{x}) := \arg \min M(\cdot) + \frac{1}{2}\mu|\cdot|^2$,
model Gradient

$$G(\hat{x}) := -\mu d(\hat{x}) \in \partial M(d(\hat{x}))$$

Affine support for the model

$$A(d) := M(d(\hat{x})) + \langle G(\hat{x}), d - d(\hat{x}) \rangle$$

Error $E(\hat{x}) := M(0) - A(0)$

How general is this algorithmic scheme?

Various Bundle Methods iteration elements

A parameter $\mu = \mu(d(\hat{x})) > 0$ and a
model function $M(d) \approx f(\hat{x} + d)$ near a center \hat{x}
search direction and **model Gradient**

$$G(\hat{x}) := -\mu d(\hat{x}) \in \partial M(d(\hat{x}))$$

Affine support for the model

$$A(d) := M(d(\hat{x})) + \langle G(\hat{x}), d - d(\hat{x}) \rangle$$

Error $E(\hat{x}) := M(0) - A(0)$

Various Bundle Methods iteration elements

A parameter $\mu = \mu(d(\hat{x})) > 0$ and a
model function $M(d) \approx f(\hat{x} + d)$ near a center x
search direction and **model Gradient**

$$G(\hat{x}) := -\mu d(\hat{x}) \in \partial M(d(\hat{x}))$$

Affine support for the model

$$A(d) := M(d(\hat{x})) + \langle G(\hat{x}), d - d(\hat{x}) \rangle$$

Error $E(\hat{x}) := M(0) - A(0)$

Note 1: $E(\hat{x}) \geq 0$ by subgradient inequality (M is convex)

Various Bundle Methods iteration elements

A parameter $\mu = \mu(d(\hat{x})) > 0$ and a
model function $M(d) \approx f(\hat{x} + d)$ near a center x
search direction and **model Gradient**

$$G(\hat{x}) := -\mu d(\hat{x}) \in \partial M(d(\hat{x}))$$

Affine support for the model

$$A(d) := M(d(\hat{x})) + \langle G(\hat{x}), d - d(\hat{x}) \rangle$$

Error $E(\hat{x}) := M(0) - A(0)$

Note 1: $E(\hat{x}) \geq 0$ by subgradient inequality (M is convex)

Note 2: $d(\hat{x})$ may provide next iterate ($y_+ = \hat{x} + d(\hat{x})$),

or we may linesearch for next iterate ($y(t) = \hat{x} + td(\hat{x})$),

or we may curvsearch ($y(\mu) = \hat{x} + d_\mu(\hat{x})$), without changing M

Stationarity result

Model conditions

[M₀] The model is convex

[M₁] The model is lower at the center: $M(0) \leq f(\hat{x})$

[M₂] (finite version) When the error is zero, the model

gradient is an f -subgradient: $E(\hat{x}) = 0 \implies G(\hat{x}) \in \partial f(\hat{x})$

Stationarity result (finite version)

Model conditions

[M₀] The model is convex

[M₁] The model is lower at the center: $M(0) \leq f(\hat{x})$

[M₂] (finite version) When the error is zero, the model gradient is an f-subgradient: $E(\hat{x}) = 0 \implies G(\hat{x}) \in \partial f(\hat{x})$

Theorem For the **progress measure**

$$D(\hat{x}) := f(\hat{x}) - M(d(\hat{x})) - \frac{\mu}{2}|d(\hat{x})|^2$$

(i) $D(\hat{x}) = f(\hat{x}) - M(0) + E(\hat{x}) + \frac{1}{2\mu}|G(\hat{x})|^2$ and, hence,

$$M_1 \implies D(\hat{x}) \geq 0.$$

(ii) If M_{1,2} hold, then either $D(\hat{x}) = 0$ or $d(\hat{x}) = 0$ implies

$0 = G(\hat{x}) \in \partial f(\hat{x})$, i.e. \hat{x} is stationary for f .

Stationarity result (asymptotic)

Model conditions

[M₀] The model is convex

[M₁] The model is lower at the center: $M(0) \leq f(\hat{x})$

[M₂] (algorithmic version) When the error $\rightarrow 0$ the model gradient is an f -subgradient in the limit:

$$E(\hat{x}^k) \rightarrow 0 \implies G(\hat{x}^\infty) \in \partial f(\hat{x}^\infty)$$

Stationarity result (asymptotic)

Model conditions

[M₀] The model is convex

[M₁] The model is lower at the center: $M(0) \leq f(\hat{x})$

[M₂] (algorithmic version) When the error $\rightarrow 0$ the model gradient is an f -subgradient in the limit:

$$E(\hat{x}^k) \rightarrow 0 \implies G(\hat{x}^\infty) \in \partial f(\hat{x}^\infty)$$

Stationarity result (asymptotic)

Model conditions

[M₀] The model is convex

[M₁] The model is lower at the center: $M(0) \leq f(\hat{x})$

[M₂] (algorithmic version) When the error $\rightarrow 0$ the model gradient is an f -subgradient in the limit:

$$E(\hat{x}^k) \rightarrow 0 \implies G(\hat{x}^\infty) \in \partial f(\hat{x}^\infty)$$

For nonconvex f , condition M₂ ensured by **selecting**
B-points via a **linesearch** ($d(\hat{x}) \neq 0$),

Linesearch

A subalgorithm detecting if $y_\ell = \hat{x}^k + t_\ell d(\hat{x}^k)$ is either

- a serious step: $\hat{x}^{k+1} = y_\ell$ and black-box information at y_ℓ enters the model; or
- a null step: black-box information at y_ℓ enters the model; or
- adjusts t_ℓ to try with $y_{\ell+1}$.

Linesearch

A subalgorithm detecting if $y_\ell = \hat{x}^k + t_\ell d(\hat{x}^k)$ is either

- a serious step: $\hat{x}^{k+1} = y_\ell$ and black-box information at y_ℓ enters the model; or
- a null step: black-box information at y_ℓ enters the model; or
- adjusts t_ℓ to try with $y_{\ell+1}$.

Main features

1. For semismooth functions it has finite termination
2. Declares serious/null step checking the progress measure and also if **\mathcal{V} -approximation is good enough**
3. Null-step condition is the weakest in the market

Null Step Condition

Only requires that at fixed center \hat{x}

- $M_+(d) \geq A(d)$
- $M_+(d) \geq L(d) := f(\hat{x}) - e_\ell + \langle g_\ell, d \rangle$

for e_ℓ and g_ℓ defined with black-box information at y_ℓ

Null Step Condition

Only requires that at fixed center \hat{x}

- $M_+(d) \geq A(d)$
- $M_+(d) \geq L(d) := f(\hat{x}) - e_\ell + \langle g_\ell, d \rangle$

for e_ℓ and g_ℓ defined with black-box information at y_ℓ

Note 1: if f convex, $g_\ell = g(y_\ell)$ and e_ℓ linearization error

Null Step Condition

Only requires that at fixed center \hat{x}

- $M_+(d) \geq A(d)$
- $M_+(d) \geq L(d) := f(\hat{x}) - e_\ell + \langle g_\ell, d \rangle$

for e_ℓ and g_ℓ defined with black-box information at y_ℓ

Note 1: if f convex, $g_\ell = g(y_\ell)$ and e_ℓ linearization error

Note 2: this does not mean M is **polyhedral**

Null Step Condition

Only requires that at fixed center \hat{x}

- $M_+(d) \geq A(d)$
- $M_+(d) \geq L(d) := f(\hat{x}) - e_\ell + \langle g_\ell, d \rangle$

for e_ℓ and g_ℓ defined with black-box information at y_ℓ

Note 1: if f convex, $g_\ell = g(y_\ell)$ and e_ℓ linearization error

Note 2: this does not mean M is **polyhedral**

Note 3: if f nonconvex, define g_ℓ and e_ℓ checking

curvature

Null Step Condition

Only requires that at fixed center \hat{x}

- $M_+(d) \geq A(d)$
- $M_+(d) \geq L(d) := f(\hat{x}) - e_\ell + \langle g_\ell, d \rangle$

for e_ℓ and g_ℓ defined with black-box information at y_ℓ

Note 1: if f convex, $g_\ell = g(y_\ell)$ and e_ℓ linearization error

Note 2: this does not mean M is **polyhedral**

Note 3: if f nonconvex, define g_ℓ and e_ℓ checking

curvature: now the bundle $B = \bigcup_i (y^i, f^i, g^i, H^i)$

includes matrices H^i

the curvature is $\kappa_\ell := \langle y_\ell - \hat{x}, H_\ell(y_\ell - \hat{x}) \rangle$

Curvature (nonconvex f , polyhedral M)

Having the bundle $B = \bigcup_i (y^i, f^i, g^i, H^i)$ and

$\kappa_\ell := \langle y_\ell - \hat{x}, H_\ell(y_\ell - \hat{x}) \rangle$, define

$$g_\ell = \begin{cases} g(y_\ell) & \text{if } \kappa_\ell \geq 0 \\ g(y_\ell) + H_\ell(y_\ell - \hat{x}) & \text{if } \kappa_\ell < 0 \end{cases}$$

Curvature (nonconvex f , polyhedral M)

Having the bundle $B = \bigcup_i (y^i, f^i, g^i, H^i)$ and

$\kappa_\ell := \langle y_\ell - \hat{x}, H_\ell(y_\ell - \hat{x}) \rangle$, define

$$g_\ell = \begin{cases} g(y_\ell) & \text{if } \kappa_\ell \geq 0 \\ g(y_\ell) + H_\ell(y_\ell - \hat{x}) & \text{if } \kappa_\ell < 0 \end{cases}$$

and

$$e_\ell = \begin{cases} \text{lin.error} & \text{if } \kappa_\ell \geq 0 \\ \max\left(\text{lin.error} - \frac{1}{2}\kappa_\ell, s|y_\ell - \hat{x}|^2\right) & \text{if } \kappa_\ell < 0 \end{cases}$$

Curvature (nonconvex f , polyhedral M)

Having the bundle $B = \bigcup_i (y^i, f^i, g^i, H^i)$ and

$\kappa_\ell := \langle y_\ell - \hat{x}, H_\ell(y_\ell - \hat{x}) \rangle$, define

$$g_\ell = \begin{cases} g(y_\ell) & \text{if } \kappa_\ell \geq 0 \\ g(y_\ell) + H_\ell(y_\ell - \hat{x}) & \text{if } \kappa_\ell < 0 \end{cases}$$

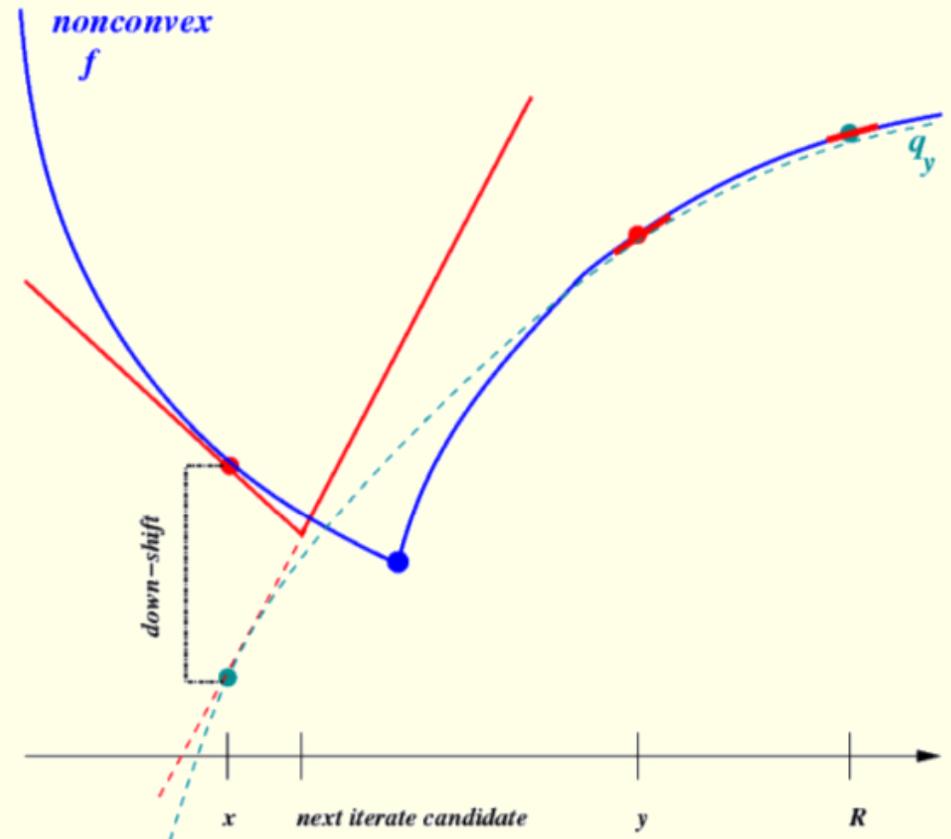
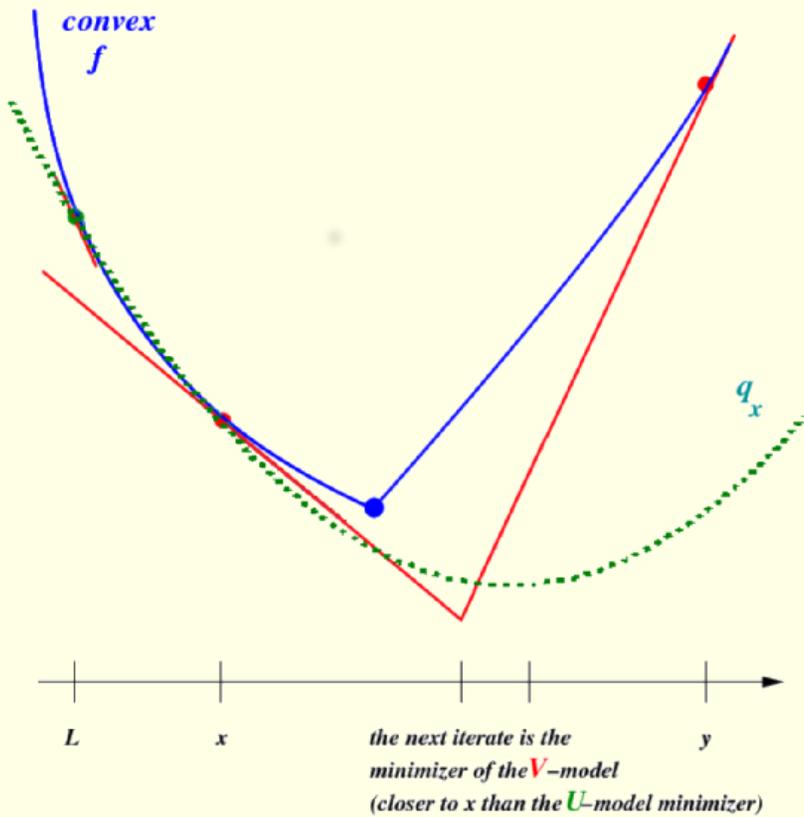
and

$$e_\ell = \begin{cases} \text{lin.error} & \text{if } \kappa_\ell \geq 0 \\ \max\left(\text{lin.error} - \frac{1}{2}\kappa_\ell, s|y_\ell - \hat{x}|^2\right) & \text{if } \kappa_\ell < 0 \end{cases}$$

Note: Safeguard s ensures $e_\ell \geq 0$ and, hence,
 $M(0) \leq f(\hat{x})$

Curvature (nonconvex f , polyhedral M)

Bundle matrices H computed via quadratic approximations



Final Comments

- Brings a model-view into the nonconvex setting
- Weakest NS condition relies heavily on conjugacy and Convex Analysis (\neq Correa and Lemaréchal NS analysis)
- The linesearch extends the one-dimensional superlinearly convergent method by Lemaréchal and Mifflin 1982

This is difficult stuff

Final Comments

- Brings a model-view into the nonconvex setting
- Weakest NS condition relies heavily on conjugacy and Convex Analysis (\neq Correa and Lemaréchal NS analysis)
- The linesearch extends the one-dimensional superlinearly convergent method by Lemaréchal and Mifflin 1982
- **This is difficult stuff**