

**Nonsmooth optimization:**  
**beyond first order methods.**

**A tutorial**  
**focusing on**  
**bundle methods**

**Claudia Sagastizábal**

**(IMECC-UniCamp, Campinas Brazil, adjunct researcher)**

**SESO 2018, Paris, May 23 and 25, 2018**

# Computational NSO: what do we mean?

For the unconstrained problem

$$\min f(x),$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex but not differentiable at some points

Algorithms defined according on **how much** information is provided by certain oracle

# Computational NSO: what do we mean?

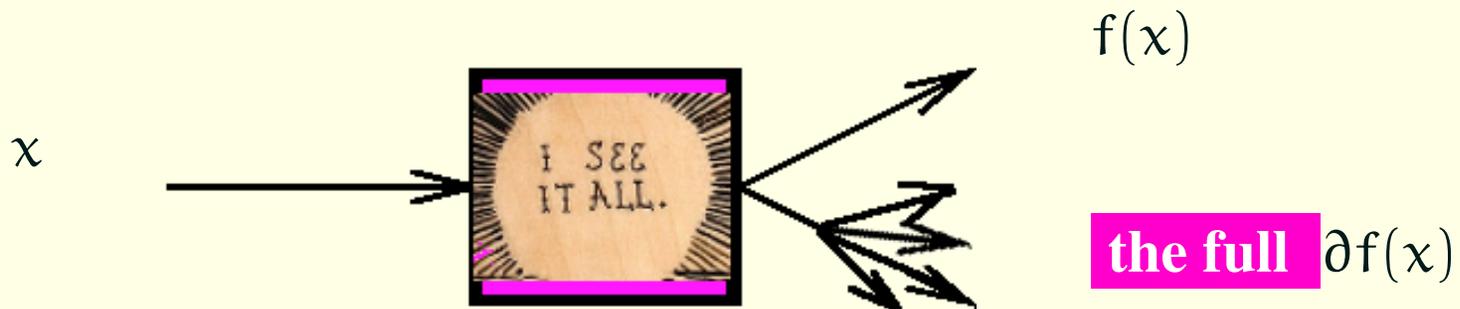
For the unconstrained problem

$$\min f(x),$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex but not differentiable at some points,

Algorithms defined according on **how much** information is provided by certain oracle

an **informative oracle**



# Computational NSO: what do we mean?

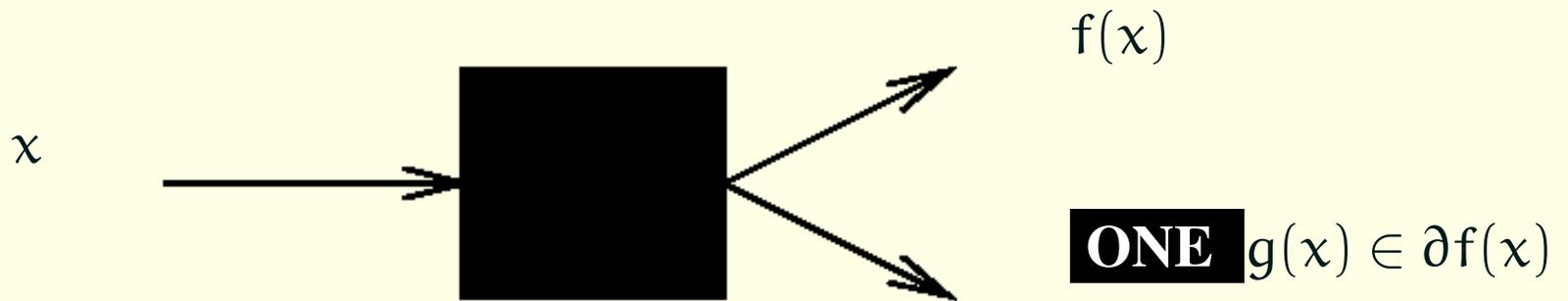
For the unconstrained problem

$$\min f(x),$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex but not differentiable at some points

Algorithms defined according on **how much** information is provided by certain oracle

a **“black box”**



# Computational NSO: what do we mean?

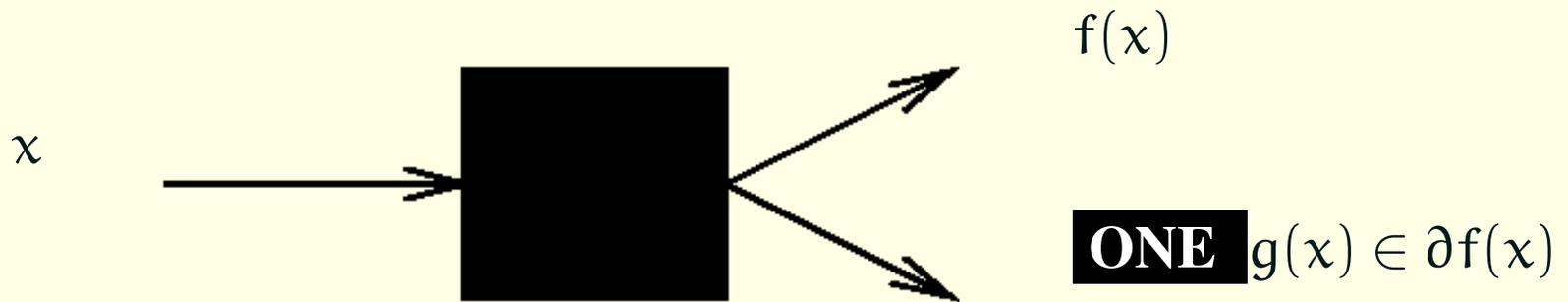
For the unconstrained problem

$$\min f(x),$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex but not differentiable at some points,

Algorithms defined according on **how much** information is provided by certain oracle

a **“black box”**



**How common are nonsmooth objective functions in optimization?**

## When does nonsmoothness appear?

- \* if the **nature** of the problem imposes a nonsmooth model; or

- \* if **sparsity** of the solution is a concern; or

- \* in problems difficult to solve,

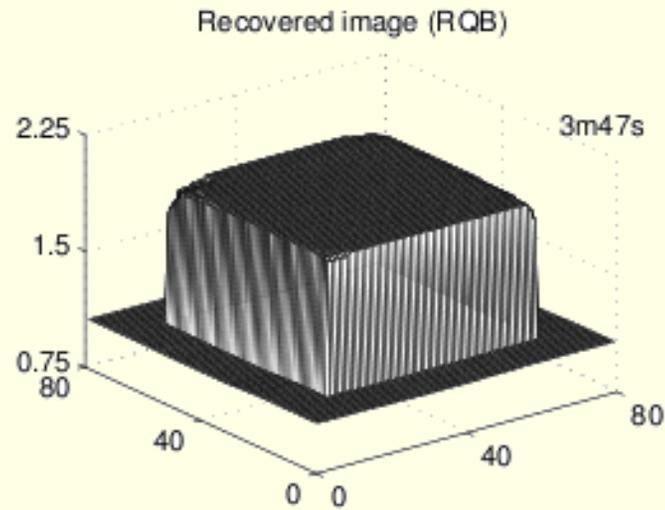
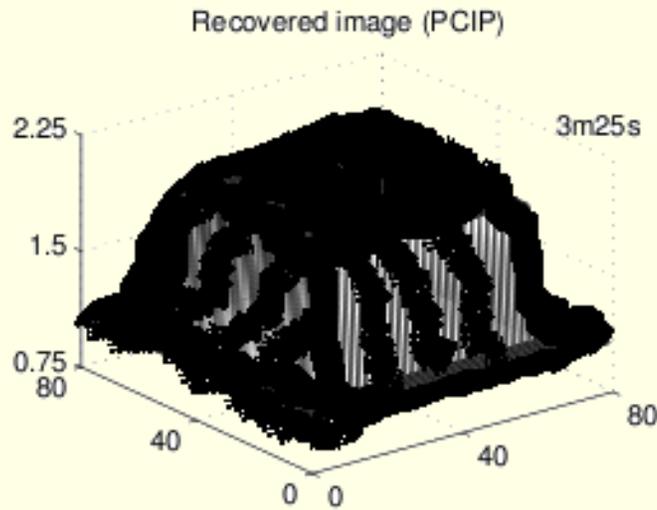
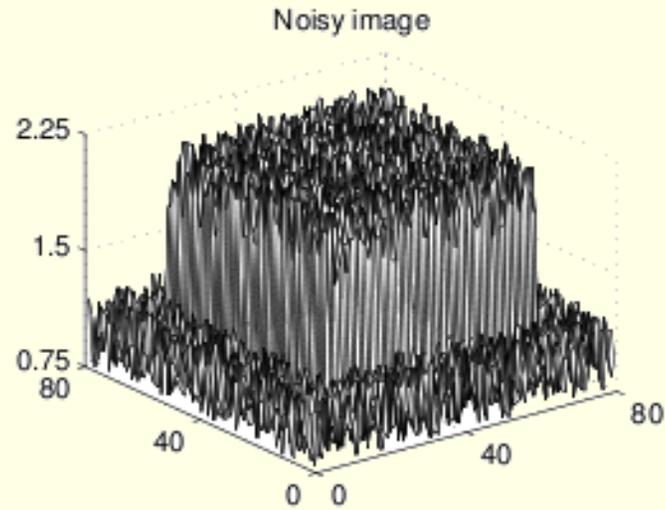
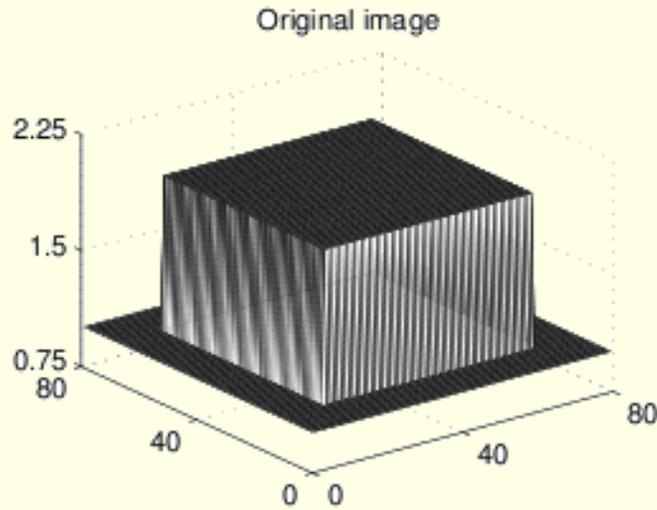
  - because they are large scale

  - because they are heterogeneous

sometimes the **solution method** induces nonsmoothness

# Example of NS model

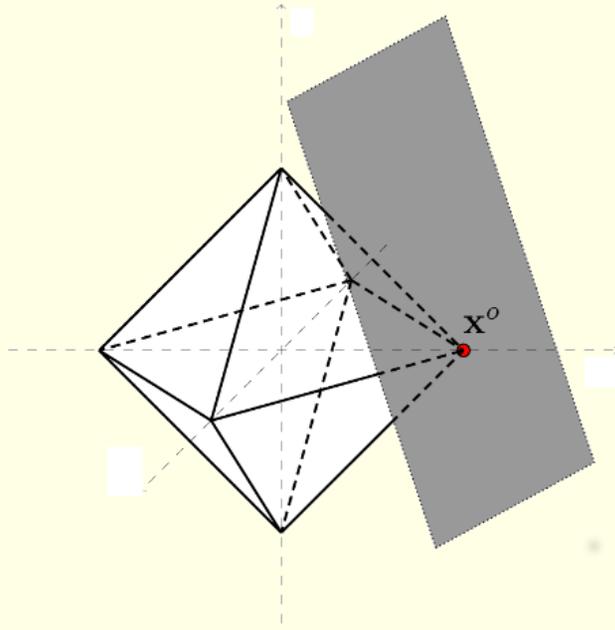
Recovery of **blocky** images ( $\ell_1$ -regularization of TV)



# Example of sparse optimization $\min\{\|x\|_1 : Ax = b\}$

**Basis pursuit:** find least 1-norm point on the affine plane

Tends to return a sparse point (sometimes, the sparsest)

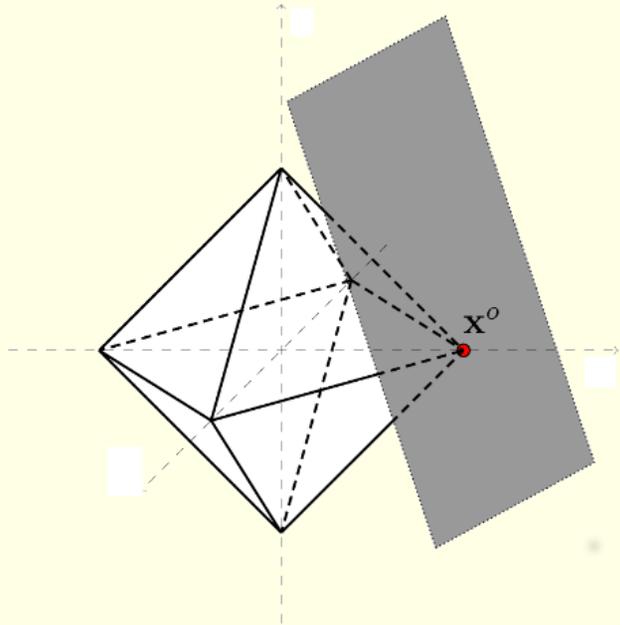


$\ell_1$  ball touches the affine plane

# Example of sparse optimization $\min\{\|x\|_1 : Ax = b\}$

**Basis pursuit:** find least 1-norm point on the affine plane

Tends to return a sparse point (sometimes, the sparsest)



$l_1$  ball touches the affine plane

**LASSO** denoises basis pursuit

$$\min \left\{ \|Ax - b\|_2^2 : \|x\|_1 \leq \tau \right\}$$

or

$$\min \left\{ \|x\|_1 + \frac{\mu}{2} \|Ax - b\|_2^2 \right\}$$

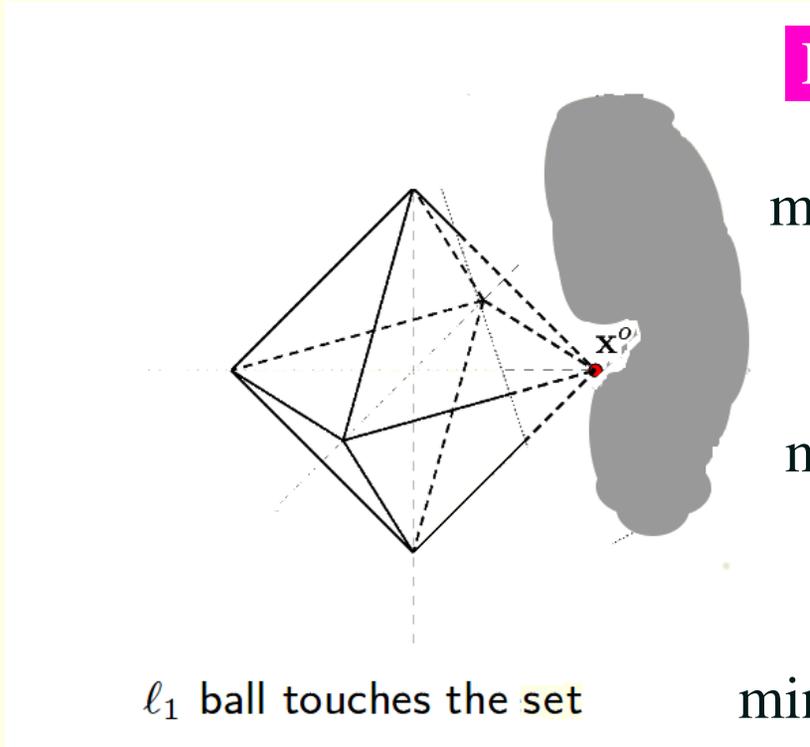
or

$$\min \left\{ \|x\|_1 : \|Ax - b\|_2^2 \leq \sigma \right\}$$

# Example of sparse optimization $\min\{\|\mathbf{x}\|_1 : \mathbf{h}(\mathbf{x}) \leq \mathbf{b}\}$

**Basis pursuit:** find least 1-norm point on a **nonlinear set**

Tends to return a sparse point (sometimes, the sparsest)



**LASSO** denoises basis pursuit

$$\min \left\{ \left\| \left( \mathbf{h}(\mathbf{x}) - \mathbf{b} \right)^+ \right\|_2^2 : \|\mathbf{x}\|_1 \leq \tau \right\}$$

or

$$\min \left\{ \|\mathbf{x}\|_1 + \frac{\mu}{2} \left\| \left( \mathbf{h}(\mathbf{x})\mathbf{x} - \mathbf{b} \right)^+ \right\|_2^2 \right\}$$

or

$$\min \left\{ \|\mathbf{x}\|_1 : \left\| \left( \mathbf{h}(\mathbf{x})\mathbf{x} - \mathbf{b} \right)^+ \right\|_2^2 \leq \sigma \right\}$$

# Lagrangian Relaxation Example

Real-life optimization problems

$$\text{(primal)} \quad \left\{ \begin{array}{l} \min \sum_{j \in J} c^j(p^j) \\ \text{for } j \in J: p^j \in \mathcal{P}^j \\ \sum_{j \in J} g^j(p^j) = \text{Dem} \end{array} \right.$$

# Lagrangian Relaxation Example

Real-life optimization problems

$$\text{(primal)} \quad \left\{ \begin{array}{l} \max \quad \sum_{j \in J} -c^j(p^j) \\ \text{for } j \in J: p^j \in \mathcal{P}^j \\ \sum_{j \in J} g^j(p^j) = \text{Dem} \end{array} \right. \quad \leftarrow x$$

# Lagrangian Relaxation Example

Real-life optimization problems

$$\text{(primal)} \quad \left\{ \begin{array}{l} \max \quad \sum_{j \in J} -c^j(p^j) \\ \text{for } j \in J: p^j \in \mathcal{P}^j \\ \sum_{j \in J} g^j(p^j) = \text{Dem} \quad \leftarrow x \end{array} \right.$$

often exhibit separable structure passing to the (dual) :

# Lagrangian Relaxation Example

Real-life optimization problems

$$\text{(primal)} \quad \left\{ \begin{array}{l} \max \quad \sum_{j \in J} -c^j(p^j) \\ \text{for } j \in J: p^j \in \mathcal{P}^j \\ \sum_{j \in J} g^j(p^j) = \text{Dem} \quad \leftarrow x \end{array} \right.$$

often exhibit separable structure passing to the (dual) :

$$\min_x f(x) := f_0(x) + \sum_{j \in J} f^j(x)$$

$$\min_x -\langle x, \text{Dem} \rangle + \sum_{j \in J} \left\{ \begin{array}{l} \max_{p^j \in \mathcal{P}^j} -c^j(p^j) + \langle x, g^j(p^j) \rangle \end{array} \right.$$

# Lagrangian Relaxation Example

Real-life optimization problems

$$\text{(primal)} \quad \left\{ \begin{array}{l} \max \quad \sum_{j \in J} -c^j(p^j) \\ \text{for } j \in J: p^j \in \mathcal{P}^j \\ \sum_{j \in J} g^j(p^j) = \text{Dem} \quad \leftarrow x \end{array} \right.$$

often exhibit separable structure passing to the (dual) :

$$\min_x f(x) := f_0(x) + \sum_{j \in J} f^j(x)$$

$$\min_x -\langle x, \text{Dem} \rangle + \sum_{j \in J} \left\{ \begin{array}{l} \max_{p^j \in \mathcal{P}^j} -c^j(p^j) + \langle x, g^j(p^j) \rangle \end{array} \right.$$

# Benders Decomposition Example

Similar situation, but now the uncoupling is done on a primal level

$$\text{(primal)} \quad \left\{ \begin{array}{l} \min \sum_{j \in J} \mathcal{I}^j(\Delta p^j) + \mathcal{C}^j(p^j) \\ \text{for } j \in J: p^j \in \mathcal{P}^j \\ \Delta p \in D \end{array} \right. \iff p^j \leq \bar{p}^j + \Delta p^j$$

# Benders Decomposition Example

Similar situation, but now the uncoupling is done on a primal level

$$\text{(primal)} \left\{ \begin{array}{l} \min \sum_{j \in J} \mathcal{I}^j(\Delta p^j) + \mathcal{C}^j(p^j) \\ \text{for } j \in J: p^j \in \mathcal{P}^j \\ \Delta p \in D \end{array} \right. \iff p^j \leq \bar{p}^j + \Delta p^j$$

$$\left\{ \begin{array}{l} \min_{\Delta p} \sum_{j \in J} \mathcal{I}^j(\Delta p^j) + \mathcal{V}^j(\Delta p^j) \\ \Delta p \in D \end{array} \right. \quad \mathcal{V}^j(\Delta p^j) := \left\{ \begin{array}{l} \min \mathcal{C}^j(p^j) \\ p^j \leq \bar{p}^j + \Delta p^j \end{array} \right.$$

$$\min f(x) := \sum_{j \in J} f^j(\Delta p^j) \quad f^j(\Delta p^j) :=$$

# Benders Decomposition Example

Similar situation, but now the uncoupling is done on a primal level

$$(\text{primal}) \quad \left\{ \begin{array}{l} \min \sum_{j \in J} \mathcal{I}^j(\Delta p^j) + \mathcal{C}^j(p^j) \\ \text{for } j \in J: p^j \in \mathcal{P}^j \\ \Delta p \in D \end{array} \right. \iff p^j \leq \bar{p}^j + \Delta p^j$$

$$\left\{ \begin{array}{l} \min_{\Delta p} \sum_{j \in J} \mathcal{I}^j(\Delta p^j) + \mathcal{V}^j(\Delta p^j) \\ \Delta p \in D \end{array} \right. \quad \mathcal{V}^j(\Delta p^j) := \left\{ \begin{array}{l} \min \mathcal{C}^j(p^j) \\ p^j \leq \bar{p}^j + \Delta p^j \end{array} \right.$$

$$\min f(x) := \sum_{j \in J} f^j(\Delta p^j) \quad \text{for } f^j(\Delta p^j) := \mathcal{I}^j(\Delta p^j) + \mathcal{V}^j(\Delta p^j)$$

## Computing $\partial f(x^k)$ : how difficult is it?

1.  $f(x) = |x|$ , for  $n = 1$
2. A linear Lasso function,  $f(x) = \|x\|_1 + \frac{\mu}{2} \|Ax - b\|_2^2$
3. A nonlinear Lasso function,  $h \in C^1$ ,  
$$f(x) = \|x\|_1 + \frac{\mu}{2} \left\| \left( h(x) - b \right)^+ \right\|_2^2$$
4. One of the local subproblems in the Lagrangian example,  
$$f^j(x^k) := \begin{cases} \max & -\mathcal{C}^j(p^j) + \langle x^k, g^j(p^j) \rangle \\ & p^j \in \mathcal{P}^j \end{cases}$$
5. One of the local subproblems in the Benders example,  
$$(\mathcal{I}^j(\Delta p^j) + \mathcal{V}^j(\Delta p^j)) = f^j(x^{k,j}) = \min \left\{ \mathcal{C}^j(p^j) : p^j \leq \bar{p}^j + x^{k,j} \right\}$$

But why would one want **ALL** of  $\partial f(x^k)$ ?

**But why would one want ALL of  $\partial f(x^k)$ ?**

Indispensible to calculate the proximal point

$$p = \text{prox}_t^f(x) \iff p = \arg \min f(y) + \frac{1}{2t} \|y - x\|_2^2$$

$$\iff \in \partial f(p) + \frac{1}{t}(p - x)$$

$$\iff$$

**But why would one want ALL of  $\partial f(x^k)$ ?**

Indispensable to calculate the proximal point

$$p = \text{prox}_t^f(x) \iff p = \arg \min f(y) + \frac{1}{2t} \|y - x\|_2^2$$

$$\iff 0 \in \partial f(p) + \frac{1}{t}(p - x)$$



**But why would one want ALL of  $\partial f(x^k)$ ?**

Indispensable to calculate the proximal point

$$p = \text{prox}_t^f(x) \iff p = \arg \min f(y) + \frac{1}{2t} \|y - x\|_2^2$$

$$\iff 0 \in \partial f(p) + \frac{1}{t}(p - x)$$

$$\iff \frac{1}{t}(x - p) \in \partial f(p)$$

But why would one want **ALL** of  $\partial f(\mathbf{x}^k)$ ?

Indispensable to calculate the proximal point

$$\mathbf{p} = \text{prox}_t^f(\mathbf{x}) \iff \mathbf{p} = \arg \min \mathbf{y} \left( f(\mathbf{y}) + \frac{1}{2t} \|\mathbf{y} - \mathbf{x}\|_2^2 \right)$$

$$\iff 0 \in \partial f(\mathbf{p}) + \frac{1}{t}(\mathbf{p} - \mathbf{x})$$

$$\iff \frac{1}{t}(\mathbf{x} - \mathbf{p}) \in \partial \mathbf{f}(\mathbf{p})$$

Without full knowledge of the subdifferential, the

**implicit** inclusion cannot be solved!

But why would one want **ALL** of  $\partial f(x^k)$ ?

Indispensable to calculate the proximal point

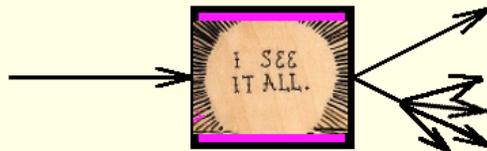
$$p = \text{prox}_t^f(x) \iff p = \arg \min f(y) + \frac{1}{2t} \|y - x\|_2^2$$

$$\iff 0 \in \partial f(p) + \frac{1}{t}(p - x)$$

$$\iff \frac{1}{t}(x - p) \in \partial f(p)$$

Without full knowledge of the subdifferential, the

**implicit** inclusion cannot be solved!



note:  $p \in x - t\partial f(p)$  akin to a subgradient method

# Proximal point algorithms (Accel. Nesterov, FISTA, AugLag)

$$\mathbf{x}^{k+1} = \text{prox}_{t_k}^f(\mathbf{x}^k)$$

$$\iff$$

$$\mathbf{x}^{k+1} = \arg \min \mathbf{y} \left( f(\mathbf{y}) + \frac{1}{2t_k} \|\mathbf{y} - \mathbf{x}^k\|_2^2 \right)$$

## Proximal point algorithms (Accel. Nesterov, FISTA, AugLag)

$$\mathbf{x}^{k+1} = \text{prox}_{t_k}^f(\mathbf{x}^k)$$

$$\iff$$

$$\mathbf{x}^{k+1} = \arg \min \mathbf{y} \left( f(\mathbf{y}) + \frac{1}{2t_k} \|\mathbf{y} - \mathbf{x}^k\|_2^2 \right)$$

- of interest if computing  $\text{prox}_{t_k}^f(\mathbf{x}^k)$  is much easier than minimizing  $f$
- stepsize  $t_k > 0$  impacts on the number of iterations

## Proximal point algorithms (Accel. Nesterov, FISTA, AugLag)

$$\mathbf{x}^{k+1} = \text{prox}_{t_k}^f(\mathbf{x}^k)$$

$$\iff$$

$$\mathbf{x}^{k+1} = \arg \min \mathbf{y} \left( f(\mathbf{y}) + \frac{1}{2t_k} \|\mathbf{y} - \mathbf{x}^k\|_2^2 \right)$$

- of interest if **computing**  $\text{prox}_{t_k}^f(\mathbf{x}^k)$  is much easier than minimizing  $f$
- stepsize  $t_k > 0$  impacts on the number of iterations

## Proximal point: calculus rules

- separable sum:

$$f(x, y) = (g(x), h(y)) \implies \\ \text{prox}_t^f(x) = \left( \text{prox}_t^g(x), \text{prox}_t^h(y) \right)$$

- scalar factor ( $\alpha \neq 0$ ) and translation ( $v \neq 0$ ):

$$f(x) = g(\alpha x + v) \implies \\ \text{prox}_t^f(x) = \frac{1}{\alpha} \left( \text{prox}_t^{\alpha^2 g}(\alpha x + v) - v \right)$$

- “perspective” ( $\alpha > 0$ ):

$$f(x) = \alpha g\left(\frac{1}{\alpha}x\right) \implies \text{prox}_t^f(x) = \alpha \text{prox}_t^{g/\alpha}\left(\frac{x}{\alpha}\right)$$

## Proximal point: special functions

- + linear term ( $v \neq 0$ ):

$$f(x) = g(x) + \langle v, x \rangle \implies \text{prox}_t^f(x) = \text{prox}_t^g(x - v)$$

- + convex quadratic term ( $t > 0$ ):

$$f(x) = g(x) + \frac{1}{2t} \|x - v\|^2 \implies$$

$$\text{prox}_t^f(x) = \text{prox}_t^{\lambda g}(\lambda x + (1 - \lambda)v) \text{ for } \lambda = \frac{t}{t + 1}$$

- composition with linear term such that  $A^\top A = \frac{1}{\alpha} I$ , ( $\alpha \neq 0$ ):

$$f(x) = g(Ax + v) \implies$$

$$\text{prox}_t^f(x) = (I - \alpha A^\top A)x + \alpha A^\top \left[ \text{prox}_t^{g/\alpha}(Ax + v) - v \right]$$

## Proximal point algorithm: convergence

If  $\arg \min f \neq \emptyset$  then

$$f(x^k) - f(\bar{x}) \leq \frac{\|x^0 - \bar{x}\|^2}{2 \sum_{i=1}^k t_i}$$

## Proximal point algorithm: convergence

If  $\arg \min f \neq \emptyset$  then

$$f(x^k) - f(\bar{x}) \leq \frac{\|x^0 - \bar{x}\|^2}{2 \sum_{i=1}^k t_i}$$

$\implies$  convergence if  $\sum t_i \rightarrow +\infty$

$\implies$  rate  $1/k$  if  $\{t_k\}$  bounded away from zero

## Proximal point algorithm: acceleration

$$\mathbf{x}^{k+1} = \text{prox}_{t_k}^f \left( \mathbf{x}^k + \theta_{k+1} \left( \frac{1}{\theta_k} - \mathbf{1} \right) (\mathbf{x}^k - \mathbf{x}^{k-1}) \right)$$

for

$$\frac{\theta_{k+1}^2}{t_{k+1}} = (1 - \theta_{k+1}) \frac{\theta_k^2}{t_k}$$

## Proximal point algorithm: acceleration

$$x^{k+1} = \text{prox}_{t_k}^f \left( x^k + \theta_{k+1} \left( \frac{1}{\theta_k} - \mathbf{1} \right) (x^k - x^{k-1}) \right)$$

for

$$\frac{\theta_{k+1}^2}{t_{k+1}} = (1 - \theta_{k+1}) \frac{\theta_k^2}{t_k}$$

$\implies$  convergence if  $\sum \sqrt{t_i} \rightarrow +\infty$

$\implies$  rate  $1/k^2$  if  $\{t_k\}$  bounded away from zero

**What if  $\text{prox}_t^f$  is not computable?**

**What if  $\text{prox}_t^f$  is not computable?**

**Use bundle methods!**

**What if  $\text{prox}_t^f$  is not computable?**

**Use bundle methods!**

**When do bundle method prove most useful?**

**What if  $\text{prox}_t^f$  is not computable?**

**Use bundle methods!**

**When do bundle method prove most useful?**

In situations

– when the objective function is not available explicitly

**and/or**

– when we do not have access to the full subdifferential

**and/or**

– when calculations need to be done with high precision

## Bundling to approximate the prox

$$\begin{aligned} \text{WANT: } p = \text{prox}_t^f(x) &\iff p = \arg \min f(y) + \frac{1}{2t} \|y - x\|_2^2 \\ &\iff 0 \in \partial f(p) + \frac{1}{t}(p - x) \\ &\iff \frac{1}{t}(x - p) \in \partial f(p) \end{aligned}$$

## Bundling to approximate the prox

$$\text{WANT: } p = \text{prox}_t^f(x) \iff p = \arg \min f(y) + \frac{1}{2t} \|y - x\|_2^2$$

$$\iff 0 \in \partial f(p) + \frac{1}{t}(p - x)$$

$$\iff \frac{1}{t}(x - p) \in \partial f(p)$$

$$\text{HAVE: } q = \text{prox}_t^{\mathbf{M}}(x) \iff q = \arg \min \mathbf{M}(y) + \frac{1}{2t} \|y - x\|_2^2$$

$$\iff 0 \in \partial \mathbf{M}(q) + \frac{1}{t}(q - x)$$

$$\iff \frac{1}{t}(x - q) \in \partial \mathbf{M}(q)$$

## Bundling to approximate the prox

$$\text{WANT: } p = \text{prox}_t^f(x) \iff p = \arg \min f(y) + \frac{1}{2t} \|y - x\|_2^2$$

$$\iff 0 \in \partial f(p) + \frac{1}{t}(p - x)$$

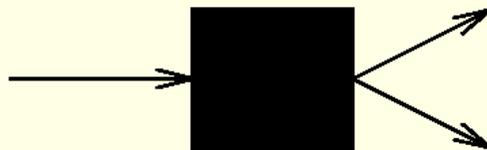
$$\iff \frac{1}{t}(x - p) \in \partial f(p)$$

$$\text{HAVE: } q = \text{prox}_t^{\mathbf{M}}(x) \iff q = \arg \min \mathbf{M}(y) + \frac{1}{2t} \|y - x\|_2^2$$

$$\iff 0 \in \partial \mathbf{M}(q) + \frac{1}{t}(q - x)$$

$$\iff \frac{1}{t}(x - q) \in \partial \mathbf{M}(q)$$

**M** is a model of  $f$  for which we do have full knowledge of the subdifferential: the **implicit** inclusion can be solved!



## Bundling to approximate the prox

$$\text{WANT: } p = \text{prox}_t^f(x) \iff p = \arg \min f(y) + \frac{1}{2t} \|y - x\|_2^2$$

$$\iff 0 \in \partial f(p) + \frac{1}{t}(p - x)$$

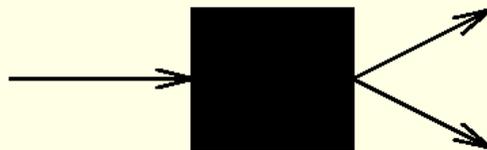
$$\iff \frac{1}{t}(x - p) \in \partial f(p)$$

$$\text{HAVE: } q = \text{prox}_t^{\mathbf{M}}(x) \iff q = \arg \min \mathbf{M}(y) + \frac{1}{2t} \|y - x\|_2^2$$

$$\iff 0 \in \partial \mathbf{M}(q) + \frac{1}{t}(q - x)$$

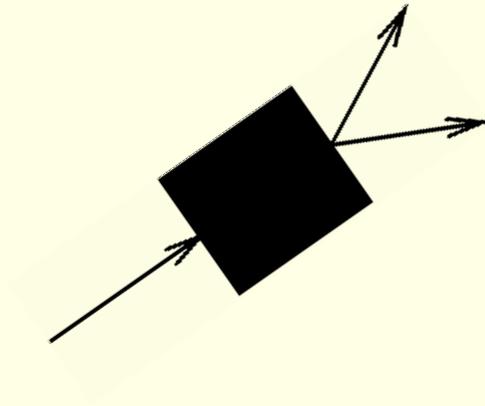
$$\iff \frac{1}{t}(x - q) \in \partial \mathbf{M}(q)$$

**M** is a model of  $f$  for which we do have full knowledge of the subdifferential: the **implicit** inclusion can be solved!



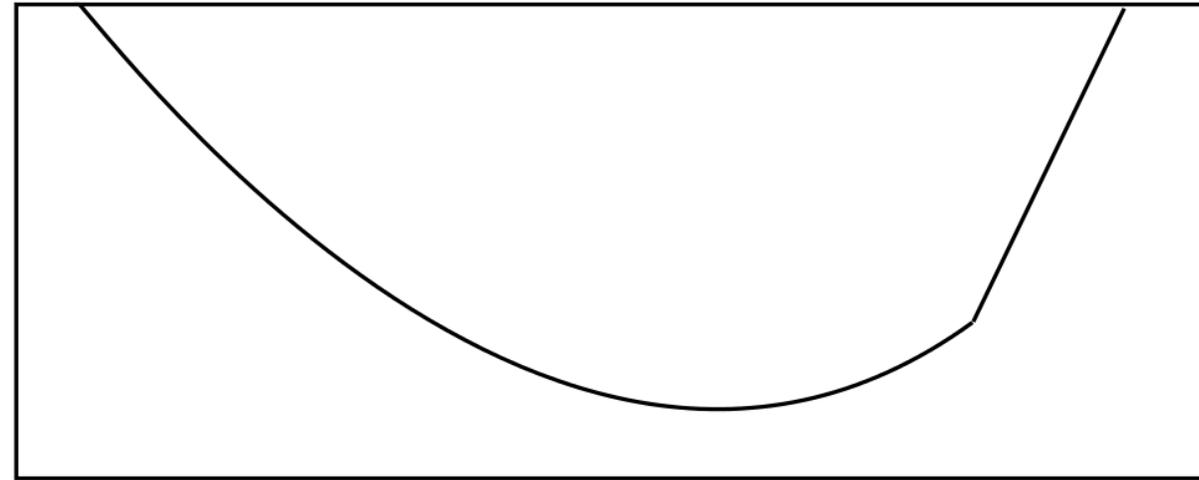
How is the model built?

# Model built with the black box



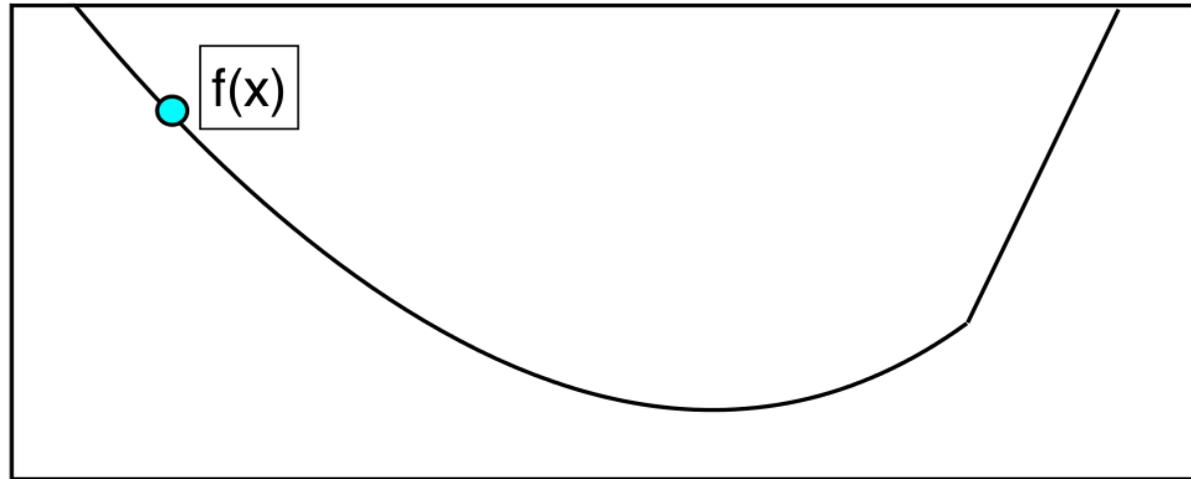
# A quick overview of Convex Analysis

An example of a convex nonsmooth function



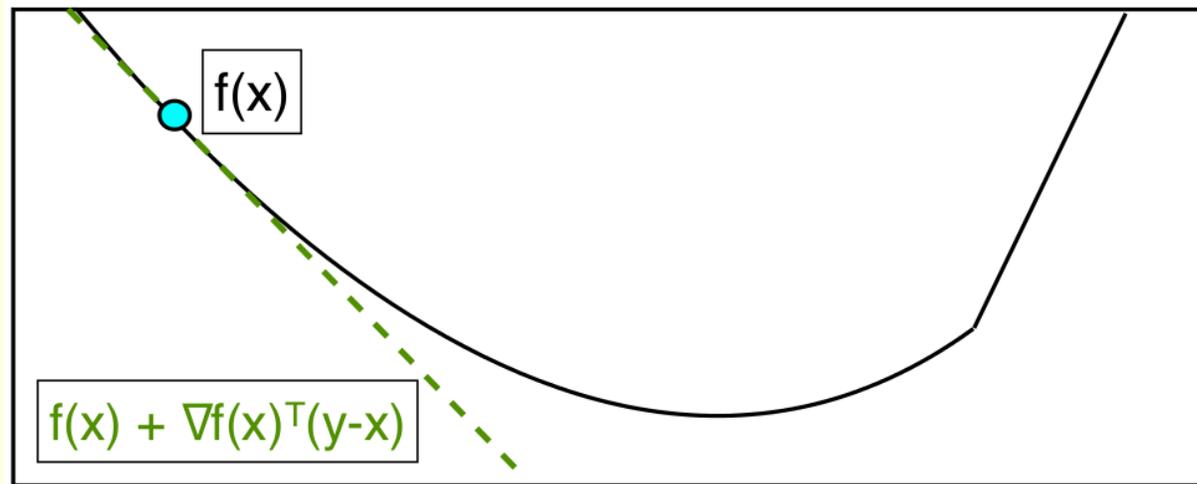
# A quick overview of Convex Analysis

An example of a convex nonsmooth function



# A quick overview of Convex Analysis

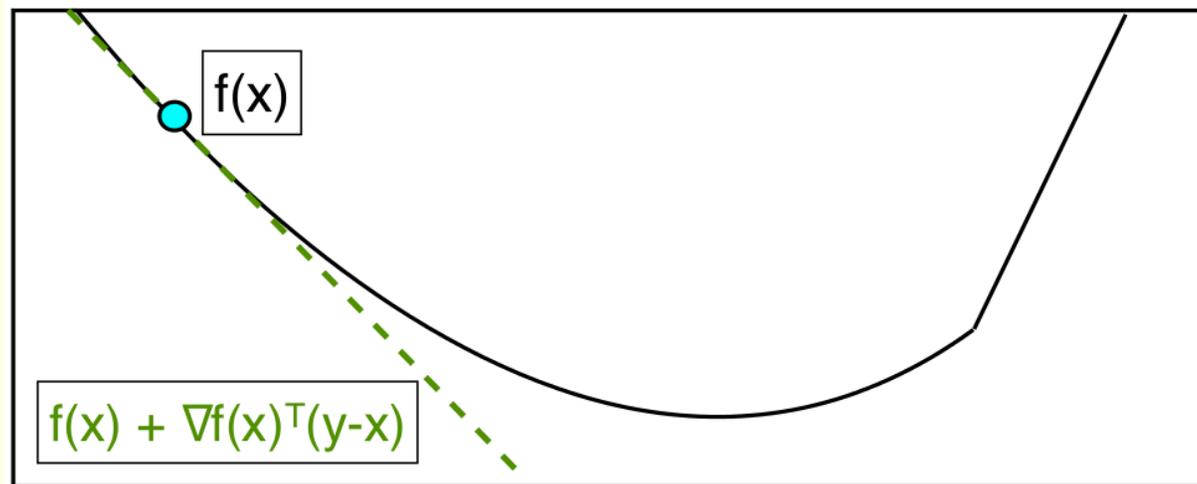
An example of a convex nonsmooth function



$\{\nabla f(x)\} = \{\text{slope of the linearization supporting } f, \text{ tangent at } x\}$

# A quick overview of Convex Analysis

An example of a convex nonsmooth function



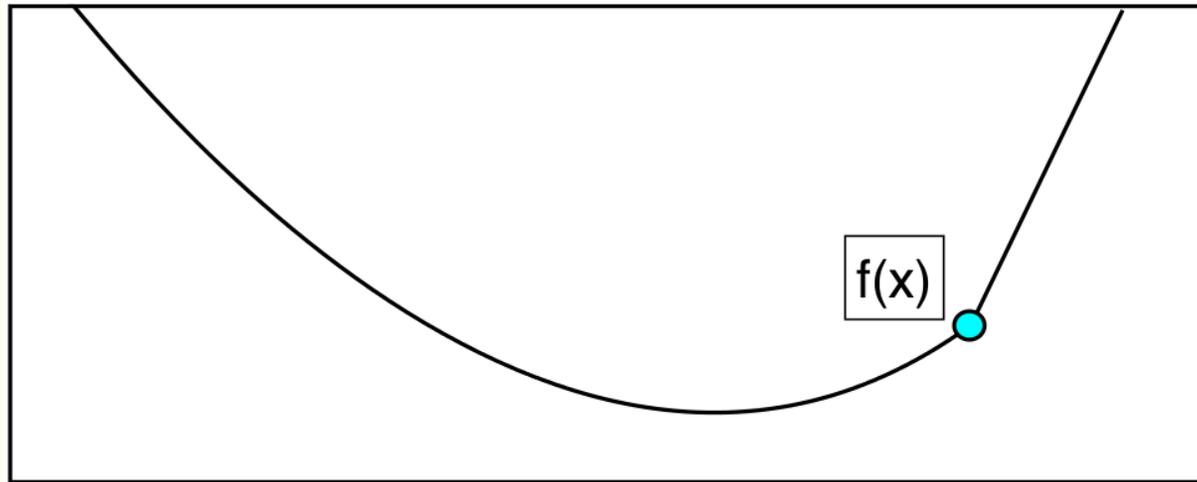
$\{\nabla f(\mathbf{x})\} = \{\text{slope of the linearization supporting } f, \text{ tangent at } \mathbf{x}\}$

By convexity,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \text{ for all } \mathbf{y}$$

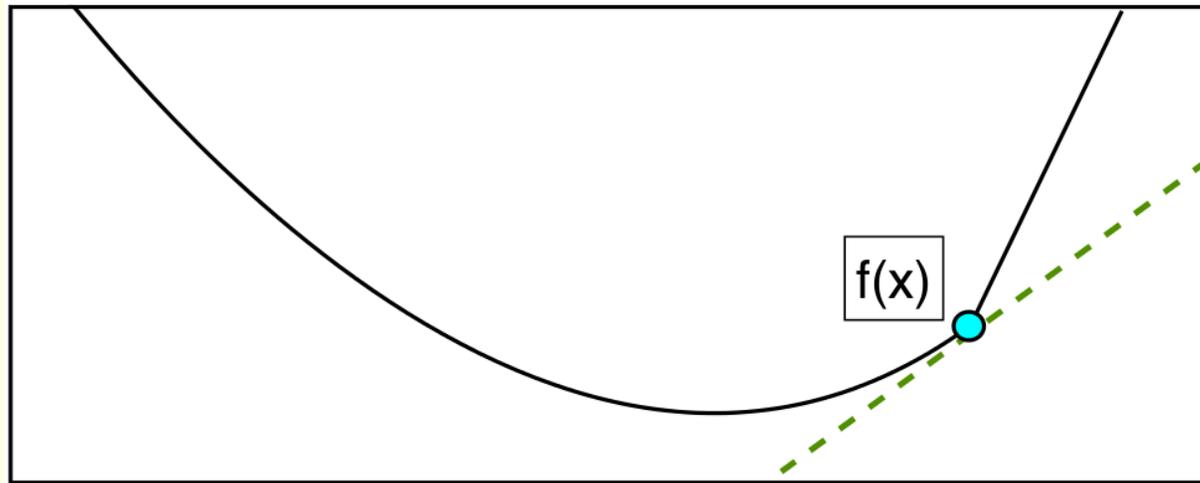
# A quick overview of Convex Analysis

An example of a convex nonsmooth function



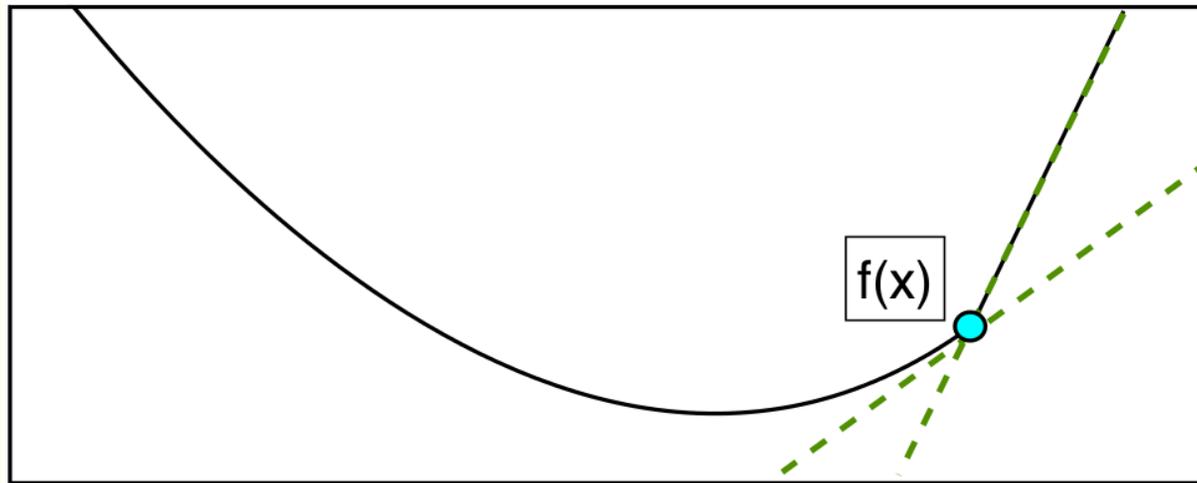
# A quick overview of Convex Analysis

An example of a convex nonsmooth function



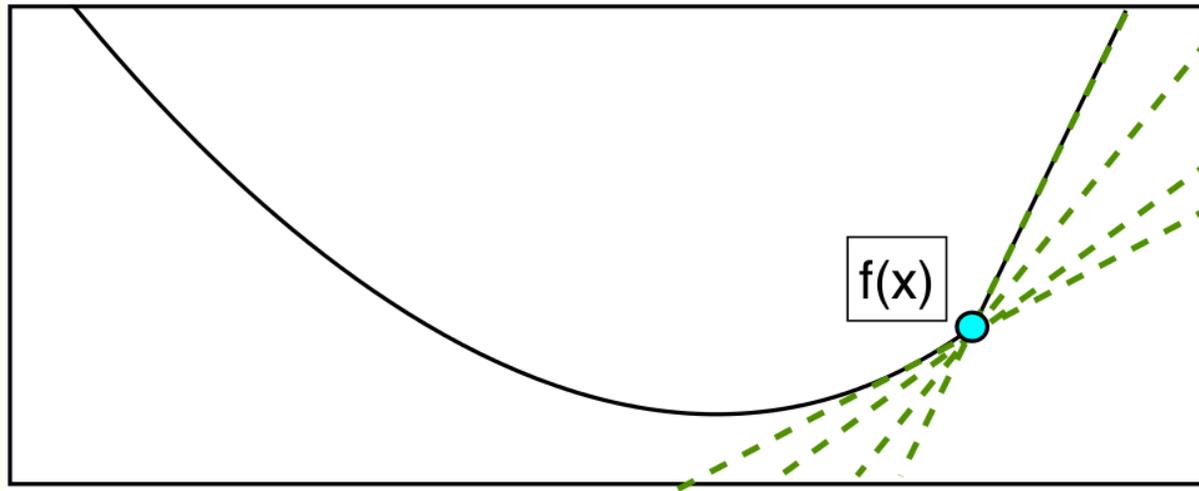
# A quick overview of Convex Analysis

An example of a convex nonsmooth function



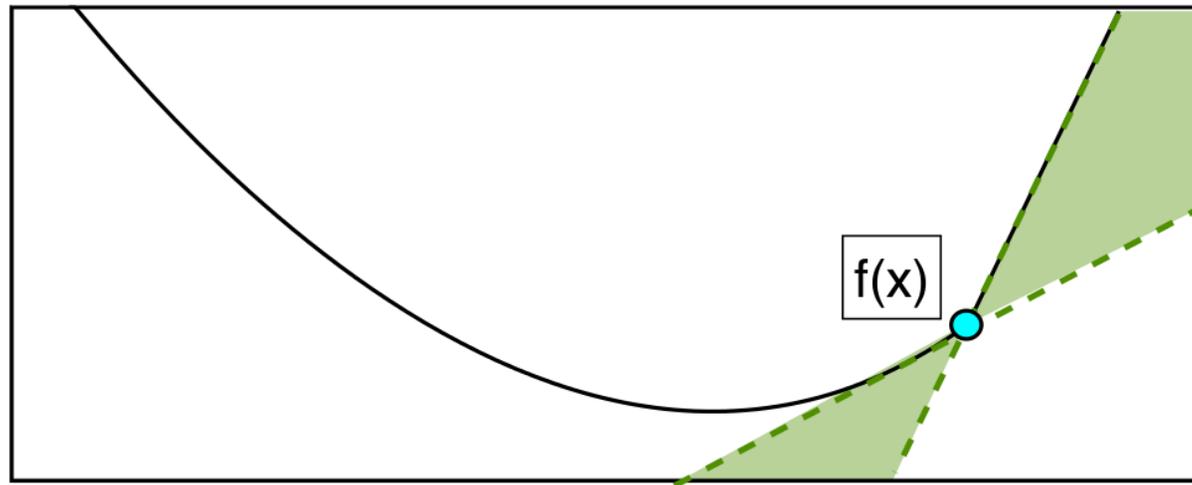
# A quick overview of Convex Analysis

An example of a convex nonsmooth function



# A quick overview of Convex Analysis

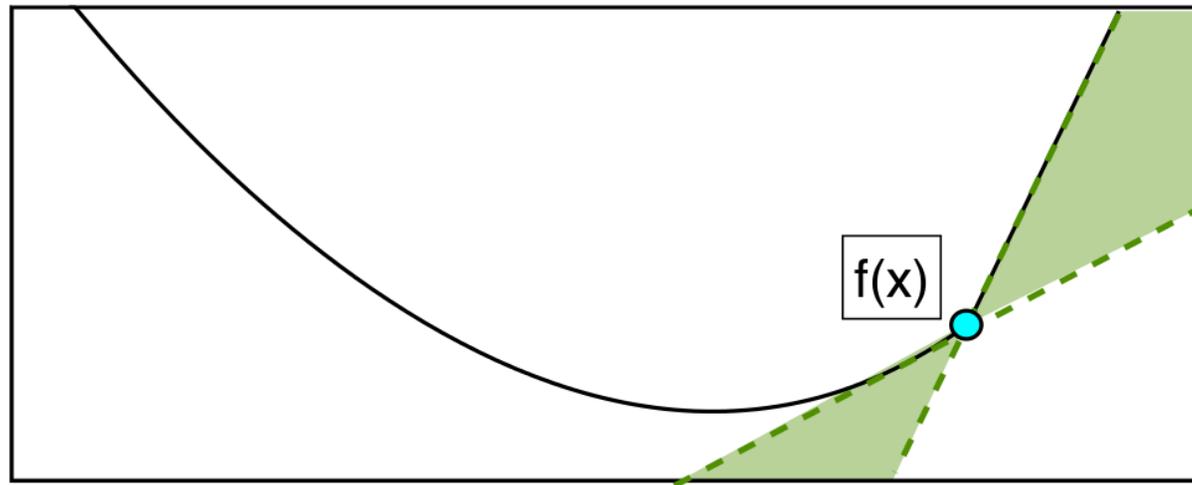
An example of a convex nonsmooth function



$$\partial f(x) = \{g \in \mathbb{R}^n : f(y) \geq f(x) + \langle g, y - x \rangle \text{ for all } y\}$$

# A quick overview of Convex Analysis

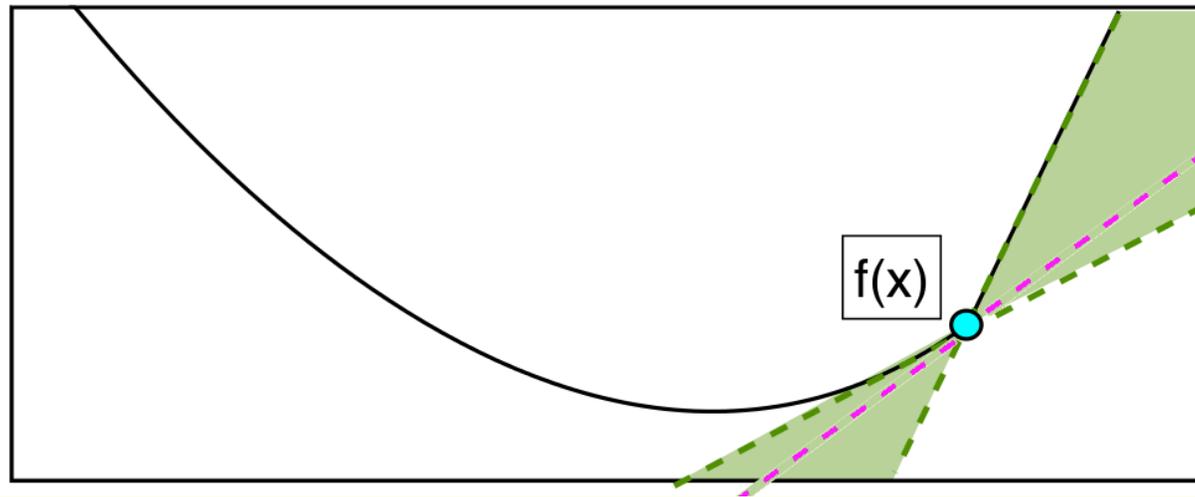
An example of a convex nonsmooth function



$$\begin{aligned}\partial f(x) &= \{g \in \mathbb{R}^n : f(y) \geq f(x) + \langle g, y - x \rangle \text{ for all } y\} \\ &= \{\text{slopes of linearizations supporting } f, \text{ tangent at } x\}\end{aligned}$$

# What can be done with the oracle output?

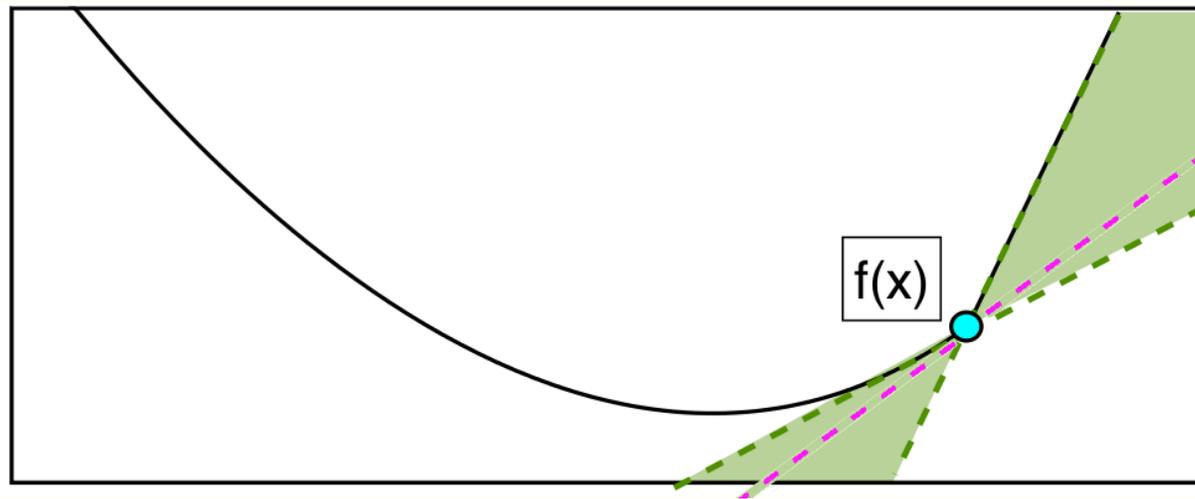
An example of a convex nonsmooth function



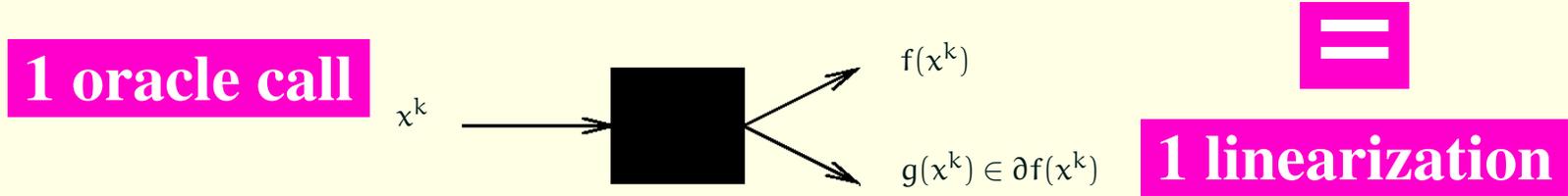
$$\begin{aligned}\partial f(x) &= \{g \in \mathbb{R}^n : f(y) \geq f(x) + \langle g, y - x \rangle \text{ for all } y\} \\ &= \{\text{slopes of linearizations supporting } f, \text{ tangent at } x\}\end{aligned}$$

# What can be done with the oracle output?

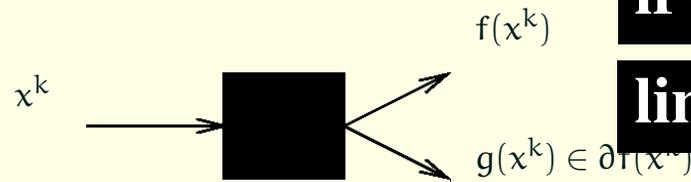
An example of a convex nonsmooth function



$$\begin{aligned} \partial f(x) &= \{g \in \mathbb{R}^n : f(y) \geq f(x) + \langle g, y - x \rangle \text{ for all } y\} \\ &= \{\text{slopes of linearizations supporting } f, \text{ tangent at } x\} \end{aligned}$$

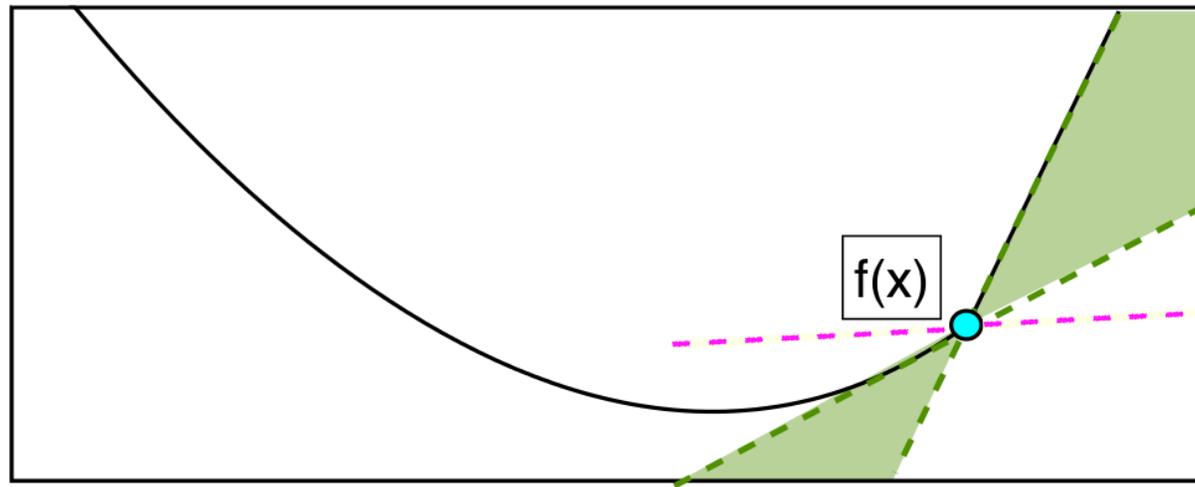


**BEWARE:**



**if oracle output is not accurate,**

**linearization can be wrong!**



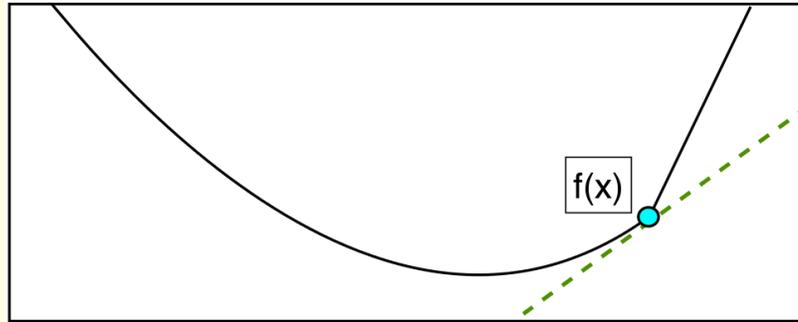
**wrong  $g(x^k)$  gives bad linearization at  $x^k$**

$$\partial f(x) = \{g \in \mathbb{R}^n : f(y) \geq f(x) + \langle g, y - x \rangle \text{ for all } y\}$$

(similarly if wrong  $f(x^k)$ , more on this later)

# How is the oracle information used?

Putting together linearizations



creates a cutting-plane **model M** for  $f$

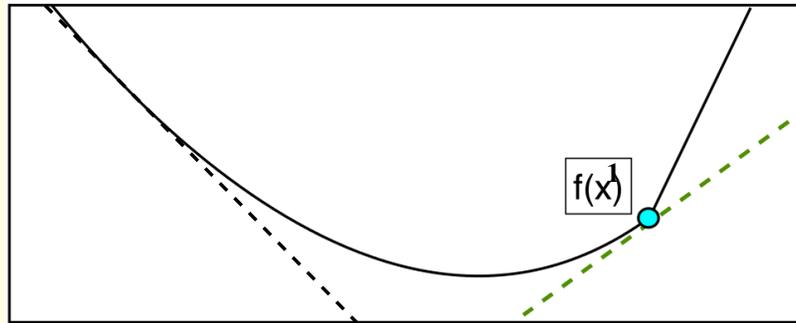
$$\begin{array}{l} f^i = f(x^i) \\ \rightarrow \blacksquare \swarrow \searrow \\ g^i = g(x^i) \end{array}$$

$\Rightarrow$

$$f^i + \langle g^i, x - x^i \rangle$$

# How is the oracle information used?

Putting together linearizations



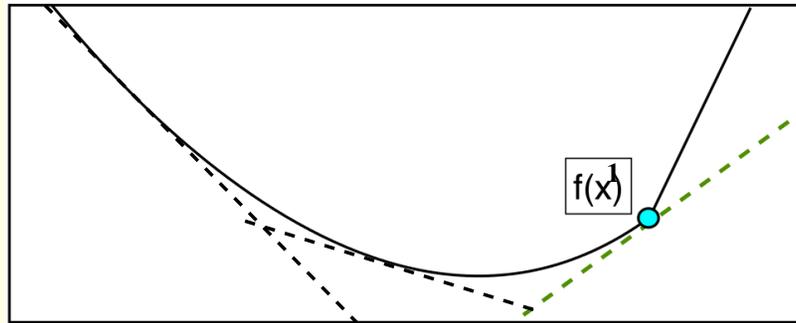
creates a cutting-plane **model M** for f

$$\begin{array}{l} f^i = f(x^i) \\ \rightarrow \blacksquare \swarrow \searrow \\ g^i = g(x^i) \end{array}$$

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

# How is the oracle information used?

Putting together linearizations



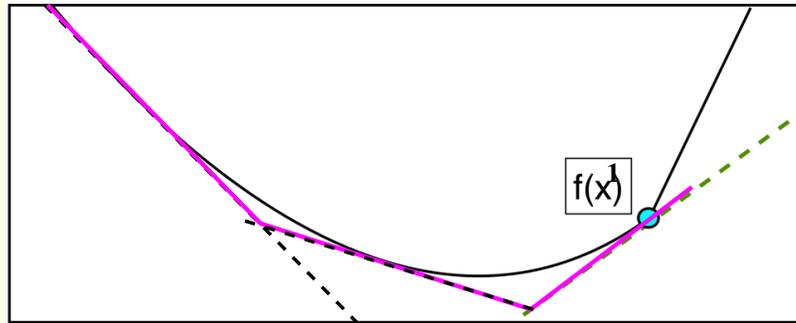
creates a cutting-plane **model M** for f

$$\begin{array}{l} f^i = f(x^i) \\ \rightarrow \blacksquare \swarrow \searrow \\ g^i = g(x^i) \end{array}$$

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

# How is the oracle information used?

Putting together linearizations



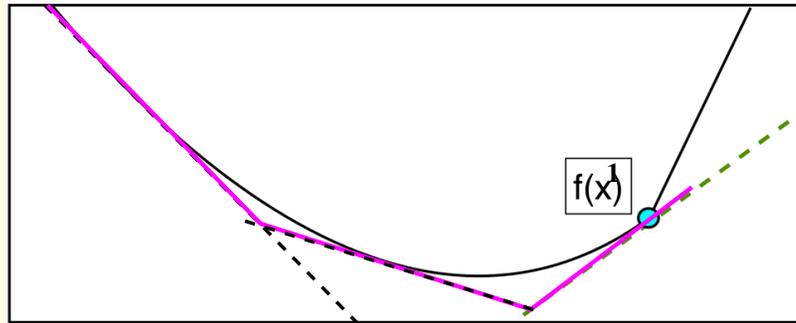
creates a cutting-plane **model M** for f

$$\begin{array}{l} f^i = f(x^i) \\ \rightarrow \blacksquare \swarrow \searrow \\ g^i = g(x^i) \end{array} \quad \Longrightarrow \quad \mathbf{M}(y) = \max_i \left\{ f^i + \langle g^i, y - x^i \rangle \right\}$$

(just one type of model, many others are possible)

# How is the oracle information used?

## Putting together linearizations



creates a cutting-plane **model M** for f

$$\begin{array}{l} f^i = f(x^i) \\ \rightarrow \blacksquare \swarrow \searrow \\ g^i = g(x^i) \end{array} \quad \Longrightarrow \quad \mathbf{M}_k(y) = \max_{i \leq k} \left\{ f^i + \langle g^i, y - x^i \rangle \right\}$$

(just one type of model, many others are possible)

## Infinite bundling yields $\text{prox}_t^f$

$$\text{WANT: } p = \text{prox}_t^f(x) \iff p = \arg \min f(y) + \frac{1}{2t} \|y - x\|_2^2$$

$$\text{HAVE: } q^k = \text{prox}_{t_k}^{M_k}(x) \iff q^k = \arg \min M_k(y) + \frac{1}{2t_k} \|y - x^k\|_2^2$$

$$\iff 0 = G^k + \frac{1}{t_k} (q^k - x^k)$$

**Theorem [CL93]** Suppose the models satisfy **for**  $G^k \in \partial M_k(q^k)$

- $M_k(y) \leq f(y)$  for all  $k$  and  $y$
- $M_{k+1}(y) \geq f(q^k) + \langle g(q^k), y - x^k \rangle$
- $M_{k+1}(y) \geq M_k(q^k) + \langle G^k, y - x^k \rangle$

If  $0 < t_{\min} \leq t_{k+1} \leq t_k$ , then

$$\lim_{k \rightarrow \infty} q^k = p \quad \text{and} \quad \lim_{k \rightarrow \infty} M_k(q^k) = f(p)$$

## Infinite bundling yields $\text{prox}_t^f$

$$\text{WANT: } p = \text{prox}_t^f(x) \iff p = \arg \min f(y) + \frac{1}{2t} \|y - x\|_2^2$$

$$\text{HAVE: } q^k = \text{prox}_{t_k}^{M_k}(x) \iff q^k = \arg \min M_k(y) + \frac{1}{2t_k} \|y - x^k\|_2^2$$

$$\iff 0 = G^k + \frac{1}{t_k} (q^k - x^k)$$

for  $G^k \in \partial M_k(q^k)$

**Theorem** [CL93] Suppose the models satisfy

- $M_k(y) \leq f(y)$  for all  $k$  and  $y$
- $M_{k+1}(y) \geq f(q^k) + \langle g(q^k), y - x^k \rangle$
- $M_{k+1}(y) \geq M_k(q^k) + \langle G^k, y - x^k \rangle$

If  $0 < t_{\min} \leq t_{k+1} \leq t_k$ , then

$$\lim_{k \rightarrow \infty} q^k = p \quad \text{and} \quad \lim_{k \rightarrow \infty} M_k(q^k) = f(p)$$

# Models for the half-and-half function

STRUCTURE	$f(x)$
<b>none</b>	$\sqrt{x^\top Ax} + x^\top Bx$
<b>sum</b>	$f_1(x) + f_2(x)$ $f_1(x) = \sqrt{x^\top Ax}$ $f_2(x) = x^\top Bx$ $f_2$ is smooth
<b>compo</b> <b>sition</b>	$(h \circ c)(x)$ $c(x) = (x, x^\top Bx) \in \mathbb{R}^{n+1}$ $h(C) = \sqrt{C_{1:n}^\top AC_{1:n}} + C_{n+1}$ $c$ $h$

# Models for the half-and-half function

STRUCTURE	$f(x)$
<b>none</b>	$\sqrt{x^\top Ax} + x^\top Bx$
<b>sum</b>	$f_1(x) + f_2(x)$ $f_1(x) = \sqrt{x^\top Ax}$ $f_2(x) = x^\top Bx$ <b><math>f_2</math> is smooth</b>
<b>composition</b>	$(h \circ c)(x)$ $c(x) = (x, x^\top Bx) \in \mathbb{R}^{n+1}$ <b><math>c</math> is smooth</b> $h(C) = \sqrt{C_{1:n}^\top AC_{1:n}} + C_{n+1}$ <b><math>h</math> is sublinear</b>

# Models for the half-and-half function

STRUCTURE	$f(x)$
<b>none</b>	$\sqrt{x^\top Ax} + x^\top Bx$ $f^k := f(x^k), g^k \in \partial f(x^k)$
<b>sum</b>	$f_1(x) + f_2(x)$ $f_1(x) = \sqrt{x^\top Ax}$ $f_2(x) = x^\top Bx$
<b>compo sition</b>	$(h \circ c)(x)$ $c(x) = (x, x^\top Bx) \in \mathbb{R}^{n+1}$ $h(C) = \sqrt{C_{1:n}^\top AC_{1:n}} + C_{n+1}$

# Models for the half-and-half function

STRUCTURE	$f(x)$
<b>none</b>	$\sqrt{x^\top Ax} + x^\top Bx$ $f^k := f(x^k), g^k \in \partial f(x^k)$
<b>sum</b>	$f_1(x) + f_2(x)$ $f_1(x) = \sqrt{x^\top Ax}$ $f_2(x) = x^\top Bx$ $f_1^k, g_1^k, f_2^k, \nabla f_2(x^k)$
<b>compo sition</b>	$(h \circ c)(x)$ $c(x) = (x, x^\top Bx) \in \mathbb{R}^{n+1}$ $h(C) = \sqrt{C_{1:n}^\top AC_{1:n}} + C_{n+1}$

# Models for the half-and-half function

STRUCTURE	$f(x)$
<b>none</b>	$\sqrt{x^\top Ax} + x^\top Bx$ $f^k := f(x^k), g^k \in \partial f(x^k)$
<b>sum</b>	$f_1(x) + f_2(x)$ $f_1(x) = \sqrt{x^\top Ax}$ $f_2(x) = x^\top Bx$ $f_1^k, g_1^k, f_2^k, \nabla f_2(x^k)$
<b>compo sition</b>	$(h \circ c)(x)$ $c(x) = (x, x^\top Bx) \in \mathbb{R}^{n+1}$ $c^k = c(x^k), c'(x^k)$ $h(C) = \sqrt{C_{1:n}^\top AC_{1:n}} + C_{n+1}$ $h^k, g^k \in \partial h(c^k)$

**NSO pitfalls**

**NSO pitfalls**

**NSO pitfalls**

**NSO pitfalls**

## Stopping test in smooth optimization

Algorithms for unconstrained smooth optimization use as optimality certificate Fermat's rule

$$0 = \nabla f(\bar{x})$$

and generate a minimizing sequence:

$$\{x^k\} \rightarrow \bar{x} \text{ such that } \nabla f(x^k) \rightarrow 0.$$

If  $f \in C^1$ , then  $\nabla f(\bar{x}) = 0$

# Stopping test in smooth optimization

Algorithms for unconstrained smooth optimization use as optimality certificate Fermat's rule

$$0 = \nabla f(\bar{x})$$

and generate a minimizing sequence:

$$\{x^k\} \rightarrow \bar{x} \text{ such that } \nabla f(x^k) \rightarrow 0.$$

If  $f \in C^1$ , then  $\nabla f(\bar{x}) = 0$  things are less straightforward if  $f$  is nonsmooth...

## What happens with the stopping test in NSO?

Algorithms for unconstrained NSO use as optimality certificate the inclusion

$$0 \in \partial f(\bar{x})$$

- As a set-valued mapping  $\partial f(x)$  is osc:

$$\left( x^k, g(x^k) \in \partial f(x^k) \right) : \begin{cases} x^k \rightarrow \bar{x} \\ g(x^k) \rightarrow \bar{g} \end{cases} \implies \bar{g} \in \partial f(\bar{x})$$

# What happens with the stopping test in NSO?

Algorithms for unconstrained NSO use as optimality certificate the inclusion

$$0 \in \partial f(\bar{x})$$

- As a set-valued mapping  $\partial f(x)$  is osc:

$$\left( x^k, g(x^k) \in \partial f(x^k) \right) : \begin{cases} x^k \rightarrow \bar{x} \\ g(x^k) \rightarrow \bar{g} \end{cases} \implies \bar{g} \in \partial f(\bar{x})$$

- As a set-valued mapping,  $\partial f(x)$  is **not** isc:  
Given  $\bar{g} \in \partial f(\bar{x})$

$$/ \quad \exists \left( x^k, g(x^k) \in \partial f(x^k) \right) : \begin{cases} x^k \rightarrow \bar{x} \\ g(x^k) \rightarrow \bar{g} \end{cases}$$

# What happens with the stopping test in NSO?

Algorithms for unconstrained NSO use as optimality certificate the inclusion

$$0 \in \partial f(\bar{x})$$

- As a set-valued mapping  $\partial f(x)$  is osc:

$$\left( x^k, g(x^k) \in \partial f(x^k) \right) : \begin{cases} x^k \rightarrow \bar{x} \\ g(x^k) \rightarrow \bar{g} \end{cases} \implies \bar{g} \in \partial f(\bar{x})$$

- As a set-valued mapping,  $\partial f(x)$  is **not** isc:  
Given  $\bar{g} \in \partial f(\bar{x})$

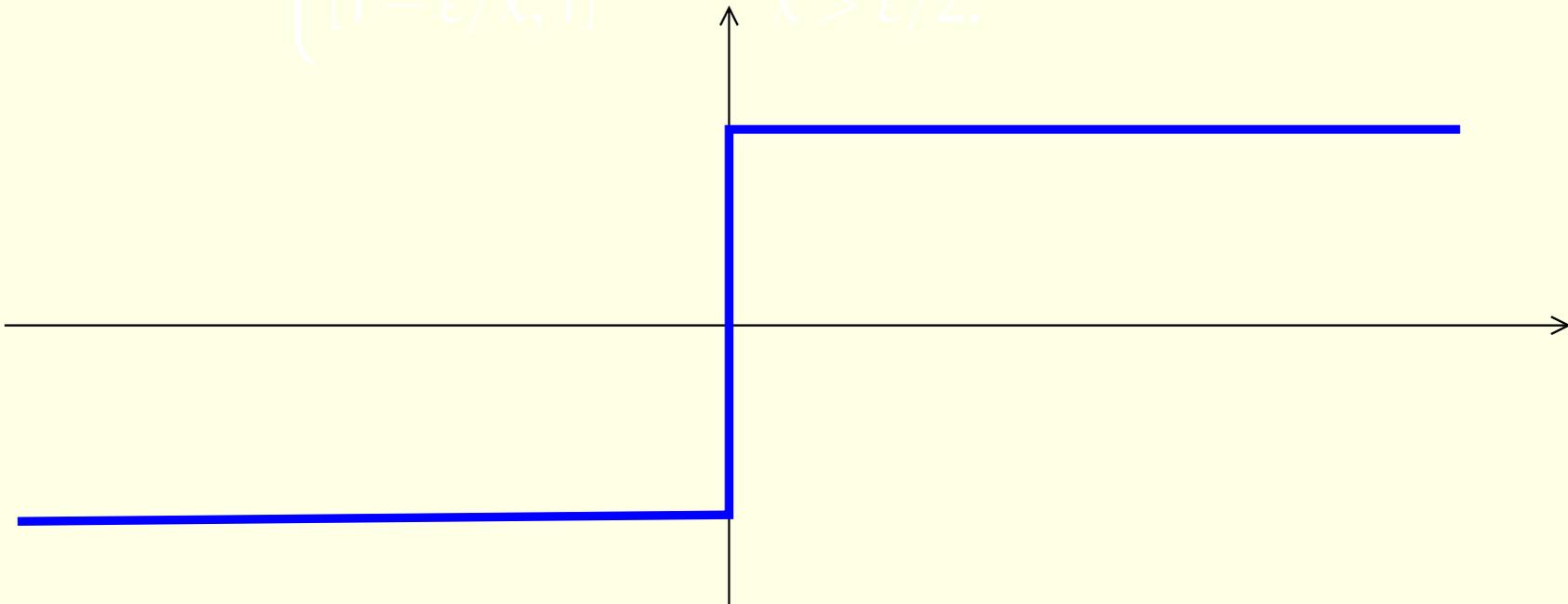
$$\not\exists \left( x^k, g(x^k) \in \partial f(x^k) \right) : \begin{cases} x^k \rightarrow \bar{x} \\ g(x^k) \rightarrow \bar{g} \end{cases}$$

# The subdifferential

For the absolute value function,  $f(x) = |x|$

$$\partial f(x) = \begin{cases} -1 & x < 0 \\ [-1, 1] & x = 0 \\ 1 & x > 0 \end{cases}$$

$$\partial_\varepsilon f(x) = \begin{cases} [-1, -1 - \varepsilon/x] & x < -\varepsilon/2, \\ [-1, 1] & -\varepsilon/2 \leq x \leq \varepsilon/2, \\ [1 - \varepsilon/x, 1] & x > \varepsilon/2. \end{cases}$$



## **What happens with the stopping test in NSO?**

We need to design a sound stopping test that does not rely on the straightforward extension of Fermat's rule.

## What happens with the stopping test in NSO?

We need to design a sound stopping test that does not rely on the straightforward extension of Fermat's rule. We use instead

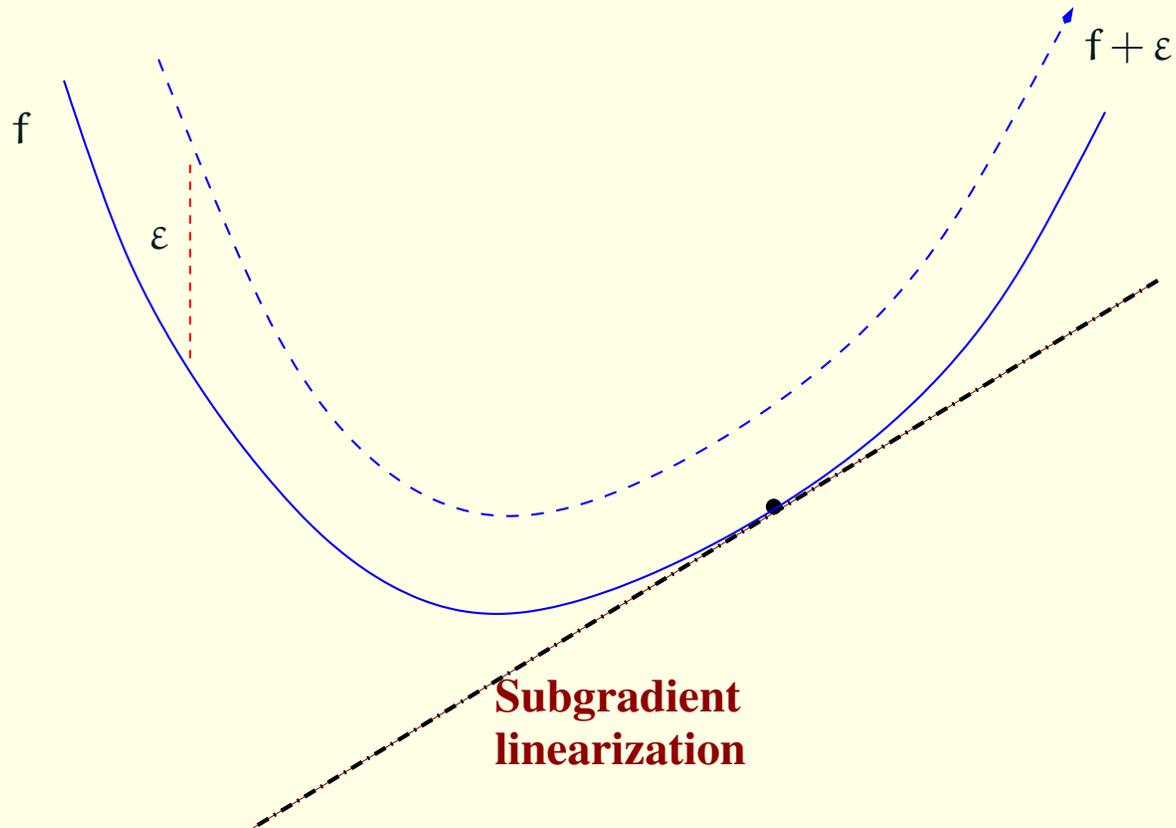
$$\bar{g} \in \partial_{\bar{\varepsilon}} f(\bar{x}) \quad \text{for } \|\bar{g}\| \text{ and } \bar{\varepsilon} \text{ small}$$

where the  $\varepsilon$ -subdifferential contains the slopes of linearizations supporting  $f$  **up to  $\varepsilon$** , tangent at  $x$ :

$$\partial_{\varepsilon} f(x) = \{g \in \mathbb{R}^n : f(y) \geq f(x) + \langle g, y - x \rangle - \varepsilon \text{ for all } y\}$$

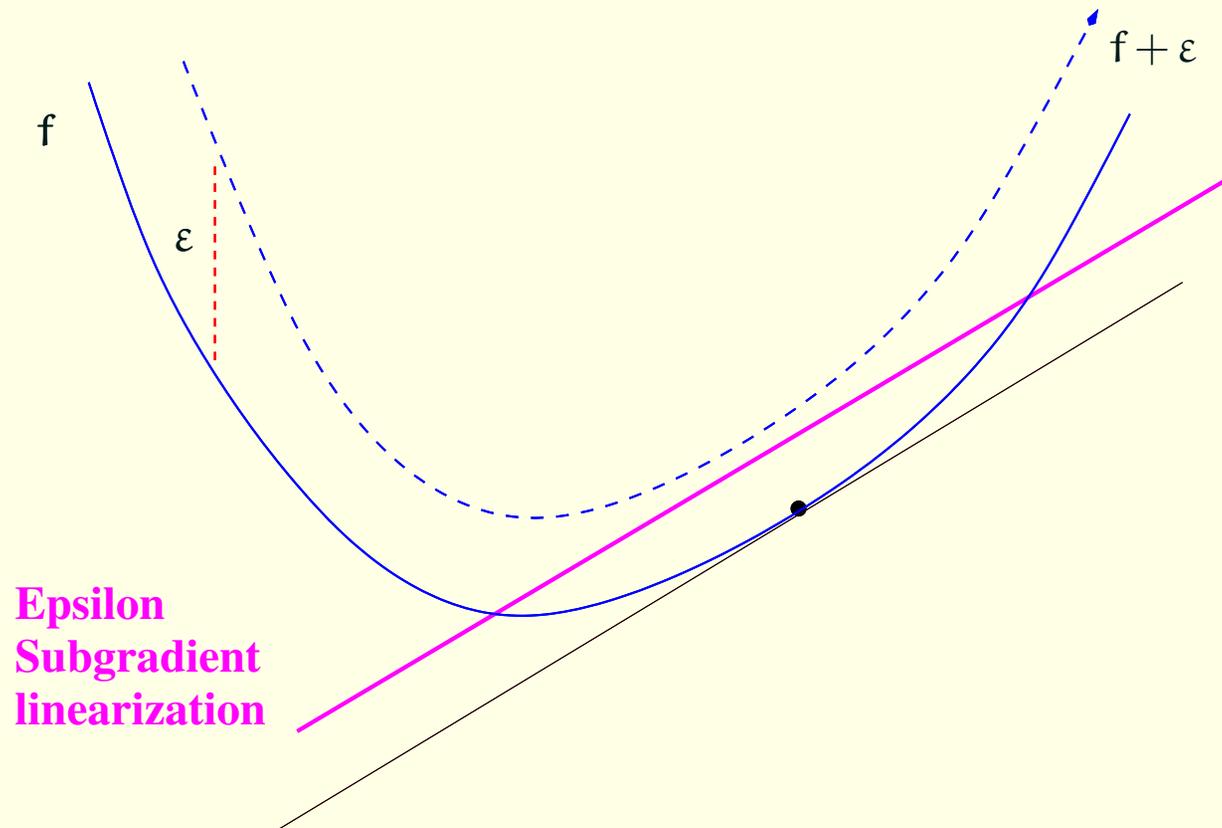
# The $\varepsilon$ -subdifferential

$$\partial_\varepsilon f(x) = \{g \in \mathbb{R}^n : f(y) \geq f(x) + \langle g, y - x \rangle - \varepsilon \text{ for all } y\}$$



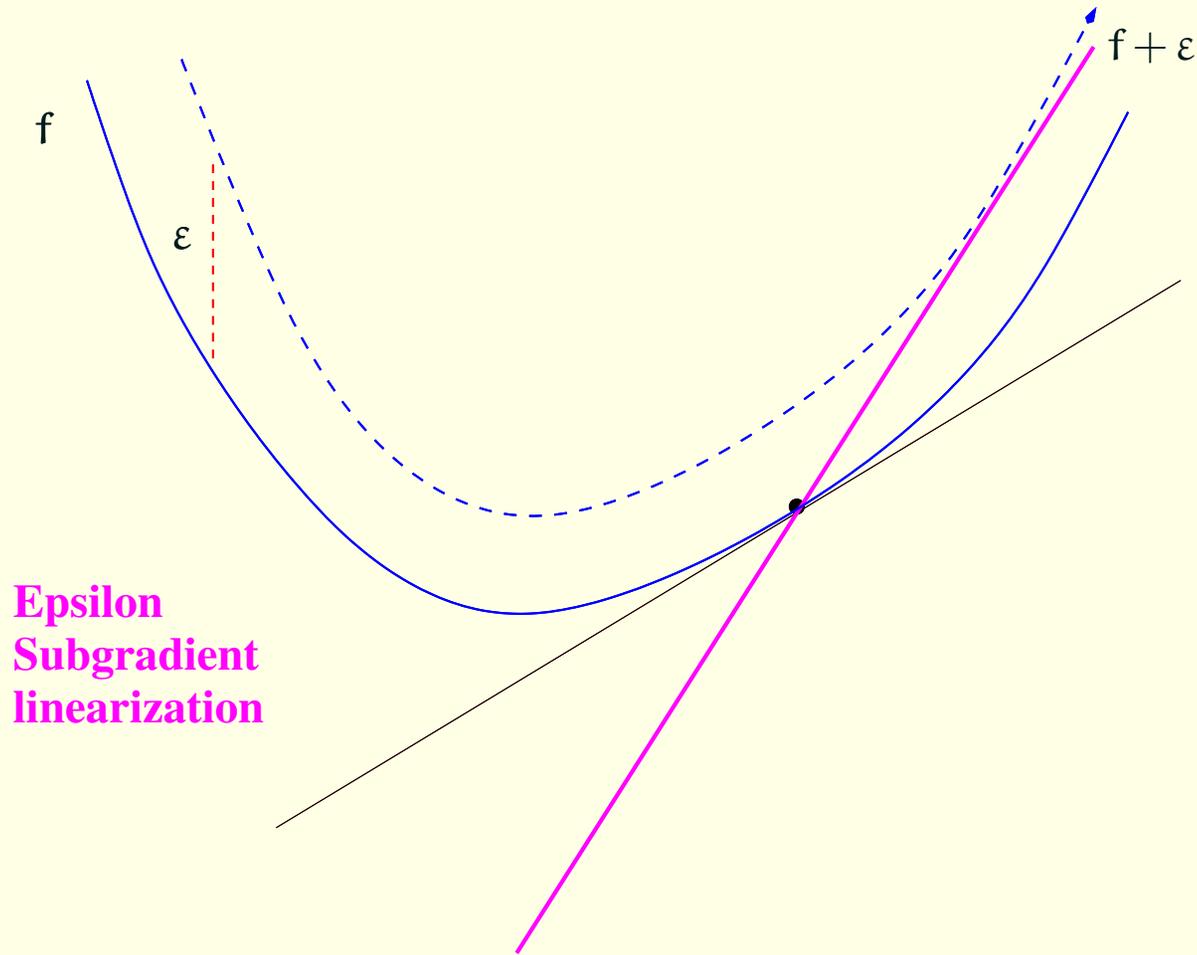
# The $\varepsilon$ -subdifferential

$$\partial_\varepsilon f(x) = \{g \in \mathbb{R}^n : f(y) \geq f(x) + \langle g, y - x \rangle - \varepsilon \text{ for all } y\}$$



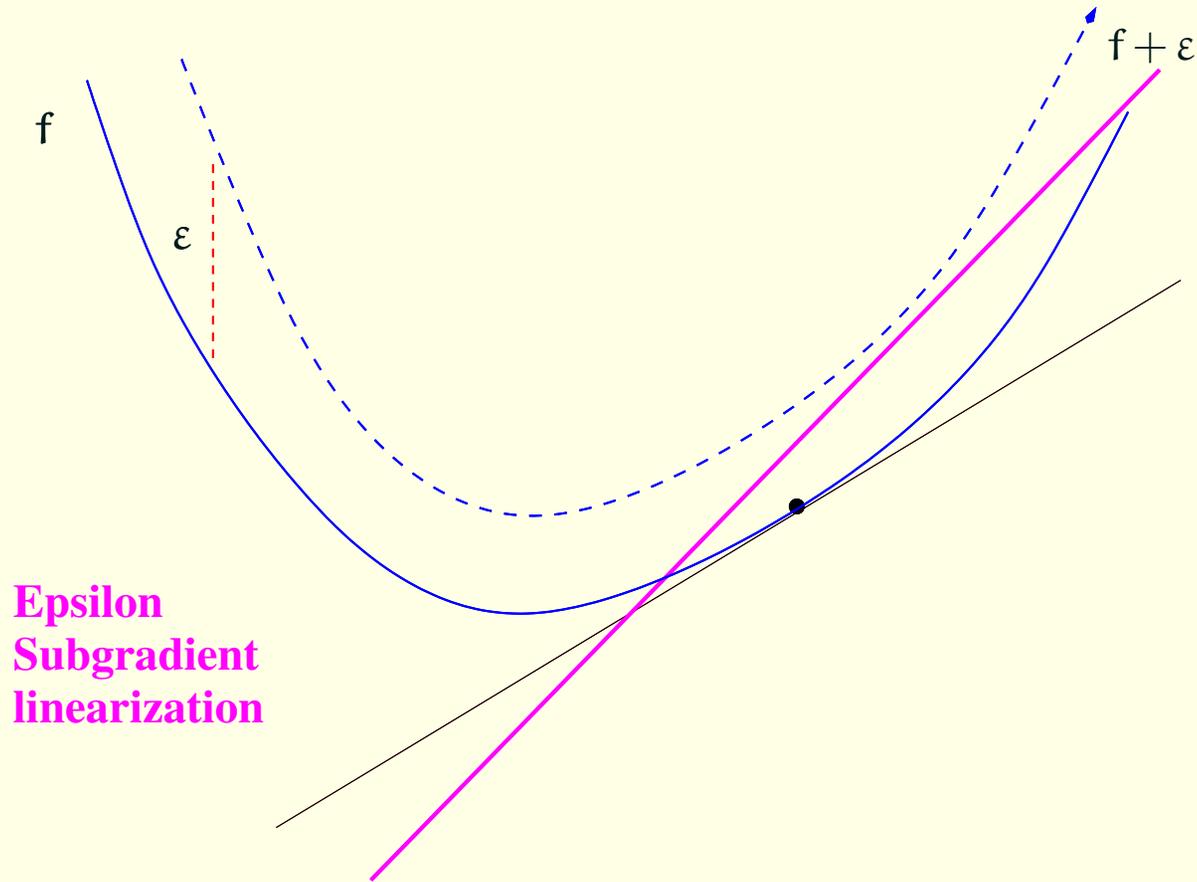
# The $\varepsilon$ -subdifferential

$$\partial_\varepsilon f(x) = \{g \in \mathbb{R}^n : f(y) \geq f(x) + \langle g, y - x \rangle - \varepsilon \text{ for all } y\}$$



# The $\varepsilon$ -subdifferential

$$\partial_\varepsilon f(x) = \{g \in \mathbb{R}^n : f(y) \geq f(x) + \langle g, y - x \rangle - \varepsilon \text{ for all } y\}$$

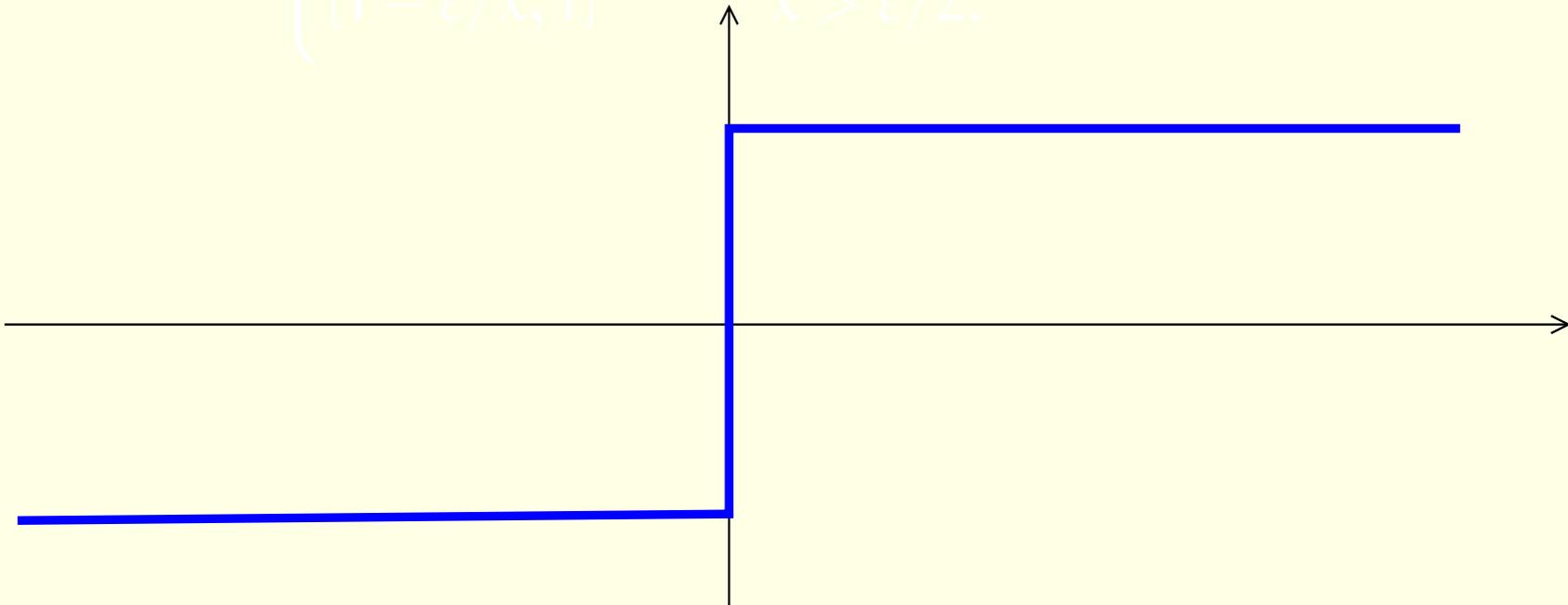


# The $\varepsilon$ -subdifferential

For the absolute value function,  $f(x) = |x|$

$$\partial f(x) = \begin{cases} -1 & x < 0 \\ [-1, 1] & x = 0 \\ 1 & x > 0 \end{cases}$$

$$\partial_\varepsilon f(x) = \begin{cases} [-1, -1 - \varepsilon/x] & x < -\varepsilon/2, \\ [-1, 1] & -\varepsilon/2 \leq x \leq \varepsilon/2, \\ [1 - \varepsilon/x, 1] & x > \varepsilon/2. \end{cases}$$

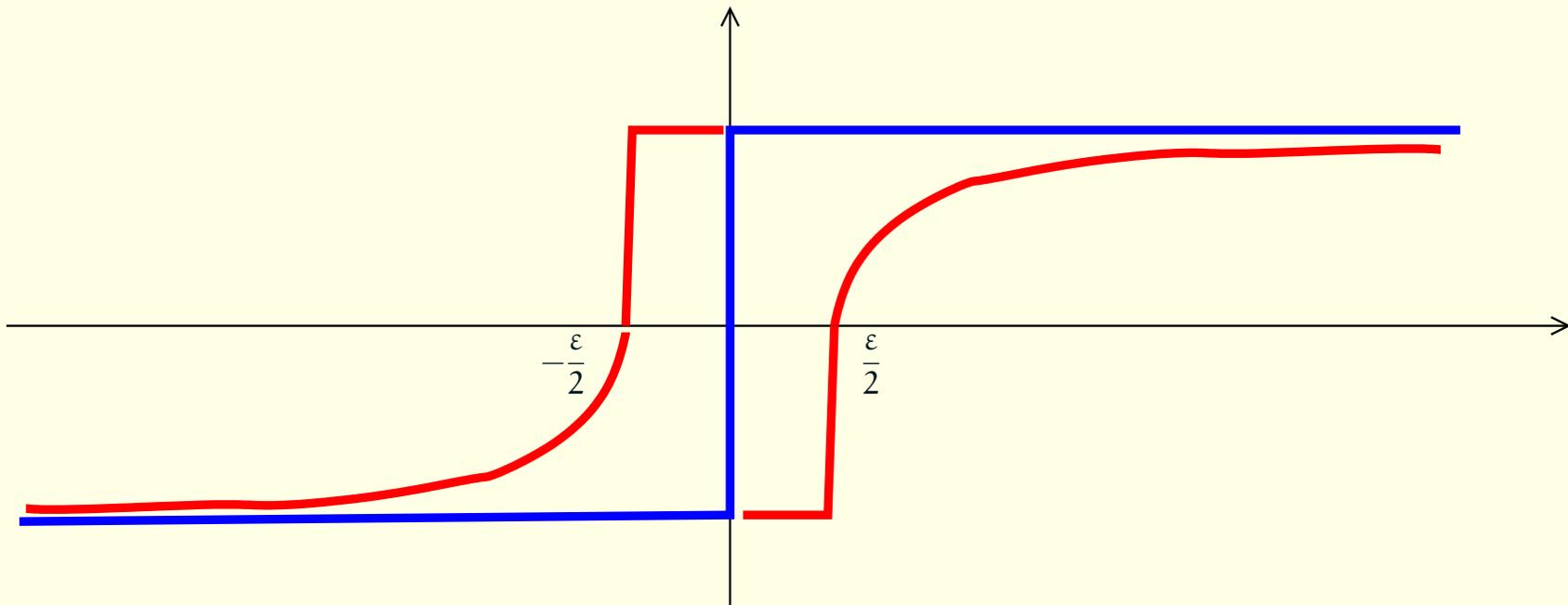


# The $\varepsilon$ -subdifferential

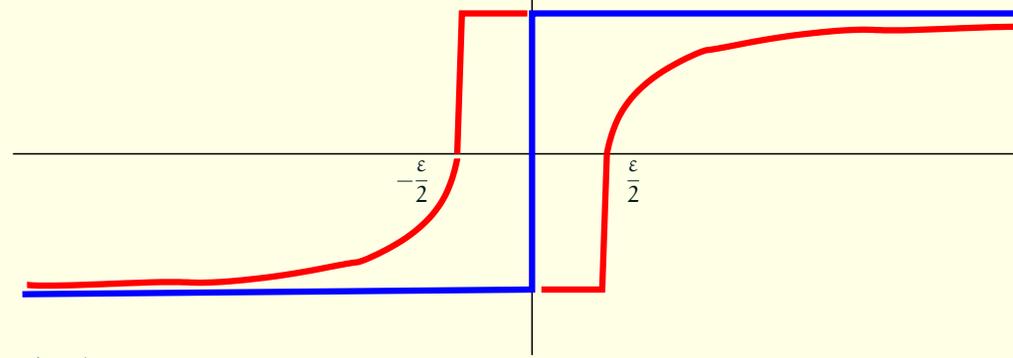
For the absolute value function,  $f(x) = |x|$

$$\partial f(x) = \begin{cases} -1 & x < 0 \\ [-1, 1] & x = 0 \\ 1 & x > 0 \end{cases}$$

$$\partial_\varepsilon f(x) = \begin{cases} [-1, -1 - \varepsilon/x] & \text{if } x < -\varepsilon/2, \\ [-1, 1] & \text{if } -\varepsilon/2 \leq x \leq \varepsilon/2, \\ [1 - \varepsilon/x, 1] & \text{if } x > \varepsilon/2. \end{cases}$$



# The $\varepsilon$ -subdifferential



- As a set-valued mapping  $\partial_\varepsilon f(x)$  is osc:

$$\left( \varepsilon^k, x^k, G(x^k) \in \partial_{\varepsilon^k} f(x^k) \right) : \begin{cases} \varepsilon^k \rightarrow \varepsilon \\ x^k \rightarrow \bar{x} \\ G(x^k) \rightarrow \bar{g} \end{cases} \implies \bar{g} \in \partial_{\bar{\varepsilon}} f(\bar{x})$$

- As a set-valued mapping,  $\partial_\varepsilon f(x)$  is isc:  
Given  $\bar{g} \in \partial_{\bar{\varepsilon}} f(\bar{x})$

$$\exists \left( \varepsilon^k, x^k, G(x^k) \in \partial_{\varepsilon^k} f(x^k) \right) : \begin{cases} \varepsilon^k \rightarrow \bar{\varepsilon} \\ x^k \rightarrow \bar{x} \\ G(x^k) \rightarrow \bar{g} \end{cases}$$

# The $\varepsilon$ -subdifferential and bundle methods

Generate iterates so that for a **subsequence**  $\{\hat{x}^k\}$

- As a set-valued mapping  $\partial_\varepsilon f(x)$  is osc:

$$\left( \varepsilon^k, \hat{x}^k, G(\hat{x}^k) \in \partial_{\varepsilon^k} f(\hat{x}^k) \right) : \begin{cases} \varepsilon^k \rightarrow \bar{\varepsilon} \\ x^k \rightarrow \bar{x} \\ G(\hat{x}^k) \rightarrow \bar{g} \end{cases} \implies \bar{g} \in \partial_{\bar{\varepsilon}} f(\bar{x})$$

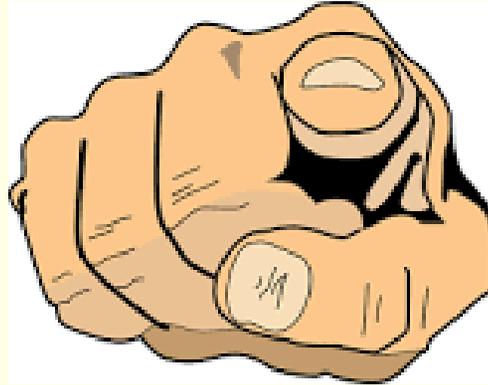
**with  $\bar{\varepsilon} = 0$  and  $\bar{g} = 0$**

- As a set-valued mapping,  $\partial_\varepsilon f(x)$  is isc: Given  $\bar{g} \in \partial_{\bar{\varepsilon}} f(\bar{x})$ :

$$\exists \left( \varepsilon^k, \hat{x}^k, G(\hat{x}^k) \in \partial_{\varepsilon^k} f(\hat{x}^k) \right) : \begin{cases} \varepsilon^k \rightarrow \bar{\varepsilon} \\ x^k \rightarrow \bar{x} \\ G(x^k) \rightarrow \bar{g} \end{cases}$$

# The $\epsilon$ -subdifferential and bundle methods

You told us



we were going to use subgradient information provided by an oracle or a black box, and now you want to use  $\epsilon$ -subgradients!



## The transportation formula

How to express subgradients at  $x^i$  as  $\varepsilon$ -subgradients at  $\hat{x}^k$ ?

$g^i \in \partial f(x^i)$  if and only if, for all  $y \in \mathbb{R}^n$

$$\begin{aligned} f(y) &\geq f(x^i) + \langle g^i, y - x^i \rangle \\ &= f(x^i) + \langle g^i, y - x \rangle + \langle g^i, x - x^i \rangle \\ &= f(\hat{x}^k) + \langle g^i, y - x \rangle - \text{Bigl}(f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - x \pm \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - e^i(\hat{x}^k) \end{aligned}$$

## The transportation formula

How to express subgradients at  $x^i$  as  $\varepsilon$ -subgradients at  $\hat{x}^k$ ?

$$g^i \in \partial f(x^i) \quad \text{if and only if, for all } y \in \mathbb{R}^n$$

$$\begin{aligned} f(y) &\geq f(x^i) + \langle g^i, y - x^i \rangle \\ &= f(x^i) + \langle g^i, y - x \rangle \pm f(\hat{x}^k) \\ &= f(\hat{x}^k) + g^{i\top} (y - x) - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + g^{i\top} (y - x \pm \hat{x}^k) - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + g^{i\top} (y - \hat{x}^k) - (f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle) \\ &= f(\hat{x}^k) + g^{i\top} (y - \hat{x}^k) - e^i(\hat{x}^k) \end{aligned}$$

$$\implies g^i \in \partial_{e^i(\hat{x}^k)} f(\hat{x}^k)$$

$$e^i(\hat{x}^k) := f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle \geq 0$$

# The transportation formula

How to express subgradients at  $x^i$  as  $\varepsilon$ -subgradients at  $\hat{x}^k$ ?

$$g^i \in \partial f(x^i) \quad \text{if and only if, for all } y \in \mathbb{R}^n$$

$$\begin{aligned} f(y) &\geq f(x^i) + \langle g^i, y - x^i \rangle \\ &= f(x^i) + \langle g^i, y - x \rangle \pm f(\hat{x}^k) \\ &= f(\hat{x}^k) + \langle g^i, y - x \rangle - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - x \pm \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - e^i(\hat{x}^k) \end{aligned}$$

$$\implies g^i \in \partial_{e^i(\hat{x}^k)} f(\hat{x}^k)$$

$$e^i(\hat{x}^k) := f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle \geq 0$$

## The transportation formula

How to express subgradients at  $x^i$  as  $\varepsilon$ -subgradients at  $\hat{x}^k$ ?

$$g^i \in \partial f(x^i) \quad \text{if and only if, for all } y \in \mathbb{R}^n$$

$$\begin{aligned} f(y) &\geq f(x^i) + \langle g^i, y - x^i \rangle \\ &= f(x^i) + \langle g^i, y - x \rangle \pm f(\hat{x}^k) \\ &= f(\hat{x}^k) + \langle g^i, y - x \rangle - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - x \pm \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - e^i(\hat{x}^k) \\ &\implies g^i \in \partial_{e^i(\hat{x}^k)} f(\hat{x}^k) \\ e^i(\hat{x}^k) &:= f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle \geq 0 \end{aligned}$$

## The transportation formula

How to express subgradients at  $x^i$  as  $\varepsilon$ -subgradients at  $\hat{x}^k$ ?

$$g^i \in \partial f(x^i) \quad \text{if and only if, for all } y \in \mathbb{R}^n$$

$$\begin{aligned} f(y) &\geq f(x^i) + \langle g^i, y - x^i \rangle \\ &= f(x^i) + \langle g^i, y - x \rangle \pm f(\hat{x}^k) \\ &= f(\hat{x}^k) + \langle g^i, y - x \rangle - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - x \pm \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - e^i(\hat{x}^k) \end{aligned}$$

$$\implies g^i \in \partial_{e^i(\hat{x}^k)} f(\hat{x}^k)$$

$$e^i(\hat{x}^k) := f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle \geq 0$$

## The transportation formula

How to express subgradients at  $x^i$  as  $\varepsilon$ -subgradients at  $\hat{x}^k$ ?

$$g^i \in \partial f(x^i) \quad \text{if and only if, for all } y \in \mathbb{R}^n$$

$$\begin{aligned} f(y) &\geq f(x^i) + \langle g^i, y - x^i \rangle \\ &= f(x^i) + \langle g^i, y - x \rangle \pm f(\hat{x}^k) \\ &= f(\hat{x}^k) + \langle g^i, y - x \rangle - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - x \pm \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - e^i(\hat{x}^k) \end{aligned}$$

$$\implies g^i \in \partial_{e^i(\hat{x}^k)} f(\hat{x}^k)$$

$$e^i(\hat{x}^k) := f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle \geq 0$$

## The transportation formula

How to express subgradients at  $x^i$  as  $\varepsilon$ -subgradients at  $\hat{x}^k$ ?

$$g^i \in \partial f(x^i) \quad \text{if and only if, for all } y \in \mathbb{R}^n$$

$$\begin{aligned} f(y) &\geq f(x^i) + \langle g^i, y - x^i \rangle \\ &= f(x^i) + \langle g^i, y - x \rangle \pm f(\hat{x}^k) \\ &= f(\hat{x}^k) + \langle g^i, y - x \rangle - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - x \pm \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - e^i(\hat{x}^k) \end{aligned}$$

$$\implies g^i \in \partial_{e^i(\hat{x}^k)} f(\hat{x}^k)$$

$$e^i(\hat{x}^k) := f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle \geq 0$$

## The transportation formula

How to express subgradients at  $x^i$  as  $\varepsilon$ -subgradients at  $\hat{x}^k$ ?

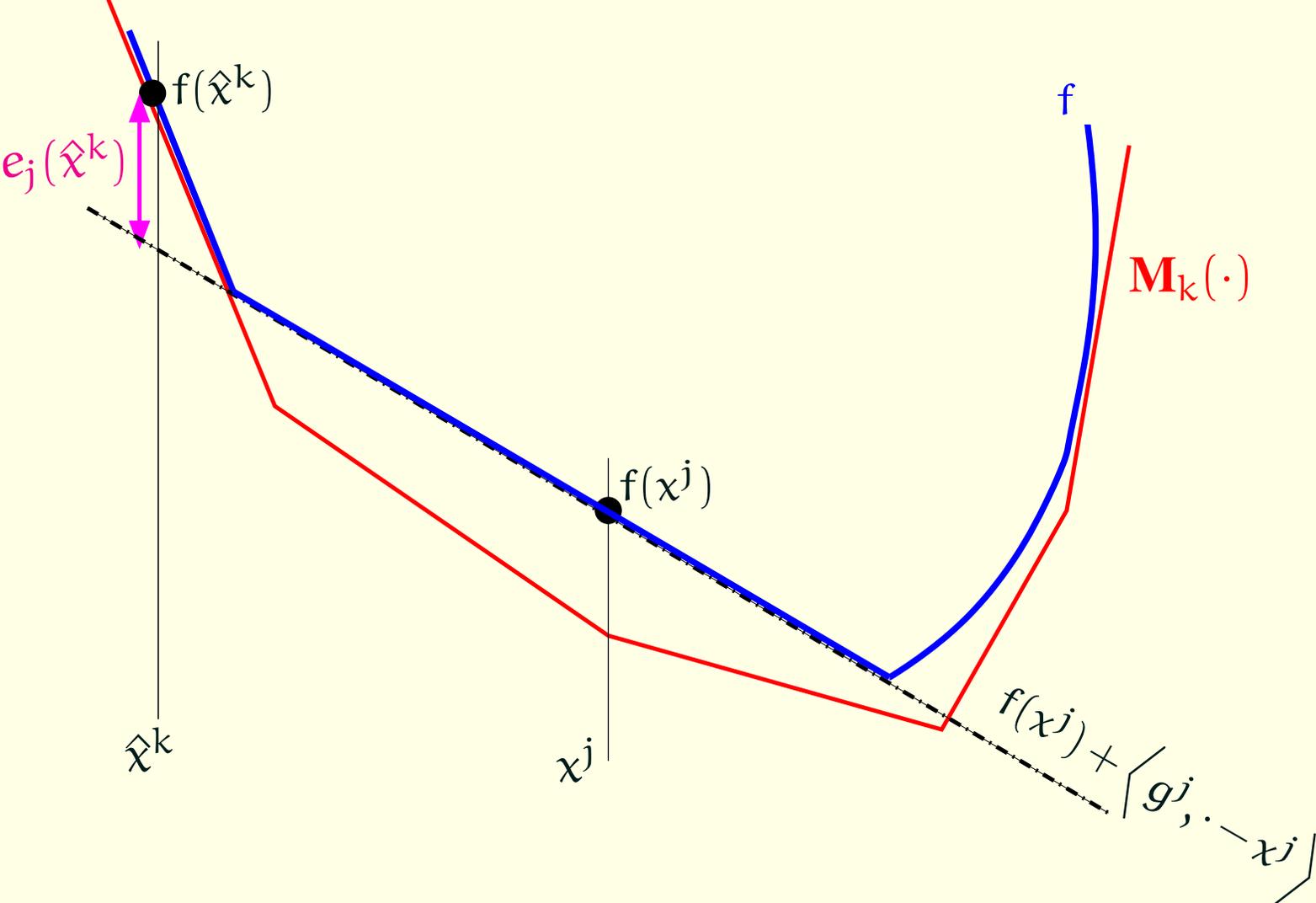
$g^i \in \partial f(x^i)$  if and only if, for all  $y \in \mathbb{R}^n$

$$\begin{aligned} f(y) &\geq f(x^i) + \langle g^i, y - x^i \rangle \\ &= f(x^i) + \langle g^i, y - x \rangle \pm f(\hat{x}^k) \\ &= f(\hat{x}^k) + \langle g^i, y - x \rangle - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - x \pm \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i)) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - (f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle) \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - e^i(\hat{x}^k) \end{aligned}$$

$$\implies g^i \in \partial_{e^i(\hat{x}^k)} f(\hat{x}^k)$$

$$e^i(\hat{x}^k) := f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle \geq 0$$

# Linearization errors



## The $\varepsilon$ -subdifferential and bundle methods

We collect the black-box

$x^i, i = 1, 2, \dots, k$ , so that at iteration  $k$  we can define a **bundle** of information, centered at a special iterate  $\hat{x}^k \in \{x^i\}$

$$\mathcal{B}^k := \left( \begin{array}{l} e^i(\hat{x}^k) = f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle \\ g^i \in \partial_{e^i(\hat{x}^k)} f(\hat{x}^k) \end{array} \right)$$

## The $\varepsilon$ -subdifferential and bundle methods

We collect the black-box

$x^i, i = 1, 2, \dots, k$ , so that at iteration  $k$  we can define a **bundle** of information, centered at a special iterate  $\hat{x}^k \in \{x^i\}$

$$\mathcal{B}^k := \left( \begin{array}{l} e^i(\hat{x}^k) = f(\hat{x}^k) - f(x^i) - \langle g^i, \hat{x}^k - x^i \rangle \\ g^i \in \partial_{e^i(\hat{x}^k)} f(\hat{x}^k) \end{array} \right)$$

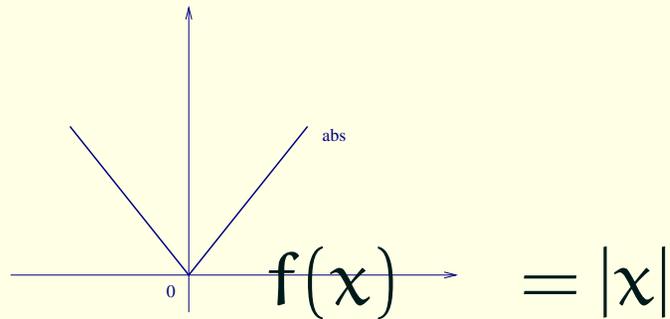
A suitable convex combination

$$\varepsilon^k := \sum_{i \in \mathcal{B}^k} \alpha^i e^i(\hat{x}^k) \text{ and } G^k := \sum_{i \in \mathcal{B}^k} \alpha^i g^i$$

will eventually satisfy the optimality condition!

## Why special NSO methods?

Smooth optimization techniques **do not work**



$$|\nabla f(x^k)| = 1, \forall x^k \neq 0 \quad \partial f(0) = [-1, 1]$$

Smooth stopping test **fails**:

$$|\nabla f(x^k)| \leq \mathbf{TOL} \quad (\Leftrightarrow |g(x^k)| \leq \mathbf{TOL})$$

## Why special NSO methods?

Smooth optimization techniques **do not work**

Smooth approximations of derivatives by finite differences **fail**

For  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  defined by  $f(x) = \max(x_1, x_2, x_3)$

$\partial f(0) = ?$

Forward finite difference  $\frac{f(x+\Delta x) - f(x)}{\Delta x}$

Central finite difference  $\frac{f(x+\Delta x) - f(x-\Delta)}{2\Delta x}$

## Why special NSO methods?

Smooth optimization techniques **do not work**

Smooth approximations of derivatives by finite differences **fail**

For  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$  defined by  $f(x) = \max(x_1, x_2, x_3)$   
 $\partial f(0) = ?$

Forward finite difference  $\frac{f(x+\Delta x) - f(x)}{\Delta x} = (1, 1, 1)$

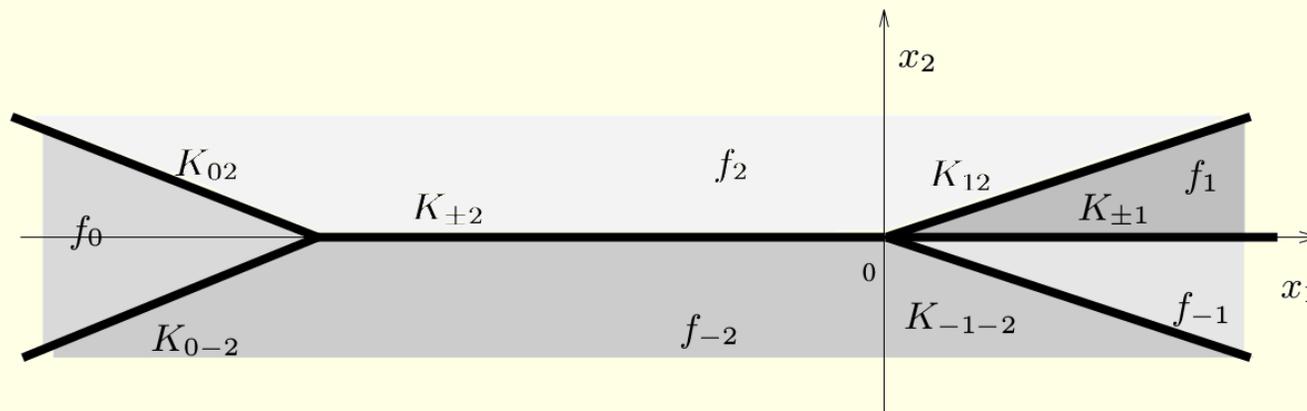
Central finite difference  $\frac{f(x+\Delta x) - f(x-\Delta)}{2\Delta x} = (\frac{1}{2}, \frac{1}{2}, \frac{1}{2})$

**none of them in the subdifferential!**

## Why special NSO methods?

Smooth optimization techniques **do not work**

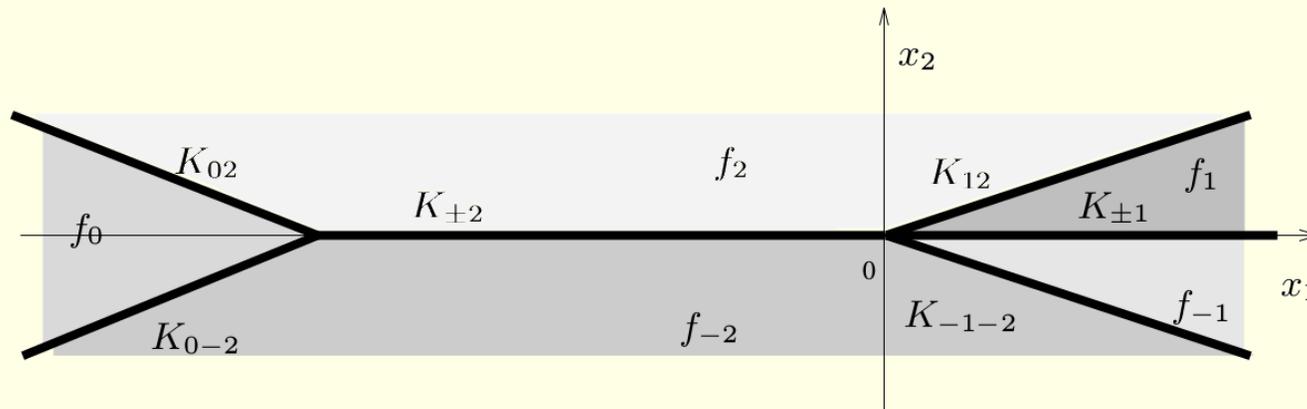
Linesearches get trapped in kinks and **fail**



# Why special NSO methods?

Smooth optimization techniques **do not work**

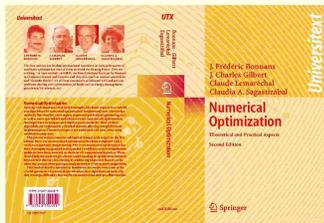
Linesearches get trapped in kinks and **fail**



Example 9.1

“Instability of steepest

descent”



## Why special NSO methods?

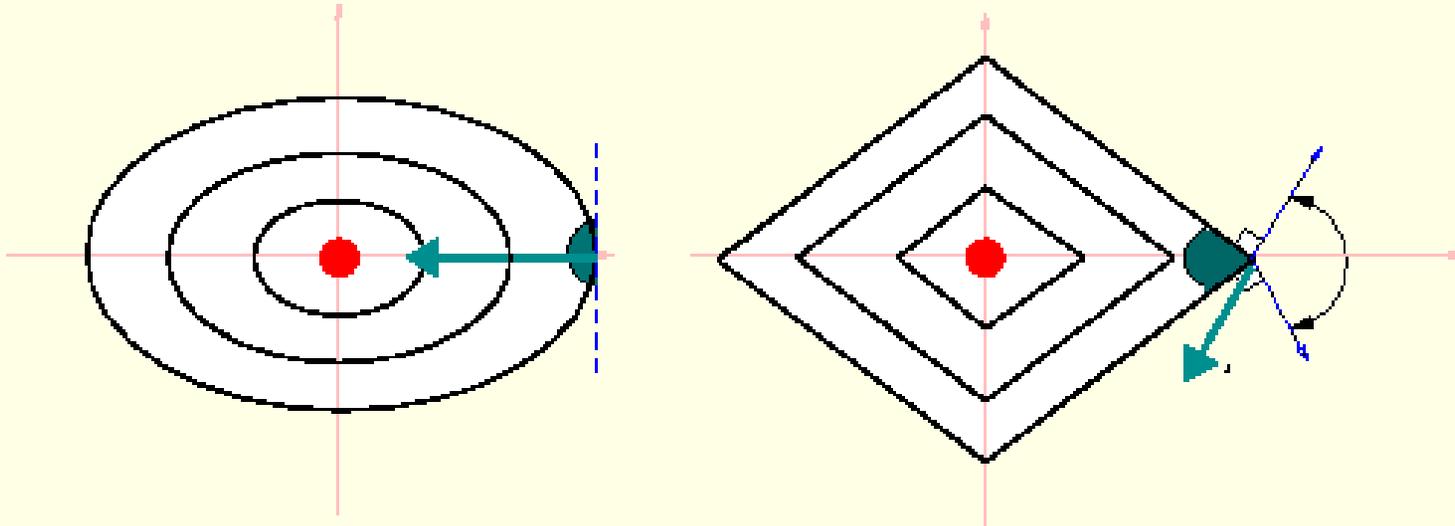
Smooth optimization techniques **do not work**

$-g(x^k)$  may **not** provide descent

## Why special NSO methods?

Smooth optimization techniques **do not work**

$-g(x^k)$  may **not** provide descent



## Why special NSO methods?

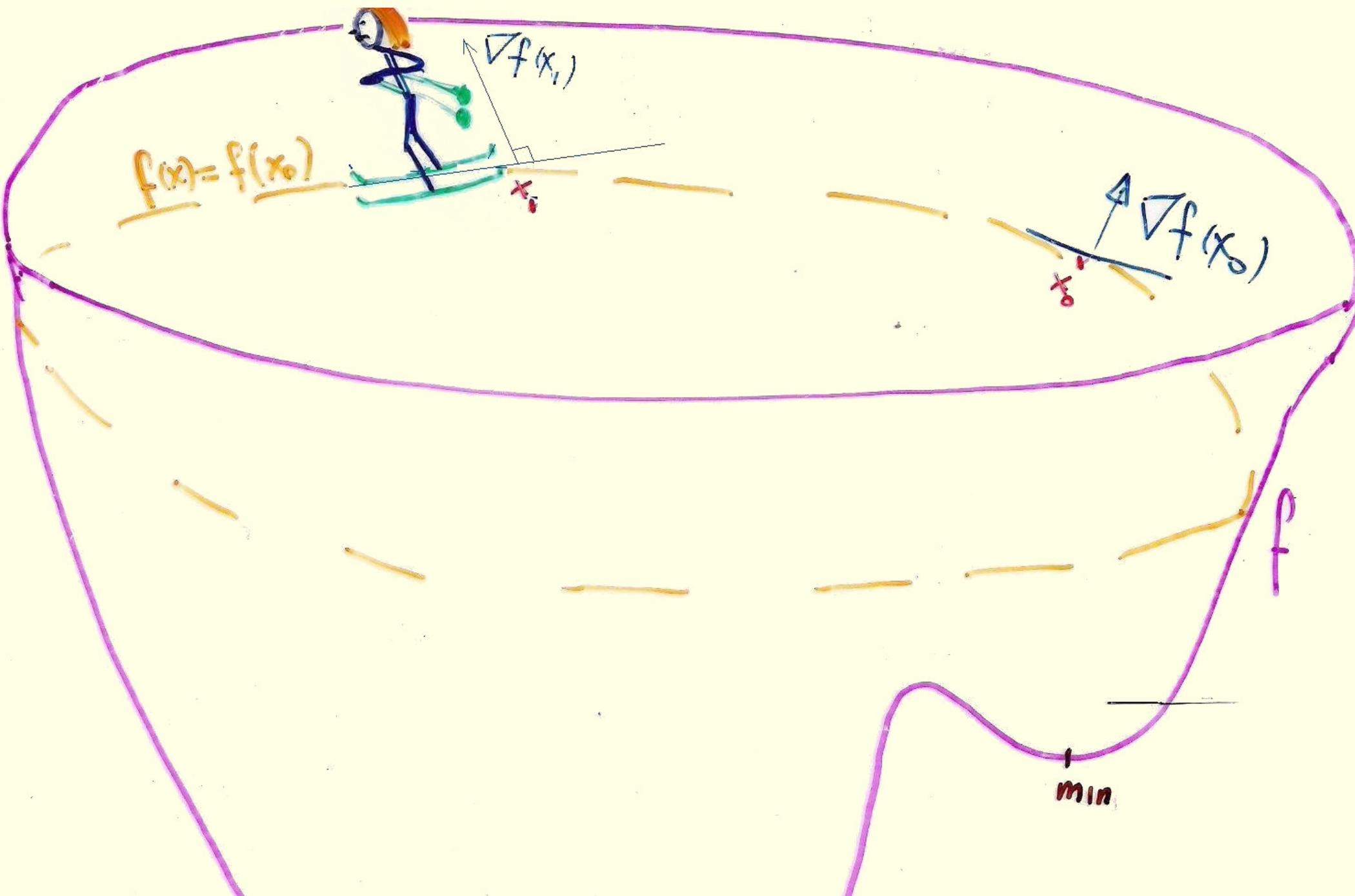
Smooth optimization techniques **do not work**

Smooth stopping test **fails**

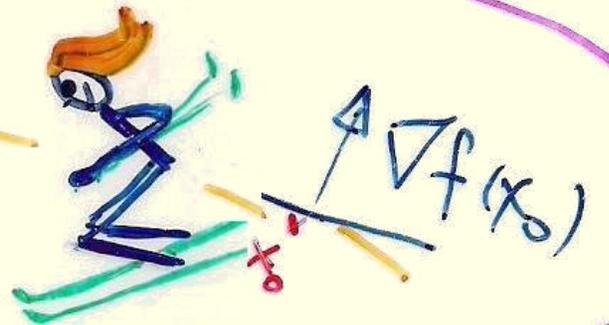
Finite difference approximations **fail**

Linesearches get trapped in kinks and **fail**

Direction opposite to a subgradient may **increase**  
the functional values



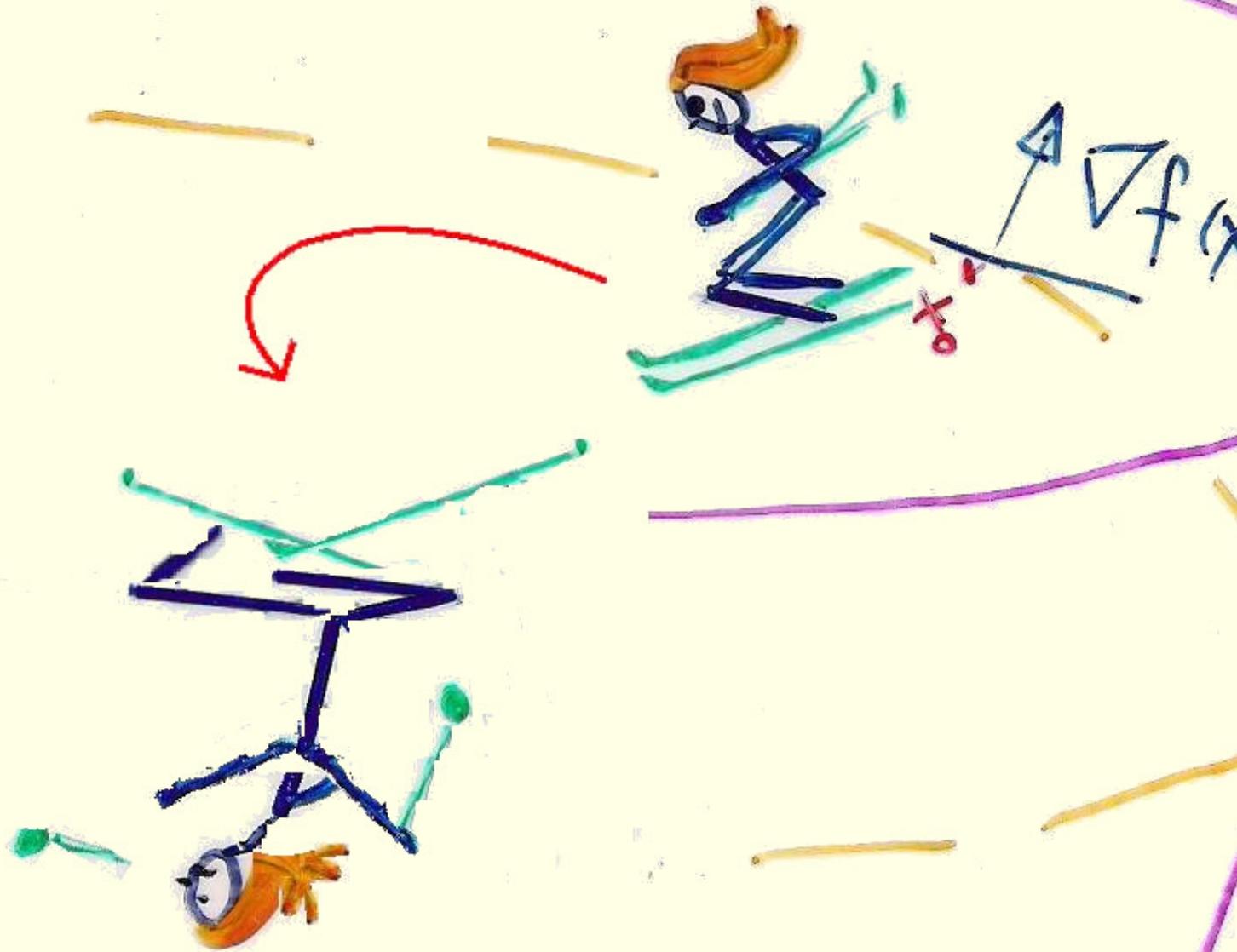
$f(x) = f(x_0)$



f

min

In NSO  
the skier  
is blind



## Bundle Methods

$$\begin{aligned} \text{WANT: } p &= \text{prox}_t^f(x) && \iff p = \arg \min f(y) + \frac{1}{2t} \|y - x\|_2^2 \\ \text{HAVE: } q^k &= \text{prox}_{t_k}^{M_k}(x) && \iff q^k = \arg \min M_k(y) + \frac{1}{2t_k} \|y - x^k\|_2^2 \\ &&& \iff 0 = G^k + \frac{1}{t_k} (q^k - x^k) \\ &&& \text{for } G^k \in \partial M_k(q^k) \\ &&& \iff G^k \in \partial_{\varepsilon_k} f(x) \\ &&& \text{for } \varepsilon_k = f(x) - M_k(q^k) - t_k \|G^k\|_2^2 \end{aligned}$$

## Bundle Methods

$$\begin{aligned} \text{WANT: } \mathbf{p} &= \text{prox}_t^f(\mathbf{x}) && \iff \mathbf{p} = \arg \min \mathbf{f}(\mathbf{y}) + \frac{1}{2t} \|\mathbf{y} - \mathbf{x}\|_2^2 \\ \text{HAVE: } \mathbf{q}^k &= \text{prox}_{t_k}^{M_k}(\mathbf{x}) && \iff \mathbf{q}^k = \arg \min M_k(\mathbf{y}) + \frac{1}{2t_k} \|\mathbf{y} - \mathbf{x}^k\|_2^2 \\ &&& \iff 0 = \mathbf{G}^k + \frac{1}{t_k} (\mathbf{q}^k - \mathbf{x}^k) \\ &&& \text{for } \mathbf{G}^k \in \partial M_k(\mathbf{q}^k) \\ &&& \iff \mathbf{G}^k \in \partial_{\varepsilon_k} f(\mathbf{x}) \\ &&& \text{for } \varepsilon_k = f(\mathbf{x}) - M_k(\mathbf{q}^k) - \frac{1}{2t_k} \|\mathbf{G}^k\|_2^2 \end{aligned}$$

Two subsequences

- Iterates giving sufficiently good approximal points
- Iterates just helping the optimization process

## Bundle Methods

$$\text{HAVE: } q^k = \text{prox}_{t_k}^{M_k}(\mathbf{x}) = \mathbf{x}^k + t_k G^k \quad G^k \in \partial_{\varepsilon_k} f(\mathbf{x})$$

for  $\varepsilon_k = f(\mathbf{x}) - M_k(q^k) - t_k \|G^k\|_2^2$

Two subsequences

- Iterates giving sufficiently good approximal points

**moving towards minimum**

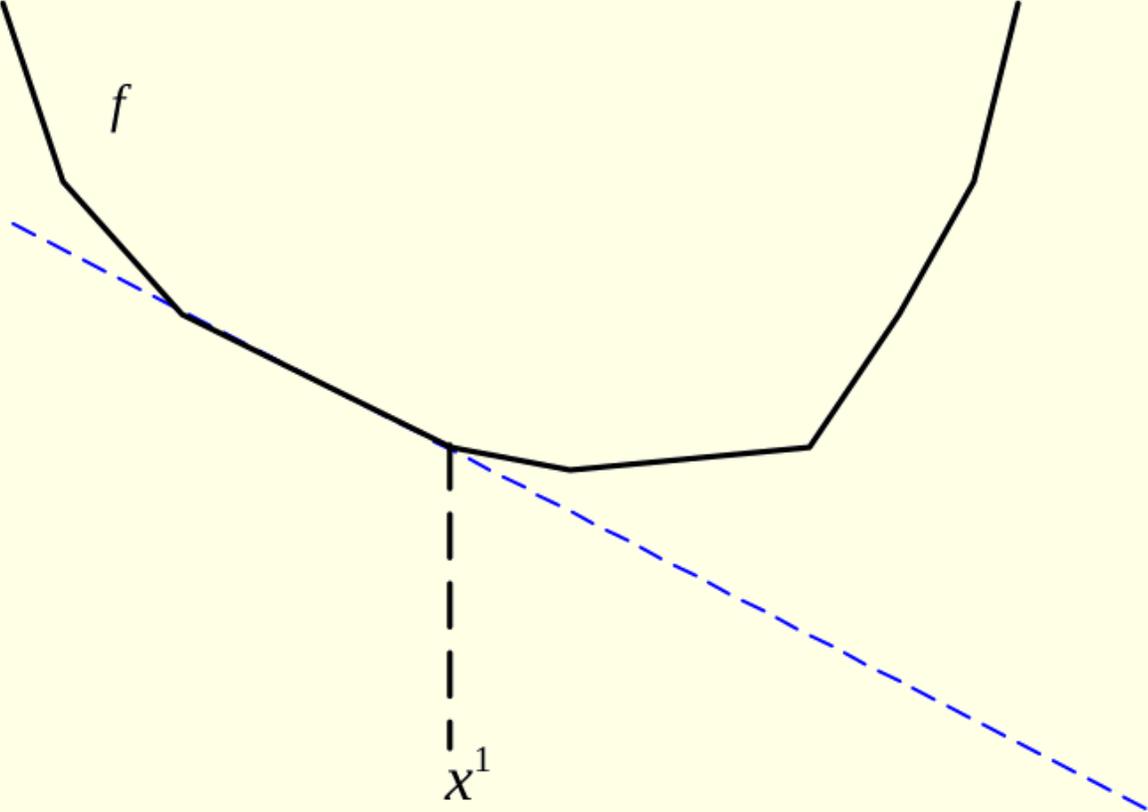
**in a manner that makes  $\delta_k := \varepsilon_k + t_k \|G^k\|_2^2 \rightarrow 0$**

**(serious)**

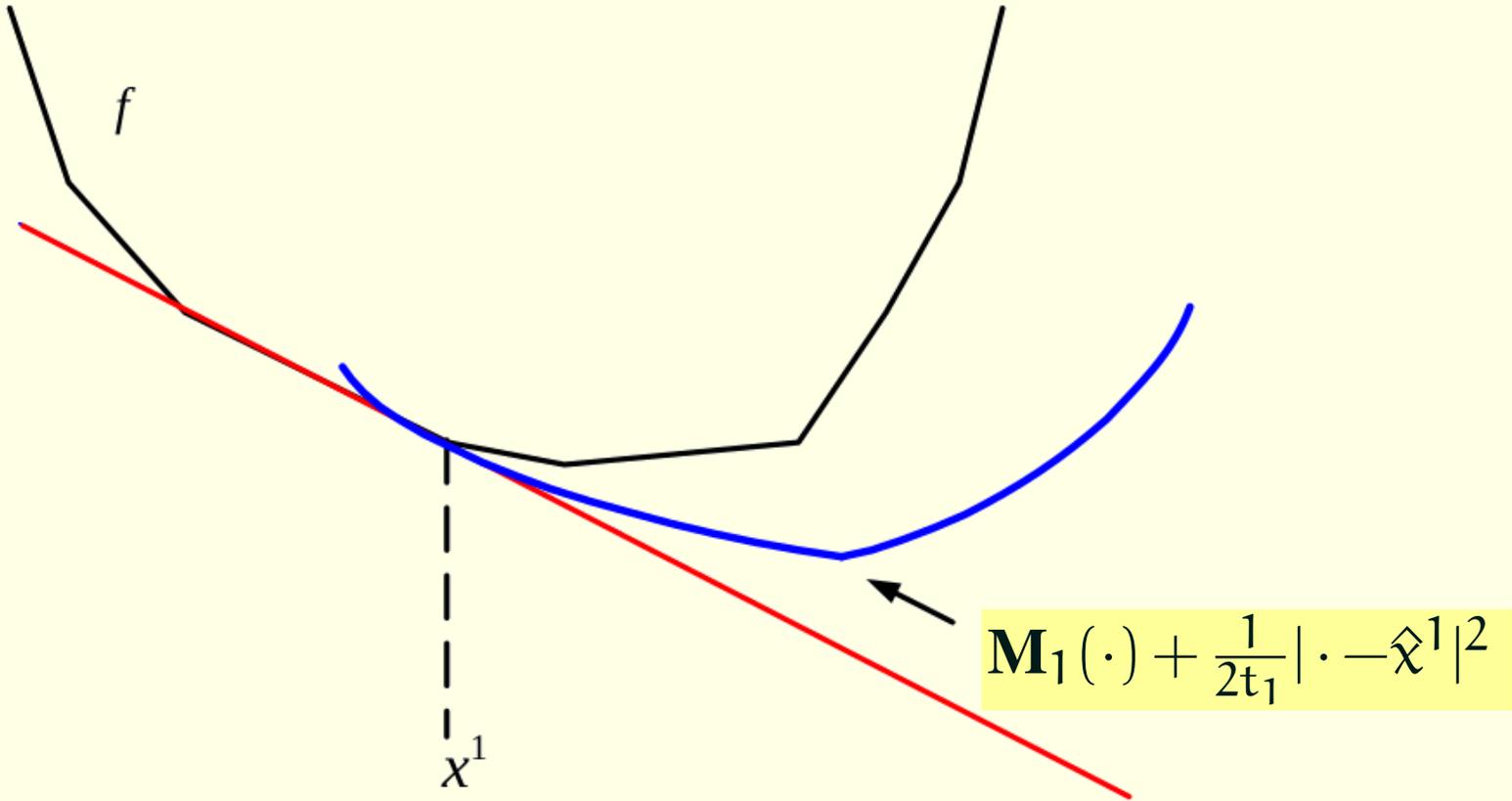
- Iterates just helping the optimization process

**CL93 eventually applies (null)**

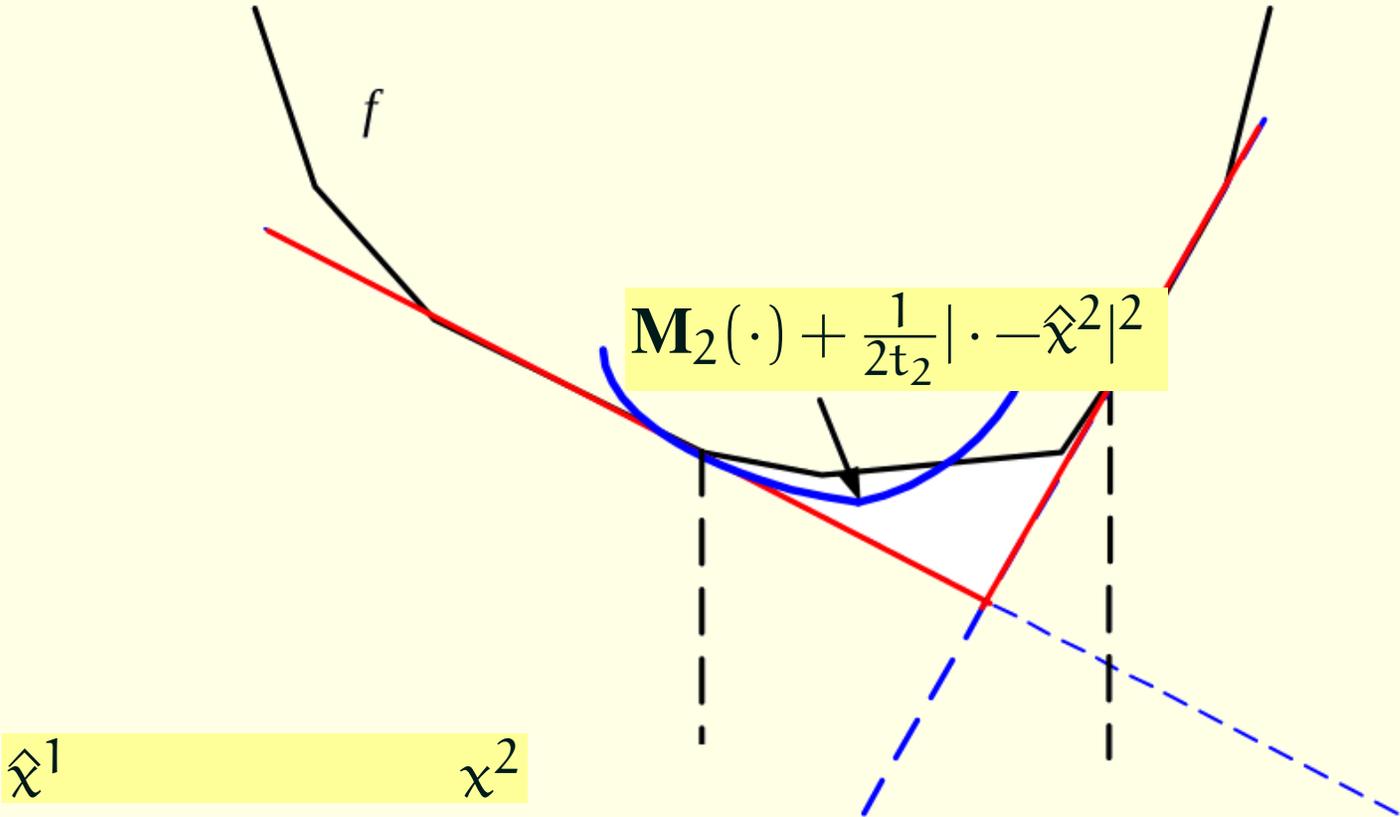
# Bundle Methods



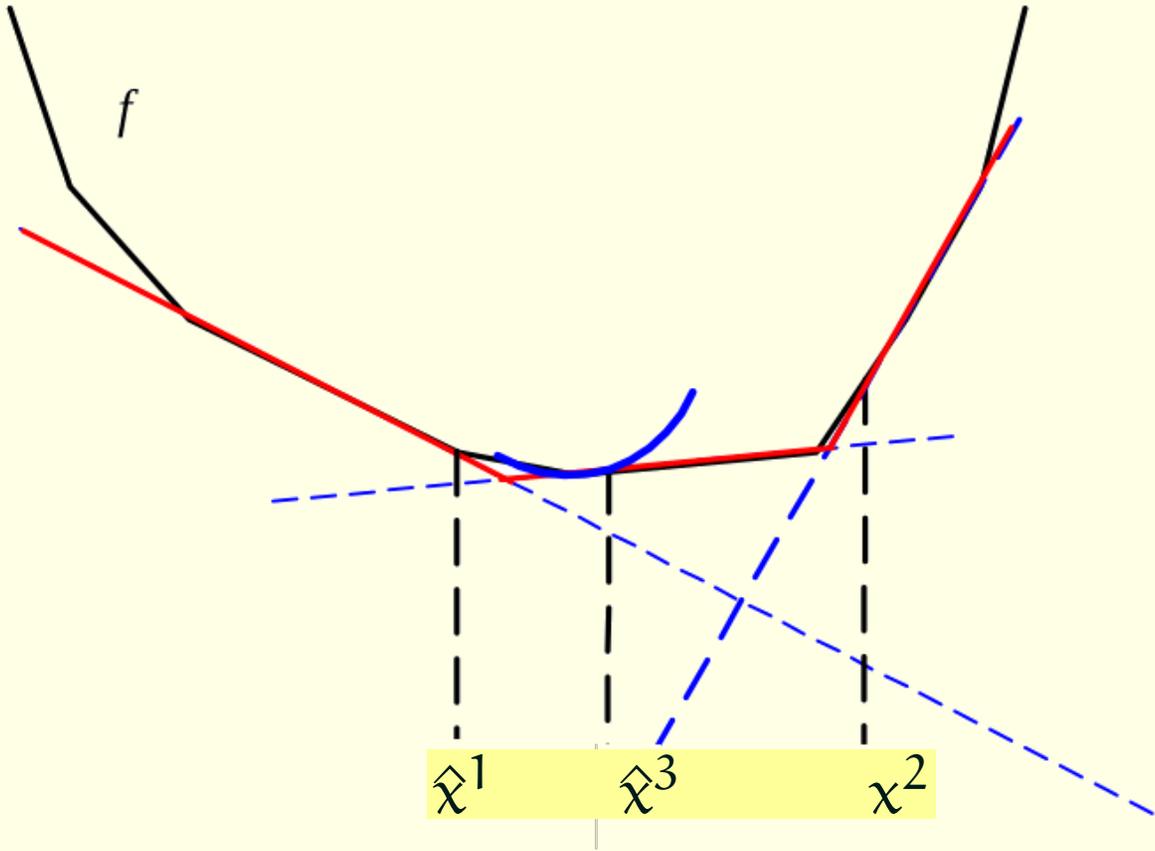
# Bundle Methods



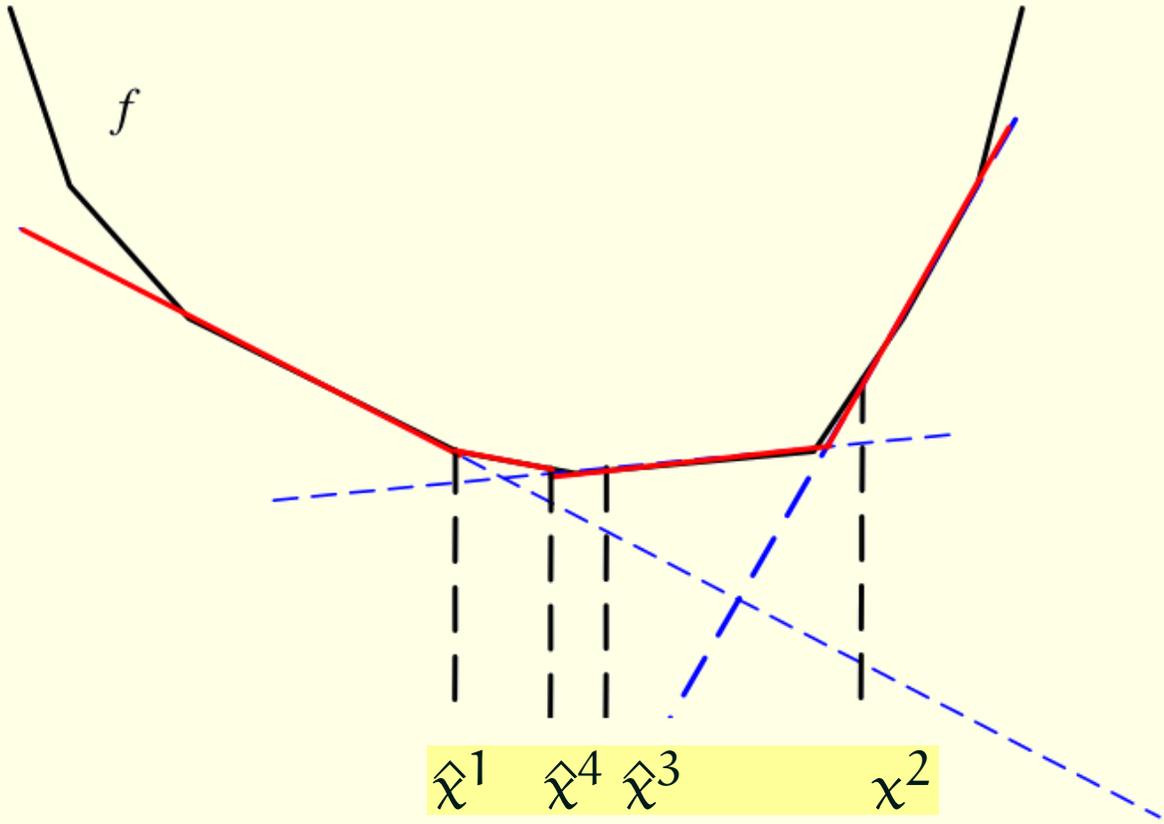
# Bundle Methods



# Bundle Methods



# Bundle Methods



## Bundle Methods

- 0 Choose  $x^1$ , set  $k = 1$ , and let  $\hat{x}^1 = x^1$ .
- 1 Compute  $x^{k+1} = \arg \min \mathbf{M}_k(x) + \frac{1}{2t_k}|x - \hat{x}^k|^2$
- 2 If  $\delta_k := f(\hat{x}^k) - \mathbf{M}_k(x^{k+1}) \leq \text{tol}$  STOP
- 3 Call the oracle at  $x^{k+1}$ . If  
 $f(x^{k+1}) \leq f(\hat{x}^k) - m\delta_k$ , set  $\hat{x}^{k+1} = x^{k+1}$  ●  
(Serious Step) Otherwise, maintain  $\hat{x}^{k+1} = \hat{x}^k$   
(Null Step)
- 4 Define  $\mathbf{M}_{k+1}$ ,  $t_{k+1}$ , make  $k = k + 1$ , and loop to 1.

## Bundle Methods: selection mechanism

$$\mathbf{M}_{k+1}(\cdot) = \max \left( \mathbf{M}_k(\cdot), f^k + \langle g^k, \cdot - x^k \rangle \right),$$

now the choice of the new model is more flexible:

$$x^{k+1} \in \arg \min \mathbf{M}_k(x) + \frac{1}{2t_k} |x - \hat{x}^k|^2$$

with  $\mathbf{M}_k(x) = \max_{i \leq k} \{f^i + \langle g^i, x - x^i \rangle\}$  is equivalent to a QP:

$$\begin{cases} \min_{r \in \mathbb{R}, x \in \mathbb{R}^n} & r + \frac{1}{2t_k} |x - \hat{x}^k|^2 \\ \text{s.t.} & r \geq f^i + \langle g^i, x - x^i \rangle \text{ for } i \leq k \end{cases}$$

A posteriori, the solution remains the same if ...

## Bundle Methods: selection mechanism

$$\mathbf{M}_{k+1}(\cdot) = \max \left( \mathbf{M}_k(\cdot), f^k + \langle g^k, \cdot - x^k \rangle \right),$$

now the choice of the new model is more flexible:

$$x^{k+1} \in \arg \min \mathbf{M}_k(x) + \frac{1}{2t_k} |x - \hat{x}^k|^2$$

with  $\mathbf{M}_k(x) = \max_{i \leq k} \{f^i + \langle g^i, x - x^i \rangle\}$  is equivalent to a QP:

$$\begin{cases} \min_{r \in \mathbb{R}, x \in \mathbb{R}^n} & r + \frac{1}{2t_k} |x - \hat{x}^k|^2 \\ \text{s.t.} & r \geq f^i + \langle g^i, x - x^i \rangle \text{ for } \mathbf{i} \leq \mathbf{k} \end{cases}$$

A posteriori, the solution remains the same if all, or

...

## Bundle Methods: selection mechanism

$$\mathbf{M}_{k+1}(\cdot) = \max \left( \mathbf{M}_k(\cdot), f^k + \langle g^k, \cdot - x^k \rangle \right),$$

now the choice of the new model is more flexible:

$$x^{k+1} \in \arg \min \mathbf{M}_k(x) + \frac{1}{2t_k} |x - \hat{x}^k|^2$$

with  $\mathbf{M}_k(x) = \max_{i \leq k} \{f^i + \langle g^i, x - x^i \rangle\}$  is equivalent to a QP:

$$\begin{cases} \min_{r \in \mathbb{R}, x \in \mathbb{R}^n} & r + \frac{1}{2t_k} |x - \hat{x}^k|^2 \\ \text{s.t.} & r \geq f^i + \langle g^i, x - x^i \rangle \text{ for active } i\text{'s} \end{cases}$$

A posteriori, the solution remains the same if all, or active, or ...

## Bundle Methods: selection mechanism

$$\mathbf{M}_{k+1}(\cdot) = \max \left( \mathbf{M}_k(\cdot), f^k + \langle g^k, \cdot - x^k \rangle \right),$$

now the choice of the new model is more flexible:

$$x^{k+1} \in \arg \min \mathbf{M}_k(x) + \frac{1}{2t_k} |x - \hat{x}^k|^2$$

with  $\mathbf{M}_k(x) = \max_{i \leq k} \{f^i + \langle g^i, x - x^i \rangle\}$  is equivalent to a QP:

$$\begin{cases} \min_{r \in \mathbb{R}, x \in \mathbb{R}^n} & r + \frac{1}{2t_k} |x - \hat{x}^k|^2 \\ \text{s.t.} & r \geq \sum_i \bar{\alpha}^i \left( f^i + \langle g^i, x - x^i \rangle \right) \end{cases} \quad \mathbf{A}$$

posteriori, the solution remains the same if all, or active, or the **optimal convex combination**

## Bundle Methods: selection mechanism

$$\mathbf{M}_{k+1}(\cdot) = \max \left( \mathbf{M}_k(\cdot), f^k + \langle g^k, \cdot - x^k \rangle \right),$$

now the choice of the new model is more flexible:

$$x^{k+1} \in \arg \min \mathbf{M}_k(x) + \frac{1}{2t_k} |x - \hat{x}^k|^2$$

with  $\mathbf{M}_k(x) = \max_{i \leq k} \{f^i + \langle g^i, x - x^i \rangle\}$  is equivalent to a QP:

$$\begin{cases} \min_{r \in \mathbb{R}, x \in \mathbb{R}^n} & r + \frac{1}{2t_k} |x - \hat{x}^k|^2 \\ \text{s.t.} & r \geq \sum_i \bar{\alpha}^i \left( f^i + \langle g^i, x - x^i \rangle \right) \end{cases}$$

A posteriori, the solution remains the same if all, or active, or the optimal convex combination are kept

## Bundle Methods: next model options

$$\mathbf{M}_{k+1}(\cdot) = \max\left(\mathbf{M}_k(\cdot), f^k + \left\langle g^k, \cdot - x^k \right\rangle\right)$$

or

$$\mathbf{M}_{k+1}(\cdot) = \max\left(\max_{\text{active}}, f^k + \left\langle g^k, \cdot - x^k \right\rangle\right)$$

or

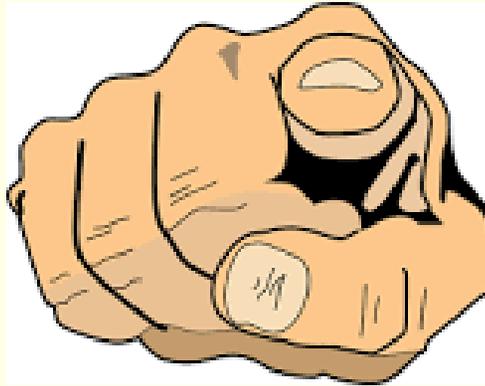
$$\mathbf{M}_{k+1}(\cdot) = \max\left(\text{aggregate}, f^k + \left\langle g^k, \cdot - x^k \right\rangle\right)$$

Same QP solution if all, or active, or the optimal convex combination

**aggregate=full Bundle Compression:** QP with only 2 constraints  
(but slows down the overall process)

# The cutting-plane model

You told us



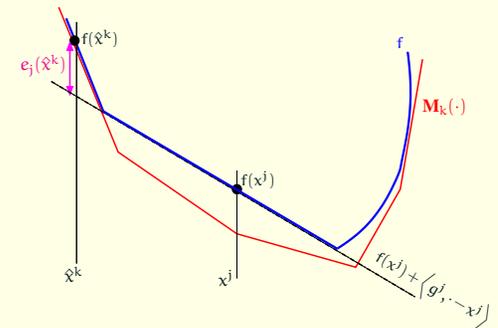
we were going to use a bundle  $\mathcal{B}_k$  composed by linearization errors and  $\varepsilon$ -subgradients at  $\hat{x}^k$ , but the model uses  $f^i$  and  $g^i \in \partial f(x^i)$



# Rewriting the cutting-plane model

The transportation formula centers the  $i$ th linearization in the serious iterate

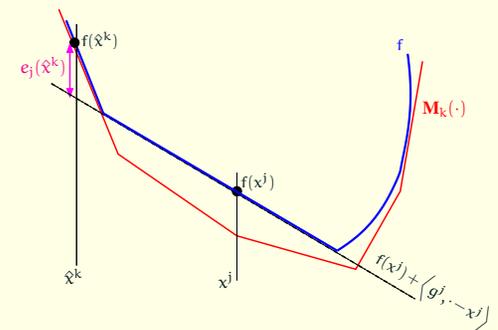
$$\begin{aligned} f(y) &\geq f(x^i) + \langle g^i, y - x^i \rangle \\ &= f(\hat{x}^k) + \langle g^i, y - \hat{x}^k \rangle - e^i(\hat{x}^k) \end{aligned}$$



# Rewriting the cutting-plane model

The transportation formula centers the  $i$ th linearization in the serious iterate

$$\begin{aligned} f(\mathbf{y}) &\geq f(\mathbf{x}^i) + \langle \mathbf{g}^i, \mathbf{y} - \mathbf{x}^i \rangle \\ &= f(\hat{\mathbf{x}}^k) + \langle \mathbf{g}^i, \mathbf{y} - \hat{\mathbf{x}}^k \rangle - e^i(\hat{\mathbf{x}}^k) \end{aligned}$$



this translates into the model as follows

$$\begin{aligned} \mathbf{M}_k(\mathbf{y}) &= \max \left\{ f(\mathbf{x}^i) + \langle \mathbf{g}^i, \mathbf{y} - \mathbf{x}^i \rangle \quad : i \in \mathcal{B}_k \right\} \\ &= \max \left\{ f(\hat{\mathbf{x}}^k) + \langle \mathbf{g}^i, \mathbf{y} - \hat{\mathbf{x}}^k \rangle - e^i(\hat{\mathbf{x}}^k) \quad : i \in \mathcal{B}_k \right\} \\ &= f(\hat{\mathbf{x}}^k) + \max \left\{ \langle \mathbf{g}^i, \mathbf{y} - \hat{\mathbf{x}}^k \rangle - e^i(\hat{\mathbf{x}}^k) \quad : i \in \mathcal{B}_k \right\} \end{aligned}$$

## Bundle Method

- 0 Choose  $x^1$ , set  $k = 1$ , and let  $\hat{x}^1 = x^1$ .
- 1 Compute  $x^{k+1} = \arg \min \mathbf{M}_k(x) + \frac{1}{2t_k}|x - \hat{x}^k|^2$
- 2 If  $\delta_k := f(\hat{x}^k) - \mathbf{M}_k(x^{k+1}) \leq \text{tol}$  **STOP**
- 3 Call the oracle at  $x^{k+1}$ . If  
 $f(x^{k+1}) \leq f(\hat{x}^k) - m\delta_k$ , set  $\hat{x}^{k+1} = x^{k+1}$

Otherwise, maintain  $\hat{x}^{k+1} = \hat{x}^k$

- 4 Define  $\mathbf{M}_{k+1}$ ,  $t_{k+1}$ , make  $k = k + 1$ , and loop to 1.

## Bundle Method

- 0 Choose  $x^1$ , set  $k = 1$ , and let  $\hat{x}^1 = x^1$ .
- 1 Compute  $x^{k+1} = \arg \min \mathbf{M}_k(x) + \frac{1}{2t_k}|x - \hat{x}^k|^2$
- 2 If  $\delta_k := f(\hat{x}^k) - \mathbf{M}_k(x^{k+1}) \leq \text{tol}$  **STOP**
- 3 Call the oracle at  $x^{k+1}$ . If  
 $f(x^{k+1}) \leq f(\hat{x}^k) - m\delta_k$ , set  $\hat{x}^{k+1} = x^{k+1}$   
(Serious Step)  $k \in \mathbf{K}_S$   
Otherwise, maintain  $\hat{x}^{k+1} = \hat{x}^k$   
(Null Step)  $k \in \mathbf{K}_N$
- 4 Define  $\mathbf{M}_{k+1}$ ,  $t_{k+1}$ , make  $k = k + 1$ , and loop to 1.

## Bundle Method

When  $k \rightarrow \infty$ , the algorithm generates two subsequences. Convergence analysis addresses the mutually exclusive situations

- either the SS subsequence is infinite

$$K_\infty := \{k \in K_S\}$$

- or there is a last SS, followed by infinitely many null steps

## Bundle Method

When  $k \rightarrow \infty$ , the algorithm generates two subsequences. Convergence analysis addresses the mutually exclusive situations

- either the SS subsequence is infinite

$$\mathbf{K}_\infty := \{\mathbf{k} \in \mathbf{K}_S\}$$

- or there is a last SS, followed by infinitely many null steps

$$\mathbf{K}_\infty := \{\mathbf{k} \in \mathbf{K}_N : \mathbf{k} \geq \text{last SS}\}$$

## Bundle Method

When  $k \rightarrow \infty$ , the algorithm generates two subsequences. Convergence analysis addresses the mutually exclusive situations

- either the SS subsequence is infinite

$$\mathbf{K}_\infty := \{\mathbf{k} \in \mathbf{K}_S\} \text{ (limit point minimizes } f)$$

- or there is a last SS, followed by infinitely many null steps

$$\mathbf{K}_\infty := \{\mathbf{k} \in \mathbf{K}_N : \mathbf{k} \geq \text{last SS}\}$$

(last SS minimizes  $f$  and  $\text{null} \rightarrow \text{last SS}$ )

## Equivalent QPs

1. Given  $t_k$ , the stepsize parameter of the proximal bundle method, with QP subproblem given by

$$(\text{PB})_k \quad \min \mathbf{M}_k(\mathbf{x}) + \frac{1}{2t_k} |\mathbf{x} - \hat{\mathbf{x}}^k|^2$$

2. Given  $\Delta_k$ , the radius parameter of the trust-region bundle method, with QP subproblem given by

$$(\text{TRB})_k \quad \begin{cases} \min & \mathbf{M}_k(\mathbf{x}) \\ \text{s.t.} & |\mathbf{x} - \hat{\mathbf{x}}^k|^2 \leq \Delta_k \end{cases}$$

3. Given  $\ell_k$ , the level parameter of the level bundle method,

with QP subproblem given by

$$(\text{LB})_k \quad \begin{cases} \min & \frac{1}{2}|\mathbf{x} - \hat{\mathbf{x}}^k|^2 \\ \text{s.t.} & \mathbf{M}_k(\mathbf{x}) \leq \ell_k \end{cases}$$

Show that

1. given  $t_k$ , there exists  $\Delta_k$  such that if  $\mathbf{x}^{k+1}$  solves  $(\text{PB})_k$ , then  $\mathbf{x}^{k+1}$  solves  $(\text{TRB})_k$ .
2. given  $\Delta_k$ , there exists  $\ell_k$  such that if  $\mathbf{x}^{k+1}$  solves  $(\text{TRB})_k$ , then  $\mathbf{x}^{k+1}$  solves  $(\text{LB})_k$ .
3. given  $\ell_k$ , there exists  $t_k$  such that if  $\mathbf{x}^{k+1}$  solves  $(\text{LB})_k$ , then  $\mathbf{x}^{k+1}$  solves  $(\text{PB})_k$ .

## Theorem $\mathbf{K}_\infty := \{\mathbf{k} \in \mathbf{K}_S\}$

Suppose the bundle method loops forever and there are infinitely many serious steps. Either the solution set of  $\min f$  is empty and  $f(\hat{\mathbf{x}}^k) \searrow -\infty$  or the following holds

- (i)  $\lim_{k \in \mathbf{K}_S} \delta_k = 0$  and  $\lim_{k \in \mathbf{K}_S} \varepsilon_k = 0$ .
- (ii) If the stepsizes are chosen so that  $\sum_{k \in \mathbf{K}_S} t_k = +\infty$  then  $\{\hat{\mathbf{x}}^k\}$  is a minimizing sequence.
- (iii) If, in addition,  $t_k \leq t^{\text{up}}$  for all  $k \in \mathbf{K}_S$ , then the subsequence  $\{\hat{\mathbf{x}}^k\}$  is bounded. In this case, any limit point  $\mathbf{x}^\infty$  minimizes  $f$  and the whole sequence converges to  $\mathbf{x}^\infty$

**Theorem**  $\mathbf{K}_\infty := \left\{ \mathbf{k} \in \mathbf{K}_N \geq \hat{\mathbf{k}} \right\}$

Suppose the bundle method loops forever and there are infinitely many null steps after a last serious one, denoted by  $\hat{\mathbf{x}}$  and generated at iteration  $\hat{\mathbf{k}}$ . Suppose stepsizes are chosen so that

$$t_{l_0} \leq t_{k+1} \leq t_k \quad \text{for all } k \in \mathbf{K}_\infty$$

The following holds

1. The sequence  $\{\mathbf{x}^{k+1}\}$  is bounded
2.  $\lim_{k \in \mathbf{K}_\infty} \mathbf{M}_k(\mathbf{x}^{k+1}) = f(\hat{\mathbf{x}})$
3.  $\hat{\mathbf{x}}$  minimizes  $f$
4.  $\lim_{k \in \mathbf{K}_\infty} \mathbf{x}^{k+1} = \hat{\mathbf{x}}$

## Model requirements

1.  $\mathbf{M}_k \leq f$

2. If  $k$  was declared a null step

a)  $\mathbf{M}_{k+1}(x) \geq f^{k+1} + \langle g^{k+1}, x - x^{k+1} \rangle$

b)  $\mathbf{M}_{k+1}(x) \geq A_k(x) = \mathbf{M}_k(x^{k+1}) + \langle G^k, x - x^{k+1} \rangle$

## Model requirements

1.  $\mathbf{M}_k \leq f$

2. If  $k$  was declared a null step

a)  $\mathbf{M}_{k+1}(x) \geq f^{k+1} + \langle g^{k+1}, x - x^{k+1} \rangle$

b)  $\mathbf{M}_{k+1}(x) \geq A_k(x) = \mathbf{M}_k(x^{k+1}) + \langle G^k, x - x^{k+1} \rangle$

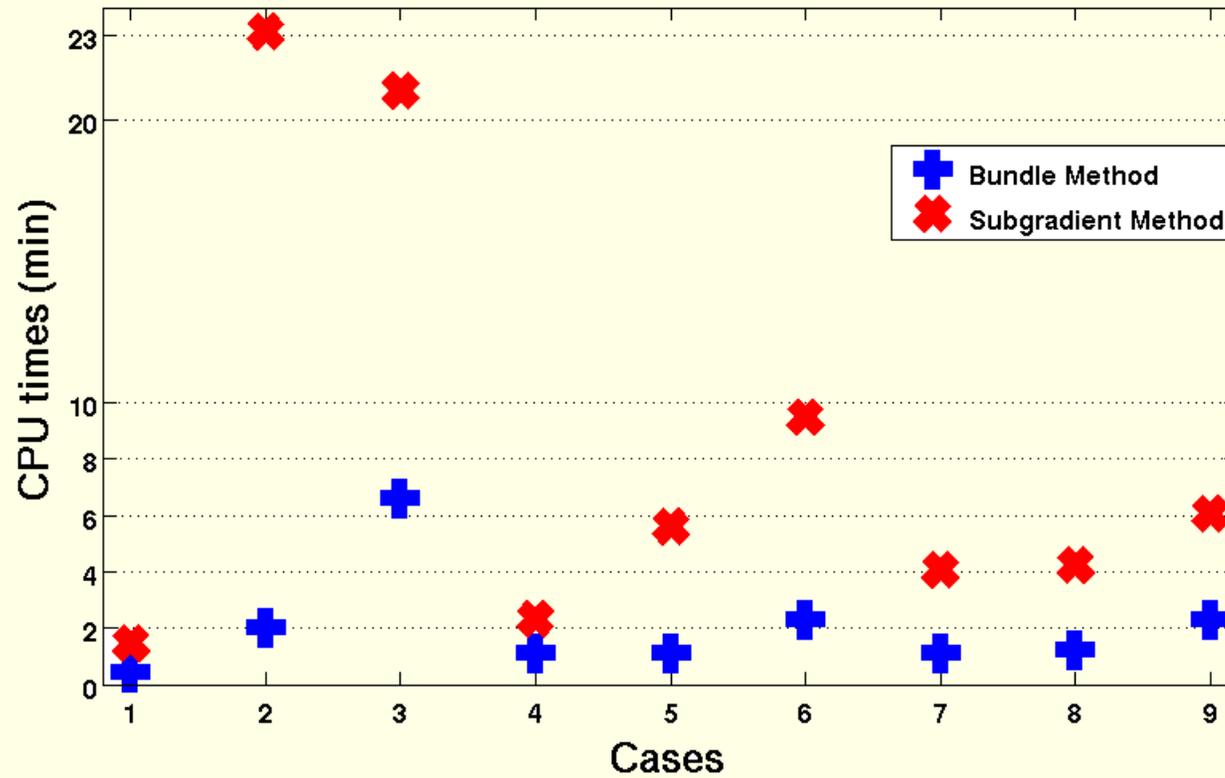
**Any model satisfying these conditions**

**that is used in the QP**

**maintains the convergence results**

# Comparing the methods: bundle and SG

Typical performance on a battery of Unit Commitment problems

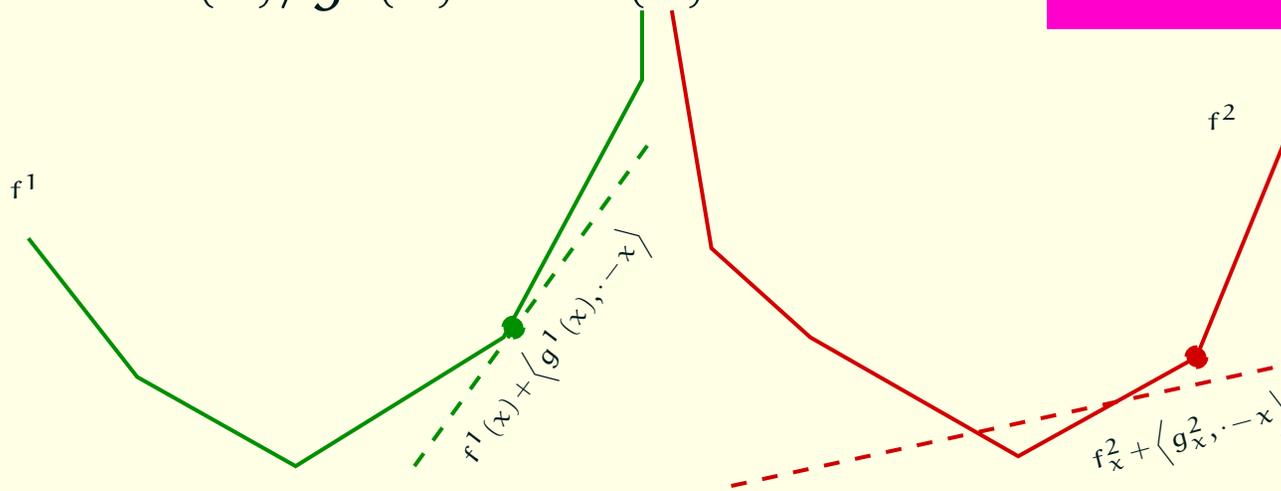


**Bundle Methods with on-demand accuracy**  
**the new generation**



# Oracle types: exact and upper

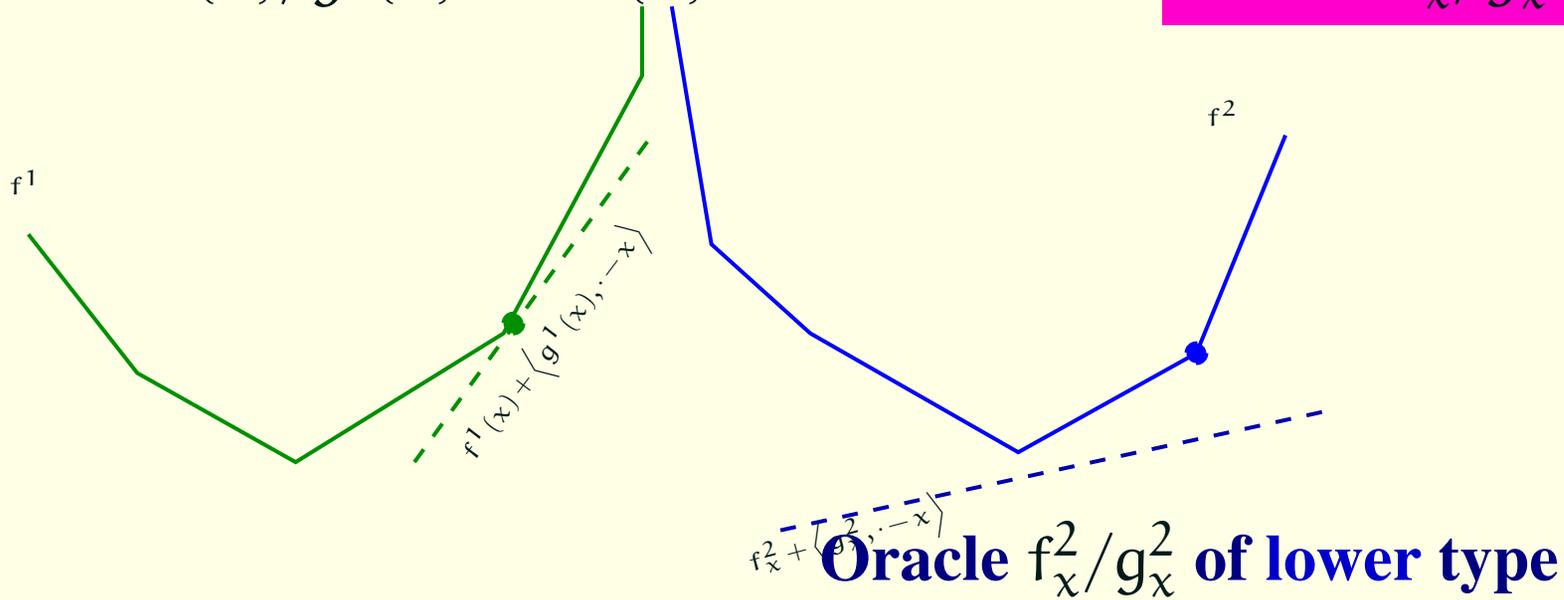
- $f^1(x)/g^1(x) \in \partial f^1(x)$  is **easy**: **exact**  $f^1(x)/g^1(x)$
- $f^2(x)/g^2(x) \in \partial f^2(x)$  is **difficult**: **inexact**  $f_x^2/g_x^2$



**Oracle  $f_x^2/g_x^2$  NOT of lower type**

# Oracle types: exact and lower

- $f^1(x)/g^1(x) \in \partial f^1(x)$  is **easy**: **exact**  $f^1(x)/g^1(x)$
- $f^2(x)/g^2(x) \in \partial f^2(x)$  is **less difficult**: **inexact**  $f^2_x/g^2_x$



**For the EM problem**  $f^j(\mathbf{x}) = \max\{-\mathcal{C}^j(q^j) + \langle \mathbf{x}, \mathbf{g}^j(q^j) \rangle : q^j \in \mathcal{P}^j\}$

By computing  $f_{\mathbf{x}^k}$  and  $\mathbf{g}_{\mathbf{x}^k}$  satisfying

$$f_{\mathbf{x}^k} = f(\mathbf{x}^k) - \eta^k \quad \text{and} \quad \mathbf{g}_{\mathbf{x}^k} \in \partial_{\eta^k} f(\mathbf{x}^k)$$

we can build

- A lower oracle
- An asymptotically exact oracle

$$\eta^k \rightarrow 0 \quad \text{as} \quad k \rightarrow \infty$$

- A partly asymptotically exact oracle

$$\eta^k \rightarrow 0 \quad \text{as} \quad K_s \ni k \rightarrow \infty$$

- An on-demand accuracy oracle

$$\eta^k \leq \bar{\eta}^k \quad \text{when} \quad f_{\mathbf{x}^k} \leq f_{\hat{\mathbf{x}}^k} - m\delta_k$$

## BM with lower inexact oracles

- $\mathbf{M}_k(\mathbf{x}) = \max\{f_{\mathbf{x}^i} + \langle \mathbf{g}_{\mathbf{x}^i}, \mathbf{x} - \mathbf{x}^i \rangle : i \in \mathcal{B}_k\}$
- $\delta^k = \varepsilon_k + t_k |\mathbf{G}^k|^2$
- SS test:  $f_{\mathbf{x}^{k+1}} \leq \hat{f}^k - m\delta^k$
- $\hat{f}^k := \max \left\{ f_{\hat{\mathbf{x}}^k}, \max \left( \mathbf{M}_j(\hat{\mathbf{x}}^k), j \geq \hat{k} \right) \right\}$

+ Oracle inaccuracy is locally bounded:

$\forall R \geq 0 \exists \eta(R) \geq 0 : |\mathbf{x}| \leq R \implies \eta \leq \eta(R)$  convergence as before, up to the accuracy on SS

## BM with lower inexact oracles

- $\mathbf{M}_k(\mathbf{x}) = \max\{f_{\mathbf{x}^i} + \langle \mathbf{g}_{\mathbf{x}^i}, \mathbf{x} - \mathbf{x}^i \rangle : i \in \mathcal{B}_k\}$
- $\delta^k = \varepsilon_k + t_k |\mathbf{G}^k|^2$
- SS test:  $f_{\mathbf{x}^{k+1}} \leq \hat{f}^k - m\delta^k$
- $\hat{f}^k := \max \left\{ f_{\hat{\mathbf{x}}^k}, \max \left( \mathbf{M}_j(\hat{\mathbf{x}}^k), j \geq \hat{k} \right) \right\}$

+ Oracle inaccuracy is locally bounded:

$\forall R \geq 0 \exists \eta(R) \geq 0 : |\mathbf{x}| \leq R \implies \eta \leq \eta(R)$  convergence as before, up to the accuracy on SS Convex proximal bundle methods in depth: a unified analysis for inexact oracles W. de Oliveira, C. Sagastizábal, C. Lemaréchal MathProg 148, pp 241-277, 2014

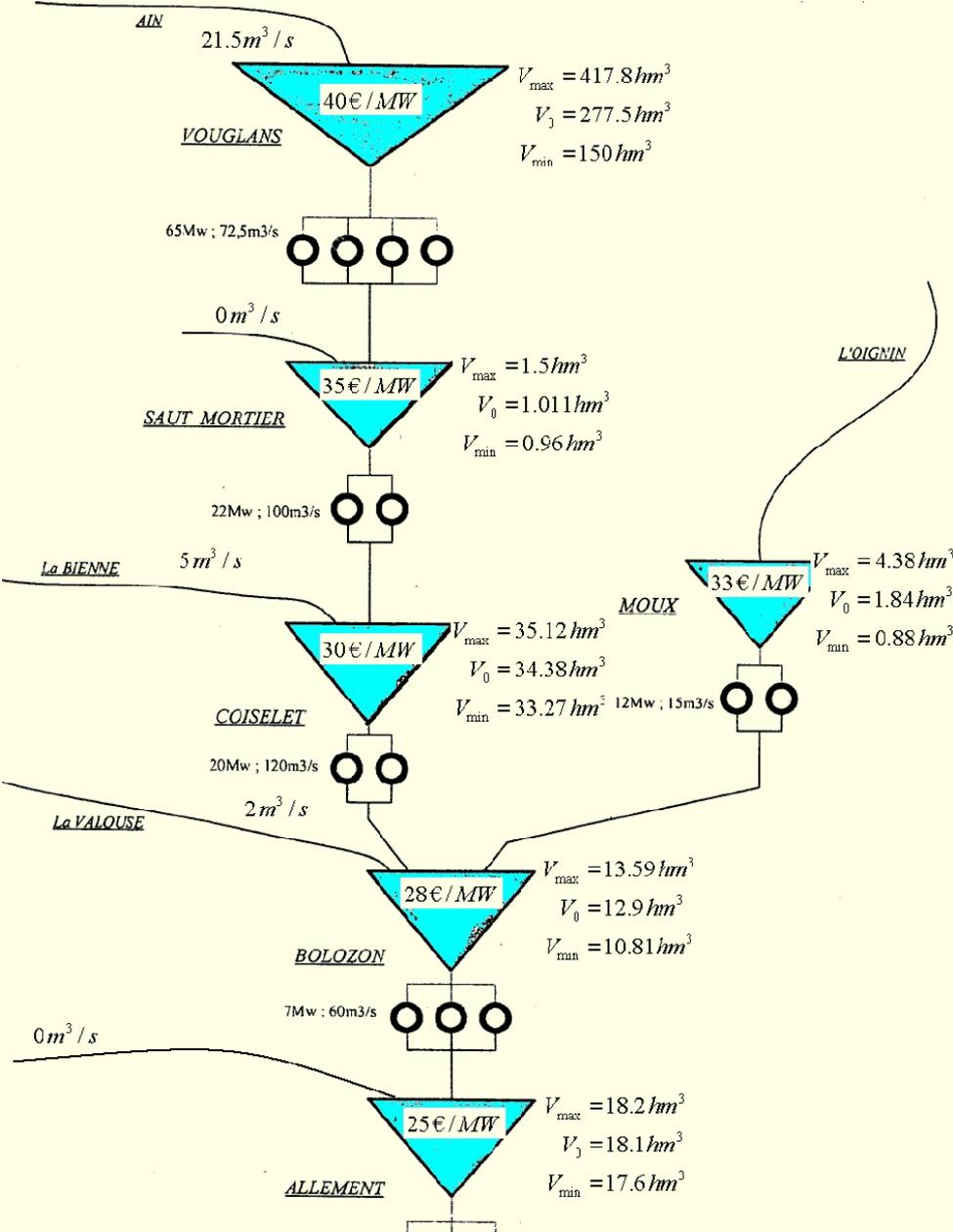
## General comments

Bundle methods are

- robust (do not oscillate, as CP methods do)
- reliable (have a stopping test, unlike SG methods)
- can deal with inaccuracy in a reasonable manner

# **Extending bundle methods**

# Constrained NSO problems: an example



# Optimal management of the hydrovalley

$$\left\{ \begin{array}{l} \max \quad \lambda^\top \mathbb{E}_\eta (\rho(\mathbf{u})) \\ \text{s.t.} \quad (\mathbf{x}, \mathbf{u}) \in \mathcal{P} \\ \mathbb{P}_\eta (A\mathbf{u} + \mathbf{a}_{\min} \leq M\eta \leq A\mathbf{u} + \mathbf{a}_{\max}) \geq p \end{array} \right.$$

# Optimal management of the hydrovalley

$$\left\{ \begin{array}{l} \max \quad \lambda^\top \mathbb{E}_\eta (\rho(\mathbf{u})) \\ \text{s.t.} \quad (\mathbf{x}, \mathbf{u}) \in \mathcal{P} \\ \mathbb{P}_\eta (A\mathbf{u} + \mathbf{a}_{\min} \leq M\eta \leq A\mathbf{u} + \mathbf{a}_{\max}) \geq p \end{array} \right.$$

**Is this a convex program?**

# Optimal management of the hydrovalley

$$\begin{cases} \max & \lambda^\top \mathbb{E}_\eta (\rho(\mathbf{u})) \\ \text{s.t.} & (\mathbf{x}, \mathbf{u}) \in \mathcal{P} \\ & \mathbb{P}_\eta (A\mathbf{u} + \mathbf{a}_{\min} \leq M\eta \leq A\mathbf{u} + \mathbf{a}_{\max}) \geq p \end{cases}$$

**Is this a convex program? YES:** the function

$$\mathbf{u} \mapsto \log \left( \mathbb{P}_\eta (A\mathbf{u} + \mathbf{a}_{\min} \leq M\eta \leq A\mathbf{u} + \mathbf{a}_{\max}) \right) \quad \text{is convex.}$$

We need to solve  $\begin{cases} \min & f(\mathbf{u}) \\ \text{s.t.} & (\mathbf{x}, \mathbf{u}) \in \mathcal{P} \quad \text{for linear } f \text{ and with} \\ & c(\mathbf{u}) \leq 0 \end{cases}$

$$c(\mathbf{u}) := \log \left( \mathbb{P}_\eta (A\mathbf{u} + \mathbf{a}_{\min} \leq M\eta \leq A\mathbf{u} + \mathbf{a}_{\max}) \right) - \log p$$

# Optimal management of the hydrovalley

$$\begin{cases} \max & \lambda^\top \mathbb{E}_\eta(\rho(\mathbf{u})) \\ \text{s.t.} & (\mathbf{x}, \mathbf{u}) \in \mathcal{P} \\ & \mathbb{P}_\eta(\mathbf{A}\mathbf{u} + \mathbf{a}_{\min} \leq \mathbf{M}\eta \leq \mathbf{A}\mathbf{u} + \mathbf{a}_{\max}) \geq p \end{cases}$$

Is this a convex program? YES: the function

$$\mathbf{u} \mapsto \log \left( \mathbb{P}_\eta(\mathbf{A}\mathbf{u} + \mathbf{a}_{\min} \leq \mathbf{M}\eta \leq \mathbf{A}\mathbf{u} + \mathbf{a}_{\max}) \right) \quad \text{is convex.}$$

We need to solve  $\begin{cases} \min & f(\mathbf{u}) \\ \text{s.t.} & (\mathbf{x}, \mathbf{u}) \in \mathcal{P} \quad \text{for linear } f \text{ and with} \\ & c(\mathbf{u}) \leq 0 \end{cases}$

$$c(\mathbf{u}) := \log \left( \mathbb{P}_\eta(\mathbf{A}\mathbf{u} + \mathbf{a}_{\min} \leq \mathbf{M}\eta \leq \mathbf{A}\mathbf{u} + \mathbf{a}_{\max}) \right) - \log p$$

**difficult to compute!**

## Need to solve the constrained problem

$$(P) \quad \begin{cases} \min & f(\mathbf{u}) \\ \text{s.t.} & (\mathbf{x}, \mathbf{u}) \in \mathcal{P} \\ & \mathbf{c}(\mathbf{u}) \leq 0 \end{cases}$$

for linear  $f$  and with inexact evaluation of  $\mathbf{c}$  and its gradient, via a black box with controllable inaccuracy (bounded by a given tolerance  $\varepsilon$ , with confidence level 99%, noting that evaluation errors can be positive or negative)

# Handling constraints in NSO

## For nonsmooth constrained problems

$$\min f(\mathbf{u}) \quad \text{s.t.} \quad c(\mathbf{u}) \leq 0$$

use the Improvement Function

$$\max_{\mathbf{u}} \{f(\mathbf{u}) - f(\hat{\mathbf{u}}), c(\mathbf{u})\}$$

(changes with each serious point  $\hat{\mathbf{u}}$  and supposes exact  $f/c$  values available)

[SagSol SiOPT, 2005 and

KarasRibSagSol MPB, 2009]

# Improvement function

Let  $(\bar{x}, \bar{u})$  be a solution to (P). The function

$$H_{\bar{u}}(\mathbf{u}) := \max_{(\mathbf{x}, \mathbf{u}) \in \mathcal{P}} \{f(\mathbf{u}) - f(\bar{u}), c(\mathbf{u})\}$$

has **perfect theoretical properties**:

If Slater condition  $(\exists (\mathbf{x}, \mathbf{u}) \in \mathcal{P} \text{ s.t. } c(\mathbf{u}) < 0)$  holds, then

$$\bar{u} \text{ solves } \min_{(\mathbf{x}, \mathbf{u}) \in \mathcal{P}} f(\mathbf{u}) \quad \text{s.t.} \quad c(\mathbf{u}) \leq 0 \quad \mathbf{(P)}$$



$$\min_{(\mathbf{x}, \mathbf{u}) \in \mathcal{P}} H_{\bar{u}}(\mathbf{u}) = H_{\bar{u}}(\bar{u}) = 0$$



$$0 \in \partial H(\bar{u}) \text{ for } H(\cdot) := H_{\bar{u}}(\cdot)$$

# Improvement function

Let  $(\bar{x}, \bar{u})$  be a solution to (P). The function

$$H_{\bar{u}}(\mathbf{u}) := \max_{(\mathbf{x}, \mathbf{u}) \in \mathcal{P}} \{f(\mathbf{u}) - f(\bar{u}), c(\mathbf{u})\}$$

has perfect theoretical properties:

## Without Slater condition

$$\bar{u} \text{ solves } \min_{(\mathbf{x}, \mathbf{u}) \in \mathcal{P}} f(\mathbf{u}) \quad \text{s.t.} \quad c(\mathbf{u}) \leq 0 \quad (\mathbf{P})$$

↓ **BUT** ↗

$$\min_{(\mathbf{x}, \mathbf{u}) \in \mathcal{P}} H_{\bar{u}}(\mathbf{u}) = H_{\bar{u}}(\bar{u}) = 0$$

↓ and also ↗

$$0 \in \partial H(\bar{u}) \text{ for } H(\cdot) := H_{\bar{u}}(\cdot)$$

# Improvement function

Let  $(\bar{x}, \bar{u})$  be a solution to (P). The function

$$H_{\bar{u}}(\mathbf{u}) := \max_{(\mathbf{x}, \mathbf{u}) \in \mathcal{P}} \{f(\mathbf{u}) - f(\bar{u}), c(\mathbf{u})\}$$

has perfect theoretical properties:

## Without Slater condition

$$\bar{u} \text{ solves } \min_{(\mathbf{x}, \mathbf{u}) \in \mathcal{P}} f(\mathbf{u}) \quad \text{s.t.} \quad c(\mathbf{u}) \leq 0 \quad (\mathbf{P})$$

↑: when  $c(\bar{u}) \leq 0$   $\bar{u}$  solves (P), otherwise it minimizes infeasibility over  $\mathcal{P}$

$$\min_{(\mathbf{x}, \mathbf{u}) \in \mathcal{P}} H_{\bar{u}}(\mathbf{u}) = H_{\bar{u}}(\bar{u}) = 0$$



$$0 \in \partial H(\bar{u}) \text{ for } H(\cdot) := H_{\bar{u}}(\cdot)$$

# **Handling nonconvex problems**

- Nonconvex proximal point mapping [PR96]

$$p_R f(x) := \operatorname{argmin}_{y \in \mathbb{R}^N} \left\{ f(y) + \frac{R}{2} |y - x|^2 \right\}$$

$x$  is the prox-center and  $R > R_x$  is the prox-parameter

**Theorem** If  $f$  is convex

- $p_R f$  is well defined **for any**  $R > 0$ .
- $p_R f$  is single valued and loc. Lip.
- $p = p_R f(x) \iff R(x - p) \in \partial f(p)$
- $x^*$  minimizes  $f \iff x^* = p_R f(x^*)$  for any  $R > 0$ .
- $x_{k+1} = p_R f(x_k)$  converges to a minimizer  $x^*$ .

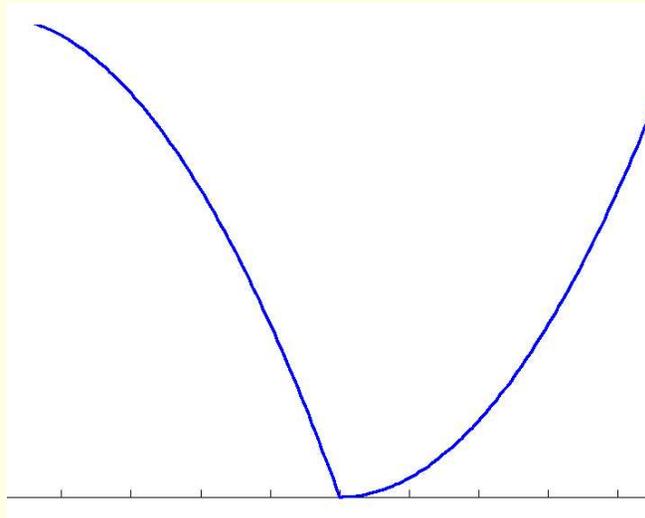
**What is  $f$  is nonconvex?**

# Nonconvex difficulties

**Proximal Bundle Methods** are the most robust and reliable (oracle) methods for convex minimization. Their success relies heavily on convexity. If  $f$  is convex:

- $x_{k+1} = p_R f(x_k)$  converges to a minimizer  $x^*$ .
- $\check{f}_k$  lies **entirely** below  $f$ .

**May no longer be true for nonconvex  $f$**

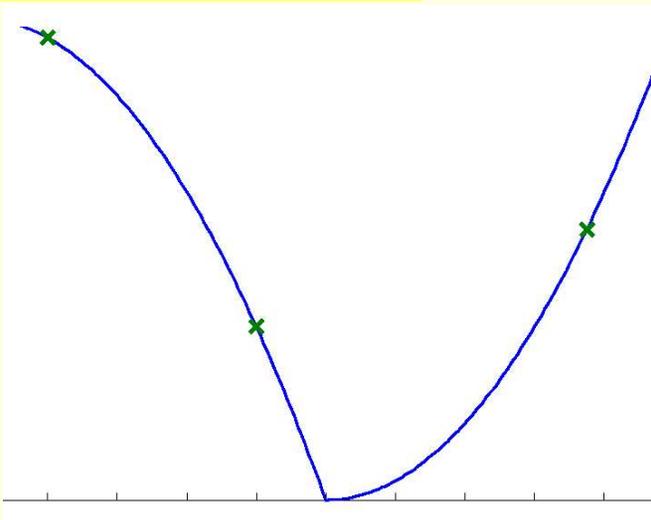


# Nonconvex difficulties

**Proximal Bundle Methods** are the most robust and reliable (oracle) methods for convex minimization. Their success relies heavily on convexity. If  $f$  is convex

- $x_{k+1} = p_R f(x_k)$  converges to a minimizer  $x^*$ .
- $\check{f}_k$  lies **entirely** below  $f$ .

**May no longer be true for nonconvex  $f$**

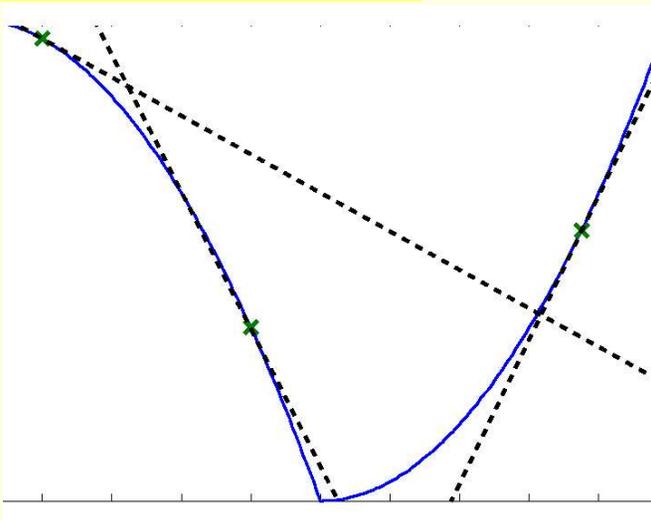


# Nonconvex difficulties

**Proximal Bundle Methods** are the most robust and reliable (oracle) methods for convex minimization. Their success relies heavily on convexity. If  $f$  is convex

- $x_{k+1} = p_R f(x_k)$  converges to a minimizer  $x^*$ .
- $\check{f}_k$  lies **entirely** below  $f$ .

**May no longer be true for nonconvex  $f$**

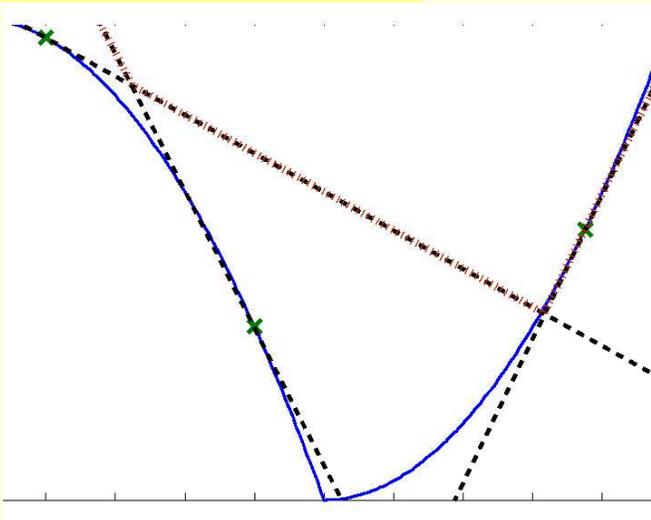


# Nonconvex difficulties

**Proximal Bundle Methods** are the most robust and reliable (oracle) methods for convex minimization. Their success relies heavily on convexity. If  $f$  is convex

- $x_{k+1} = p_R f(x_k)$  converges to a minimizer  $x^*$ .
- $\check{f}_k$  lies **entirely** below  $f$ .

**May no longer be true for nonconvex  $f$**

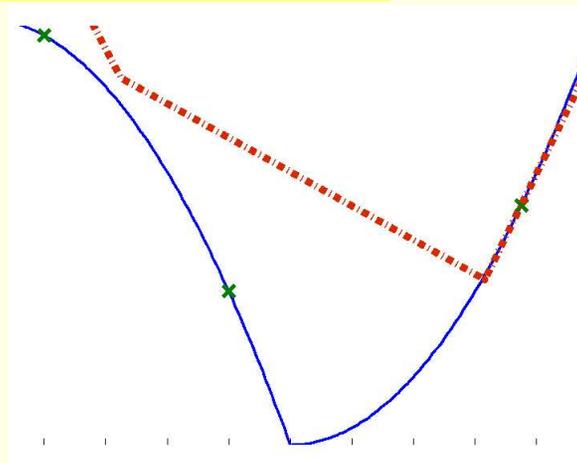


# Nonconvex difficulties

**Proximal Bundle Methods** are the most robust and reliable (oracle) methods for convex minimization. Their success relies heavily on convexity. If  $f$  is convex

- $x_{k+1} = p_R f(x_k)$  converges to a minimizer  $x^*$ .
- $\check{f}_k$  lies **entirely** below  $f$ .

**May no longer be true for nonconvex  $f$**



**How this difficulty has been addressed?**

Take each plane in the model:  $f_i + \langle g_i, \cdot - y_i \rangle$  and rewrite it, centered at  $x_k$ :

$$f(x_k) - \left[ f(x_k) - \left( f_i + \langle g_i, x_k - y_i \rangle \right) \right] + \langle g_i, \cdot - x_k \rangle$$

$$f(x_k) - e_{k,i}^f + \langle g_i, \cdot - x_k \rangle$$

$$\Rightarrow \check{f}_k(y) = \max \left\{ f(x_k) - e_{k,i}^f + \langle g_i, y - x_k \rangle \right\}$$

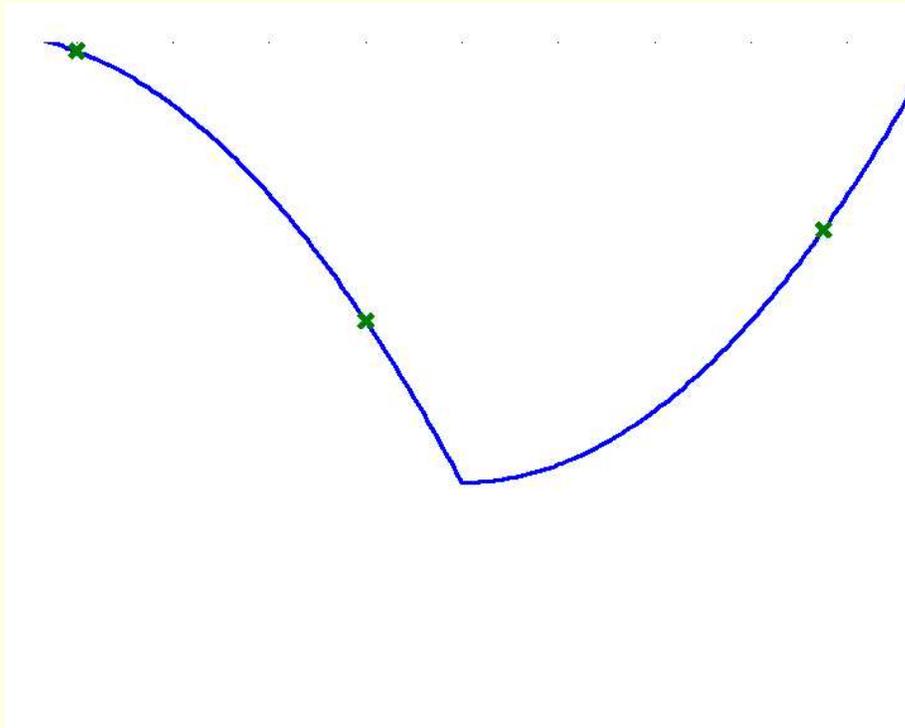
**Good:**  $e_{k,i}^f$  positive  $\Rightarrow$  convergence **Good:** If  $f$  convex  $\Rightarrow e_{k,i}^f$  positive. **BAD:** If  $f$  nonconvex,  $e_{k,i}^f$  may be **negative**

## Nonconvex bundle methods

**fix negative linearization errors**, replacing  $\check{f}_k$  by:

$$\check{f}_k^{\mathbf{FIX}}(y) = \max \left\{ f(x_k) - |e_{k,i}^f| + \langle g_i, y - x_k \rangle \right\}$$

[Mif77, Lem80, Kiw85, Luk98]

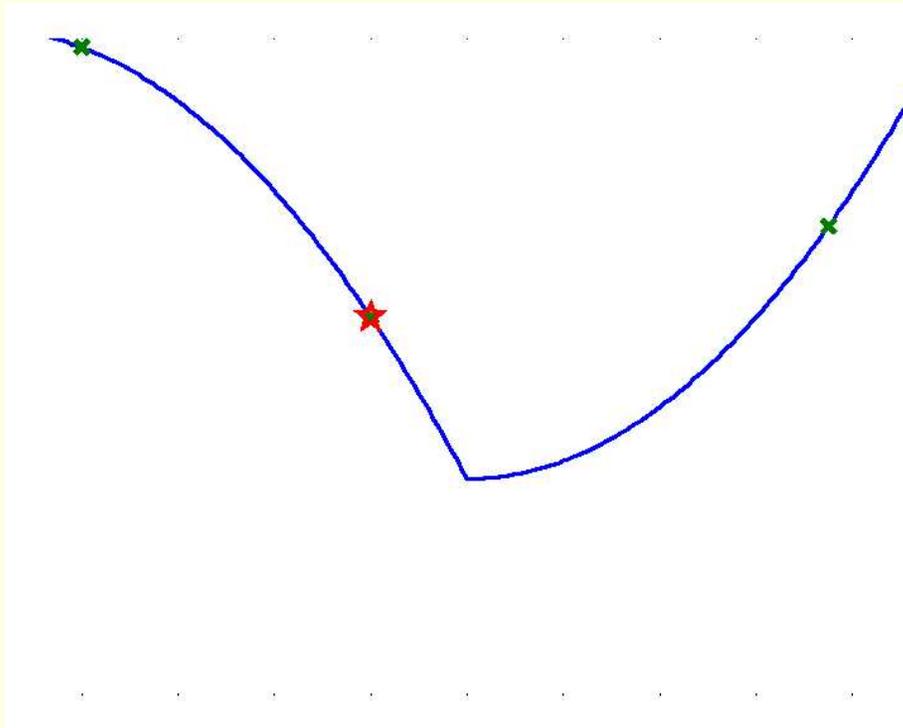


## Nonconvex bundle methods

**fix negative linearization errors**, replacing  $\check{f}_k$  by:

$$\check{f}_k^{\mathbf{FIX}}(y) = \max \left\{ f(x_k) - |e_{k,i}^f| + \langle g_i, y - x_k \rangle \right\}$$

[Mif77, Lem80, Kiw85, Luk98]

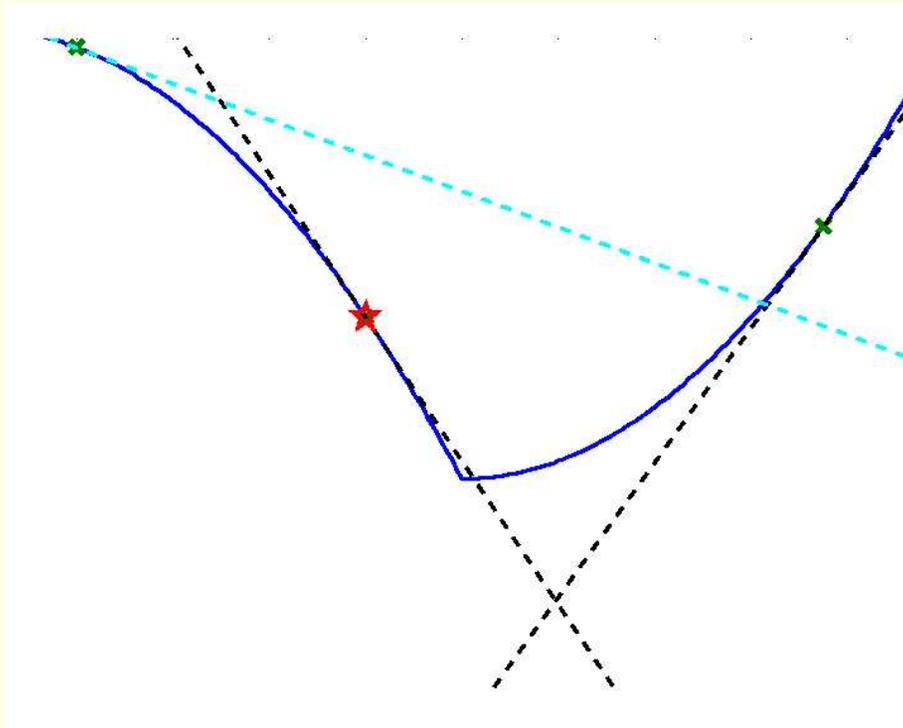


## Nonconvex bundle methods

**fix negative linearization errors**, replacing  $\check{f}_k$  by:

$$\check{f}_k^{\mathbf{FIX}}(y) = \max \left\{ f(x_k) - |e_{k,i}^f| + \langle g_i, y - x_k \rangle \right\}$$

[Mif77, Lem80, Kiw85, Luk98]

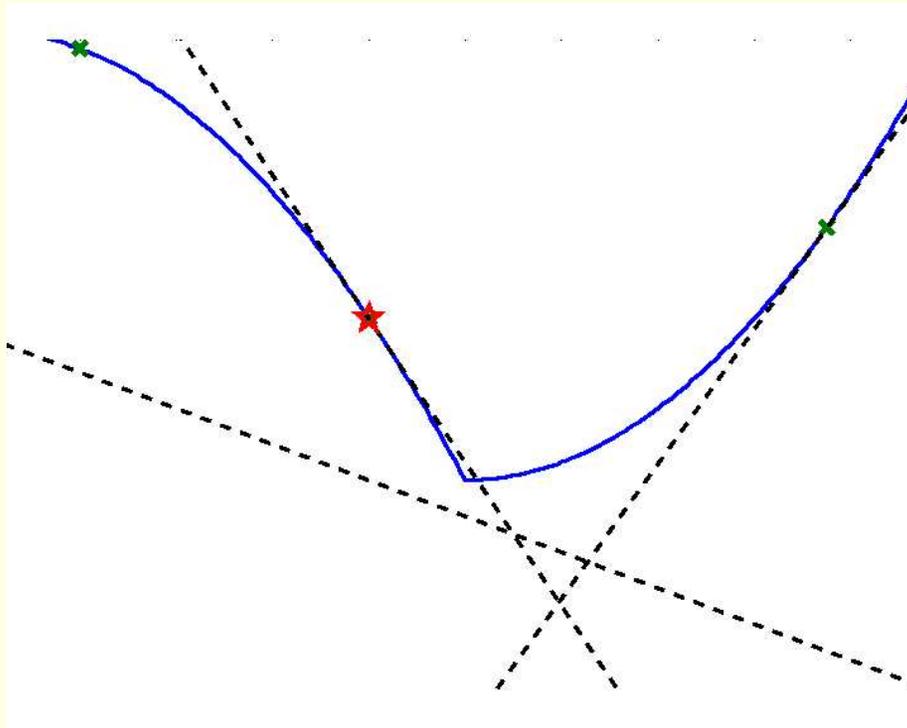


## Nonconvex bundle methods

**fix negative linearization errors**, replacing  $\check{f}_k$  by:

$$\check{f}_k^{\mathbf{FIX}}(y) = \max \left\{ f(x_k) - |e_{k,i}^f| + \langle g_i, y - x_k \rangle \right\}$$

[Mif77, Lem80, Kiw85, Luk98]

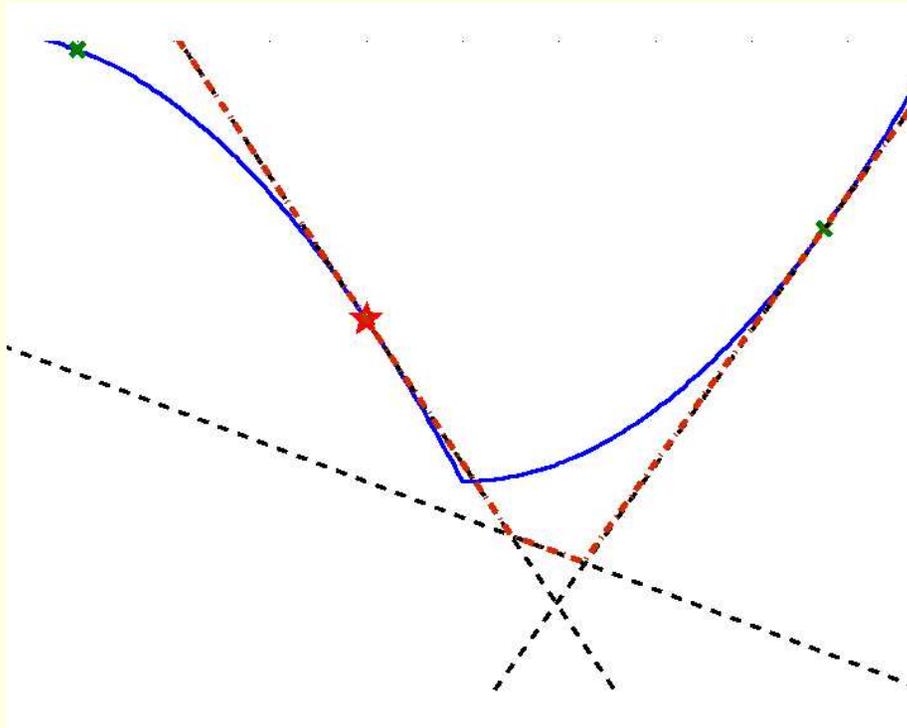


## Nonconvex bundle methods

**fix negative linearization errors**, replacing  $\check{f}_k$  by:

$$\check{f}_k^{\mathbf{FIX}}(y) = \max \left\{ f(x_k) - |e_{k,i}^f| + \langle g_i, y - x_k \rangle \right\}$$

[Mif77, Lem80, Kiw85, Luk98]

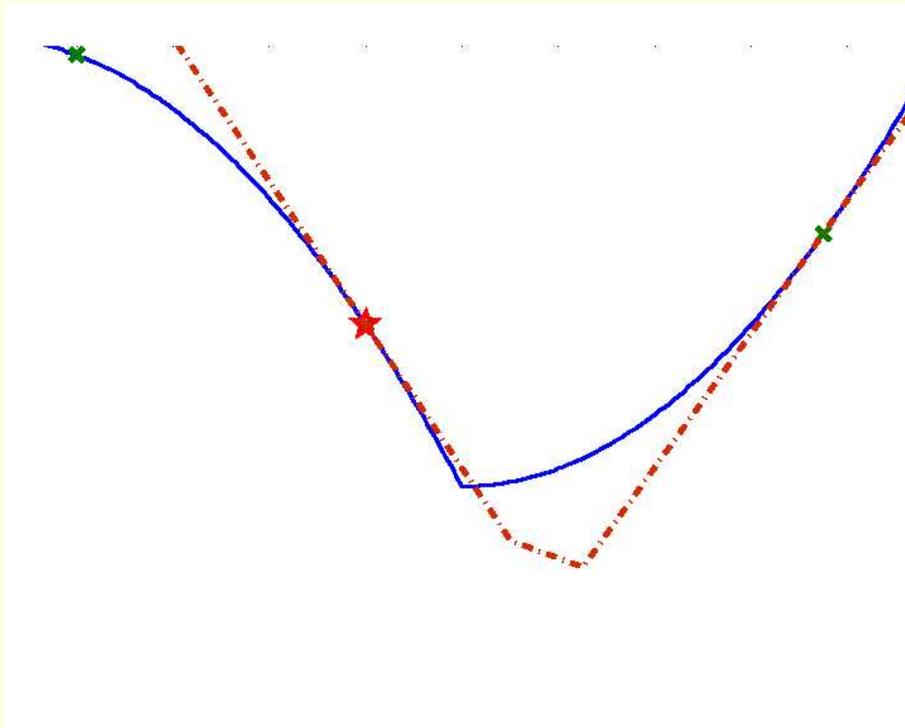


## Nonconvex bundle methods

**fix negative linearization errors**, replacing  $\check{f}_k$  by:

$$\check{f}_k^{\mathbf{FIX}}(y) = \max \left\{ f(x_k) - |e_{k,i}^f| + \langle g_i, y - x_k \rangle \right\}$$

[Mif77, Lem80, Kiw85, Luk98]



# A new method

A different approach (ours) is based on the following trick

Take  $\eta, \mu > 0 : R = \eta + \mu$  and note

$$\begin{aligned} p_R f(x_k) &= \min_w \left\{ f(w) + R \frac{1}{2} |w - x_k|^2 \right\} \\ &= \min_w \left\{ f(w) + (\eta + \mu) \frac{1}{2} |w - x_k|^2 \right\} \\ &= \min_w \left\{ f(w) + \eta \frac{1}{2} |w - x_k|^2 + \mu \frac{1}{2} |w - x_k|^2 \right\} \\ &= \min_w \left\{ F_\eta(w) + \mu \frac{1}{2} |w - x_k|^2 \right\} \\ &= p_\mu(F_\eta)(x_k) \end{aligned}$$

$$\Rightarrow p_R f = p_\mu(F_\eta)$$

# Redistributed Proximal Bundle Method

At  $\ell^{\text{th}}$ -iteration, for  $k = k(\ell)$ , given  $R_k$ ,  $x_k$  and a bundle

$$\mathcal{B} = \{y_i, f_i, g_i, i \in I_\ell\}$$

0. Split  $R_k$  into  $\eta_\ell$  and  $\mu_\ell$ .

1. Model  $F_{\eta_\ell}$   $\check{F}_{\eta_\ell, \ell}(y) = \max_{i \in \mathcal{B}} \{F_{\eta_\ell i} + \langle g_{\eta_\ell i}, y - y_i \rangle\}$

2. Minimize **the penalized model**

$$y_{\ell+1} = \arg \min \{ \check{F}_{\eta_\ell, \ell}(y) + \frac{\mu_\ell}{2} |y - x_k|^2 \}$$

3. Descent test If  $y_{\ell+1}$  good:  $x_{k+1} \leftarrow y_{\ell+1}$ , define  $R_{k+1}$

**serious step**

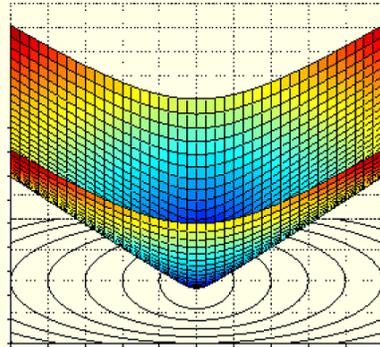
If  $y_{\ell+1}$  bad:

**null step**

4. Update bundle  $\mathcal{B} \leftarrow \mathcal{B} \cup \{y_{\ell+1}, f_{\ell+1}, g_{\ell+1}\}$

# $\mathcal{V}\mathcal{U}$ quasi-Newton bundle

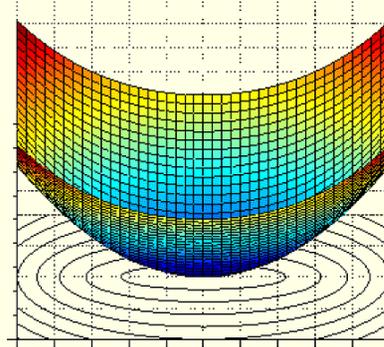
For  $\mathbf{x} \in \mathbb{R}^n$ , given matrices  $A \succeq 0$ ,  $B \succ 0$ ,  $f(\mathbf{x}) = \sqrt{\mathbf{x}^\top A \mathbf{x}} + \mathbf{x}^\top B \mathbf{x}$  has a unique minimizer at  $\bar{\mathbf{x}} = 0$ . On  $\mathcal{N}(A)$  the function is not differentiable, and the first term vanishes:  $f|_{\mathcal{N}(A)}$  looks smooth.



$\mathcal{R}(A)$

$\mathcal{V}$

parallel to  $\partial f(\bar{\mathbf{x}})$



$\mathcal{N}(A)$

$\mathcal{U}$

$\mathcal{U}$  perpendicular to  $\mathcal{V}$

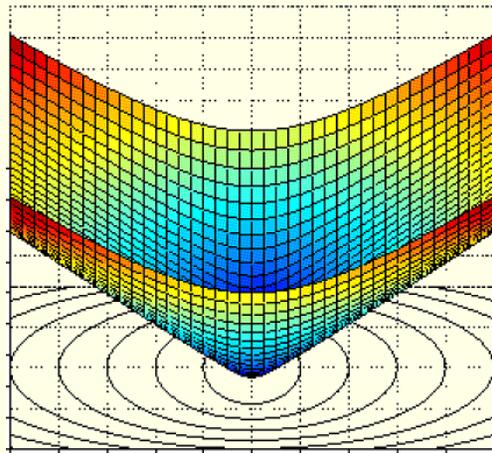
$\mathcal{V}$  is parallel to  $\mathcal{N}(A)$ , the “ridge” of nonsmoothness



# $\mathcal{V}\mathcal{U}$ -Algorithm:

(Mifflin&Sagastizábal, MathProg 05) Recall that

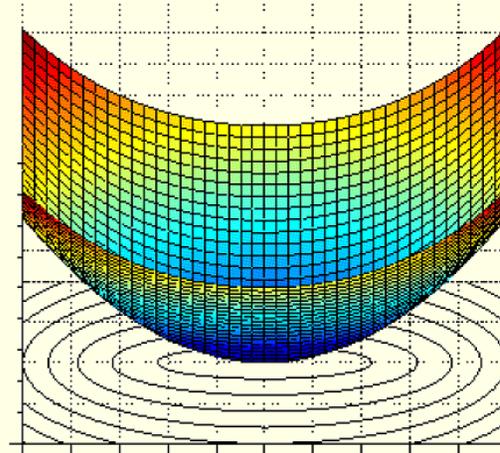
$f|_{\mathcal{V}||\mathcal{N}(A)}$  is nice: **the key is the two QP-solves**



$\mathcal{R}(A)$

$\mathcal{V}$

**2 bundle QPs**



$\mathcal{N}(A)$

$\mathcal{U}$

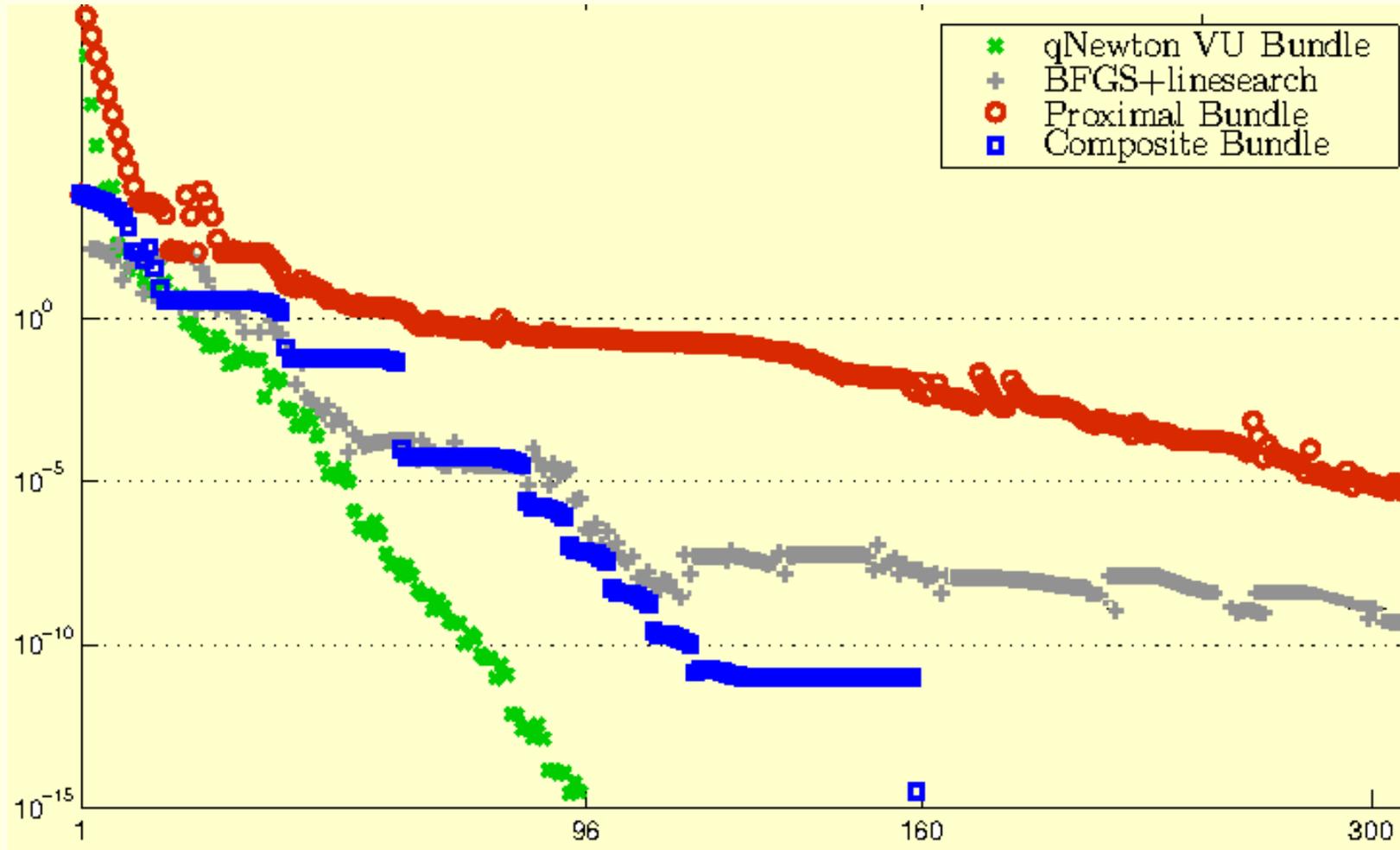
**Newton-move**

Two successive bundle QPs identify the “ridge” of nonsmoothness

Solve a **2nd QP** to create an approximation of  $\mathcal{V}$  based on  $\partial \check{f}(\hat{p}^k)$

# VU-Algorithm:

superlinear “serious” subsequence (Mifflin&Sag, MathProg 05)



## To learn more

### **Bundle methods history**

R. MIFFLIN, C. SAGASTIZÁBAL, Documenta Math, 2012. A Science Fiction Story in Nonsmooth Optimization Originating at IIASA

[https://www.math.uni-bielefeld.de/documenta/vol-ismmp/44\\_mifflin-robert.pdf](https://www.math.uni-bielefeld.de/documenta/vol-ismmp/44_mifflin-robert.pdf)

### **(exact) Bundle books**

J.F. BONNANS, J.C. GILBERT, C. LEMARÉCHAL, AND C. SAGASTIZÁBAL, Numerical Optimization: Theoretical and Practical Aspects, Springer, 2nd ed., 2006.

J.B. HIRIART-URRUTY AND C. LEMARÉCHAL, Convex Analysis and Minimization Algorithms II, no. 306 in Grund. der math. Wissenschaften, Springer, 2nd ed., 1996.

### **Inexact Bundle theory**

(next page)

### **Inexact Bundle theory**

M. HINTERMÜLLER, A proximal bundle method based on approximate subgradients, COAp 20 (2001), pp. 245–266.

M. V. SOLODOV, On Approximations with Finite Precision in Bundle Methods for Nonsmooth Optimization. JoTA 119.1 (2003), pp. 151–165

K.C. KIWIEL, A proximal bundle method with approximate subgradient linearizations, SiOpt 16 (2006), pp. 1007–1023.

W. DE OLIVEIRA, C. SAGASTIZÁBAL, AND C. LEMARÉCHAL, Convex proximal bundle methods in depth: a unified analysis for inexact oracles, MathProg 148 (2014), pp. 241–277.

### **Inexact Bundle variants with applications**

G. EMIEL AND C. SAGASTIZÁBAL, Incremental-like bundle methods with application to energy planning, COAp 46 (2010), pp. 305–332.

W. DE OLIVEIRA, C. SAGASTIZÁBAL, AND S. SCHEIMBERG, Inexact bundle methods for two-stage stochastic programming, SiOpt 21 (2011), pp. 517–544.

(next page)

W. VAN ACKOOIJ AND C. SAGASTIZÁBAL, Constrained bundle methods for upper inexact oracles with application to joint chance constrained energy problems, *SiOpt* 24 (2014), pp. 733–765.

W. DE OLIVEIRA AND C. SAGASTIZÁBAL, Level bundle methods for oracles with on-demand accuracy, *OMS* 29 (2014), pp. 1180 –1209

W. DE OLIVEIRA AND C. SAGASTIZÁBAL, Bundle methods in the xxi century: A birds'-eye view, *Pesquisa Operacional* 34 (2014), pp. 647 – 670.

W. DE OLIVEIRA AND M. SOLODOV, A doubly stabilized bundle method for nonsmooth convex optimization, *MathProg* 156(1), pp. 126–159, 2016.

Any doubts or questions?

Just e-mail me