



# BUNDLE METHODS FOR STOCHASTIC PROGRAMS

## PROXIMAL BUNDLE METHOD

Wellington de Oliveira

BAS Lecture 25, June 9, 2016, IMPA

## GENERAL FORMULATION

In this part of the course we will focus on efficient optimization methods to solve convex programs of the form

$$\min f(x) \quad \text{s.t.} \quad x \in X,$$

with

- ▶  $f : \mathfrak{R}^n \rightarrow \mathfrak{R}$  a convex but **nonsmooth function**
- ▶  $X \subset \mathfrak{R}^n$  a convex set (e.g.  $X = \{x \in \mathfrak{R}_+^n : Ax = b\}$ ,  $X = \mathfrak{R}^n$ )

This formulation covers many practical optimization problems, for instance

- ▶ Two-stage stochastic programming problems
- ▶ Multistage stochastic programming problems

## TWO-STAGE STOCHASTIC LINEAR PROGRAMMING

In two-stage stochastic linear programming problems with finitely many scenarios  $\xi^i = (q^i, T^i, W^i, h^i)$  we wish to solve the high dimensional LP

$$\left\{ \begin{array}{l} \min \quad c^\top x + \sum_{i=1}^N p_i [q^{i\top} y^i] \\ \text{s.t.} \quad Ax = b, x \geq 0 \\ \quad \quad T^i x + W^i y^i = h^i, y^i \geq 0, \quad i = 1, \dots, N \end{array} \right.$$

## TWO-STAGE STOCHASTIC LINEAR PROGRAMMING

In two-stage stochastic linear programming problems with finitely many scenarios  $\xi^i = (q^i, T^i, W^i, h^i)$  we wish to solve the high dimensional LP

$$\begin{cases} \min & c^\top x + \sum_{i=1}^N p_i [q^{i\top} y^i] \\ \text{s.t.} & Ax = b, x \geq 0 \\ & T^i x + W^i y^i = h^i, y^i \geq 0, \quad i = 1, \dots, N \end{cases}$$

## TWO-STAGE DECOMPOSITION

$$\min f(x) \quad \text{s.t.} \quad x \in X, \quad \text{with} \quad f(x) := c^\top x + \sum_{i=1}^N p_i Q(x, \xi^i),$$

$$Q(x, \xi) = \begin{cases} \min & q^\top y \\ \text{s.t.} & Wy = h - Tx \\ & y \geq 0. \end{cases} \quad \text{and} \quad X := \{x \in \mathbb{R}_+^n : Ax = b\}$$

We know that  $g = c - \sum_{i=1}^N p_i T^{i\top} \pi^i \in \partial f(x)$ , where  $\pi^i$  is a dual solution of  $Q(x, \xi^i)$

## MULTISTAGE STOCHASTIC LINEAR PROGRAMS

$$\min_{\substack{A_1 x_1 = b_1 \\ x_1 \geq 0}} c_1^\top x_1 + \mathbb{E} \left[ \min_{\substack{B_2 x_1 + A_2 x_2 = b_2 \\ x_2 \geq 0}} c_2^\top x_2 + \mathbb{E} \left[ \cdots + \mathbb{E} \left[ \min_{\substack{B_T x_{T-1} + A_T x_T = b_T \\ x_T \geq 0}} c_T^\top x_T \right] \right] \right]$$

- Some elements of the data  $\xi = (c_t, B_t, A_t, b_t)$  depend on uncertainties.

By assuming finitely many scenarios and dualizing the nonanticipativity constraints (that can be written as  $Gx = 0$ ) we get

# MULTISTAGE STOCHASTIC LINEAR PROGRAMS

(See Lecture 17)

## DUAL PROBLEM

$$\min_u f(u), \quad \text{with} \quad f(u) := - \sum_{i=1}^N D^i(u)$$

$$D^i(u) := \begin{cases} \min_{x^i} & p_i \sum_{t=1}^T (c_t^i)^\top x_t^i + u^\top G^i x^i \\ \text{s.t.} & A_1 x_1 = b_1 \\ & B_t^i x_{t-1}^i + A_t^i x_t^i = b_t^k, \quad t = 2, \dots, T \\ & x_t^i \geq 0. \end{cases}$$

Computing  $f(u)$  for each given  $u$  amounts to solving  $N$  LPs.

We know that  $g = -Gx(u) \in \partial f(u)$ , where  $x(u) = (x^1(u), \dots, x^N(u))$  and  $x^i(u)$  is a solution of  $D^i(u)$

Let's stick with the more compact and general formulation

$$\min f(x) \quad \text{s.t.} \quad x \in X,$$

with  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  a convex but **nonsmooth function** and  $X \subset \mathbb{R}^n$  a convex set.

We'll assume the availability of an oracle providing us with first-order information on  $f$ :

$$x \longrightarrow \boxed{\text{Oracle}} \longrightarrow \begin{cases} \text{function value} & f(x) \\ \text{subgradient} & g \in \partial f(x) \end{cases}$$

**In stochastic programming, the oracle should be smart enough to use parallel computing:**

- ▶ the oracle consists of solving  $N$  optimization subproblems to compute  $f(x)$  and a subgradient  $g$
- ▶ **most of time dedicate to minimize  $f$  is spent in the oracle!**

Therefore, subgradient and (pure) cutting-plane methods are not very efficient<sup>1</sup>...

---

<sup>1</sup>These methods require, in general, many oracle calls.

## CUTTING-PLANE METHOD

Consider the problem

$$\min_{x \in X} f(x)$$

and suppose that  $X$  is a compact set.

### ALGORITHM

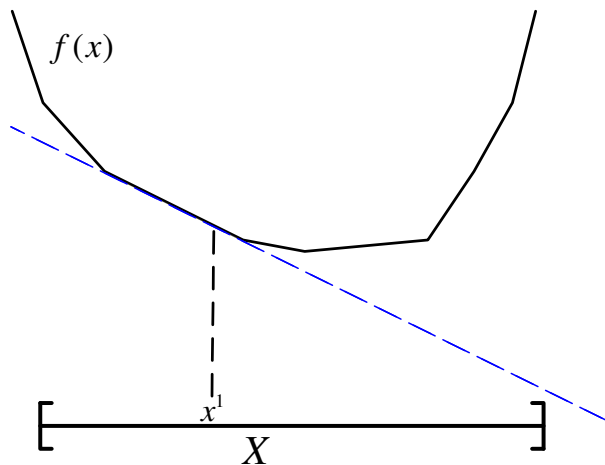
1. Given  $x_0 \in X$ , call the oracle to compute  $f(x_0)$  and  $g_0 \in \partial f(x_0)$ . Set  $f_0^{\text{up}} = f(x_0)$  and  $k = 0$
2. (iterate) Find  $x_{k+1} = \arg \min_{x \in X} \check{f}_k(x)$ . Let  $f_k^{\text{low}} = \check{f}_k(x_{k+1})$ .
3. (stopping test) If  $f_k^{\text{up}} - f_k^{\text{low}}$  is small enough, stop.
4. (oracle) Compute  $f(x_{k+1})$ ,  $g_{k+1} \in \partial f(x_{k+1})$  and set  $f_{k+1}^{\text{up}} = \min\{f(x_{k+1}), f_k^{\text{up}}\}$ .
5. (loop) Set  $k \leftarrow k + 1$  and go back to Step 2.

### CUTTING-PLANE MODEL

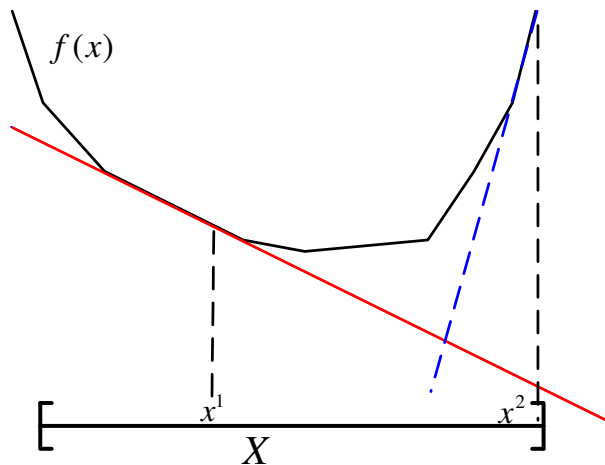
$$\check{f}_k(\cdot) = \max_{j=1, \dots, k} \{f(x_j) + g_j^\top(\cdot - x_j)\}$$



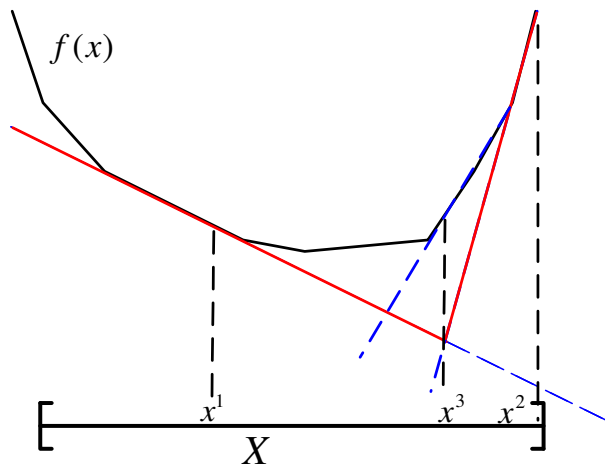
## CUTTING-PLANE METHOD



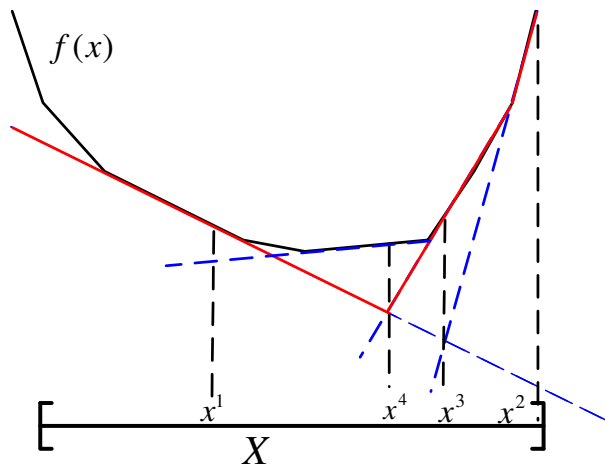
# CUTTING-PLANE METHOD



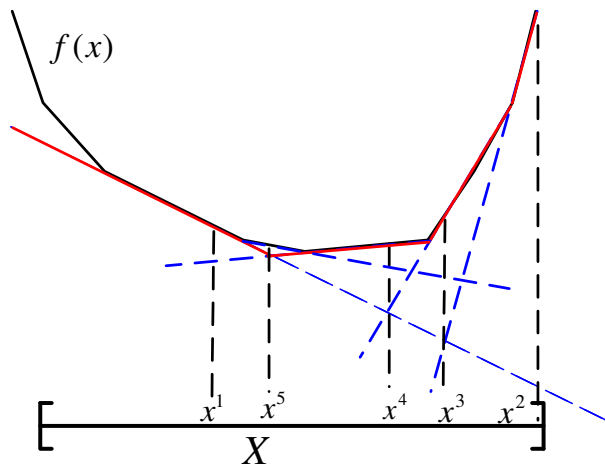
# CUTTING-PLANE METHOD



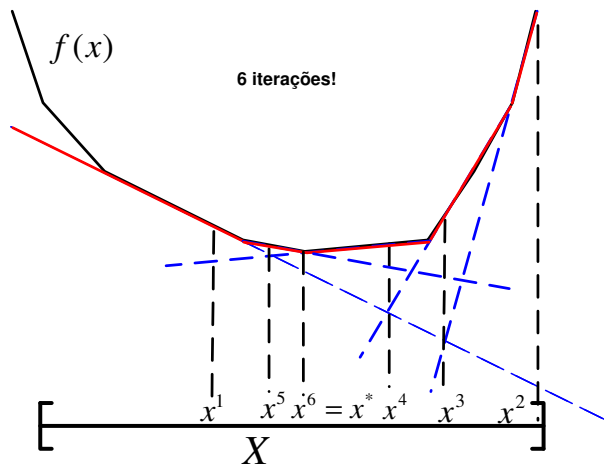
# CUTTING-PLANE METHOD



# CUTTING-PLANE METHOD



# CUTTING-PLANE METHOD



## CUTTING-PLANE METHOD

The method requires solving a LP at each iteration

$$x_{k+1} = \arg \min_{x \in X} \check{f}_k(x), \quad \check{f}_k(\cdot) = \max_{j=1, \dots, k} \{f(x_j) + g_j^\top(\cdot - x_j)\}$$

that is equivalent to

$$\begin{cases} \min_{x, r} & r \\ \text{s.t.} & f(x_j) + g_j^\top(x - x_j) \leq r, \quad j = 1, \dots, k \\ & x \in X, r \in \mathfrak{R}. \end{cases}$$

A new constraint is added at each iteration!

# CUTTING-PLANE METHOD

## PROS × CONS

- 👍 only computes a single subgradient per iteration
- 👍 easy to code
- 👍 easy and reliable stopping test
  
- 👎  $f(x_{k+1}) \not\leq f(x_k)$  (it is not a descent method)
- 👎 instable and has low convergence rate
- 👎 requires compactness of the feasible set
- 👎 doesn't exploit good starting points
- 👎 subproblem becomes heavier and heavier...

The *Regularized Decomposition Method* (1986) for 2-SLP address some of the above drawbacks.

Regularized Decomposition Method is just a particular case of (proximal) Bundle Methods!



# BUNDLE METHODS

## MAIN INGREDIENTS

- (I) a convex model  $f_k^M \leq f$  (eg. cutting-plane model)
- (II) a stability center  $\hat{x}_k$  (eg.: the best point so far)
- (III) a parameter  $t_k$  (or  $f_k^{\text{lev}}$ ) to be updated at every iteration

The next trial point  $x_{k+1}$  of a bundle method depends on the above 3 ingredients, whose organization define different methods:

**PROXIMAL BUNDLE METHOD** ( $t_k > 0$ )

$$x_{k+1} := \arg \min \left\{ f_k^M(x) + \frac{1}{2t_k} \|x - \hat{x}_k\|^2 : x \in X \right\}.$$

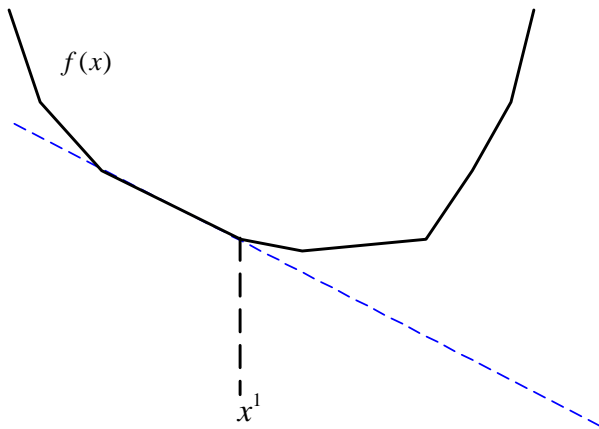
**LEVEL BUNDLE METHOD** ( $f_k^{\text{lev}} \in \mathfrak{R}$ )

$$x_{k+1} := \arg \min \left\{ \frac{1}{2} \|x - \hat{x}_k\|^2 : f_k^M(x) \leq f_k^{\text{lev}}, x \in X \right\}.$$

Today we focus on proximal bundle method!

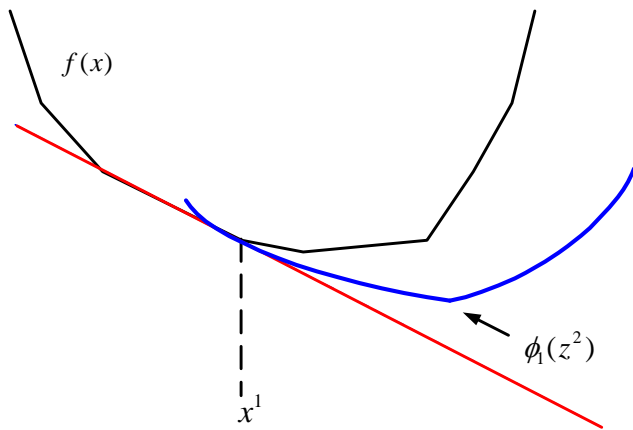
## PROXIMAL BUNDLE METHOD

$$f^M \equiv \tilde{f}, \quad x_{k+1} := \arg \min \left\{ \tilde{f}_k(x) + \frac{1}{2t_k} \|x - \hat{x}_k\|^2 : x \in X \right\}$$



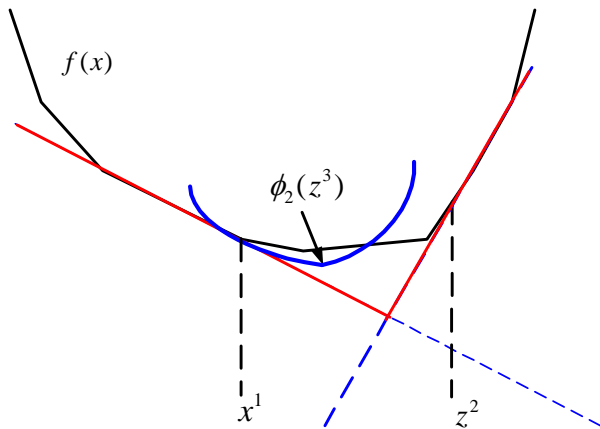
## PROXIMAL BUNDLE METHOD

$$f^M \equiv \tilde{f}, \quad x_{k+1} := \arg \min \left\{ \tilde{f}_k(x) + \frac{1}{2t_k} \|x - \hat{x}_k\|^2 : x \in X \right\}$$



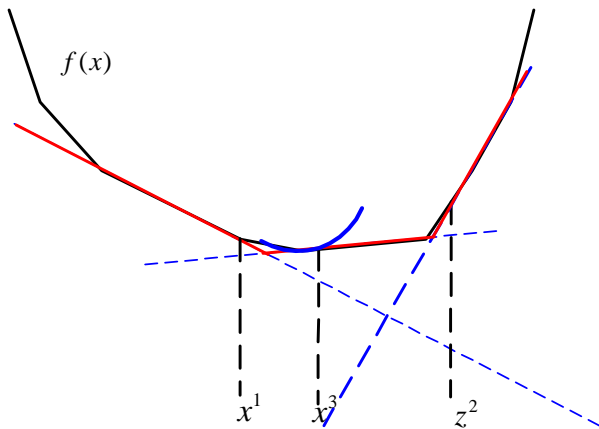
## PROXIMAL BUNDLE METHOD

$$f^M \equiv \tilde{f}, \quad x_{k+1} := \arg \min \left\{ \tilde{f}_k(x) + \frac{1}{2t_k} \|x - \hat{x}_k\|^2 : x \in X \right\}$$



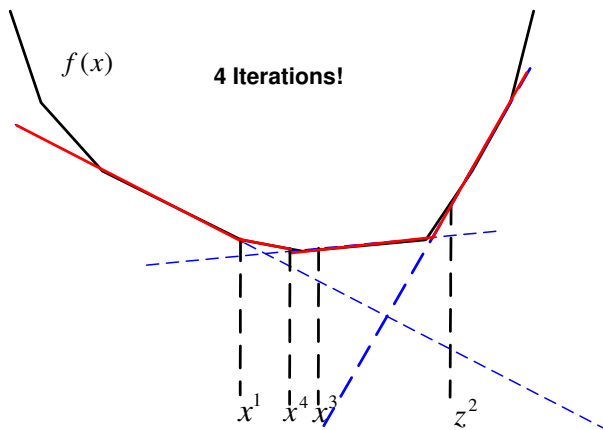
## PROXIMAL BUNDLE METHOD

$$f^M \equiv \tilde{f}, \quad x_{k+1} := \arg \min \left\{ \tilde{f}_k(x) + \frac{1}{2t_k} \|x - \hat{x}_k\|^2 : x \in X \right\}$$



## PROXIMAL BUNDLE METHOD

$$f^M \equiv \tilde{f}, \quad x_{k+1} := \arg \min \left\{ \tilde{f}_k(x) + \frac{1}{2t_k} \|x - \hat{x}_k\|^2 : x \in X \right\}$$



# PROXIMAL BUNDLE METHOD

## PROS × CONS

- 👍 only computes a single subgradient per iteration
- 👍 easy and reliable stopping test
- 👍 stable
- 👍 does not require  $X$  to be compact
- 👍 it is a descent method
- 👍 exploit good-quality initial points
- 👍 subproblem defining  $x_{k+1}$  can be kept small

## PROXIMAL BUNDLE METHOD

## PROS × CONS

- 👍 only computes a single subgradient per iteration
- 👍 easy and reliable stopping test
- 👍 stable
- 👍 does not require  $X$  to be compact
- 👍 it is a descent method
- 👍 exploit good-quality initial points
- 👍 subproblem defining  $x_{k+1}$  can be kept small
  
- 👎 convergence analysis is more involving...



# PROXIMAL BUNDLE METHOD

Let's consider a more economical model:

$$f_k^M(x) := \max_{j \in \mathcal{B}_k} \{f(x_j) + g_j^\top(x - x_j)\}$$

- ▶ The cutting-plane method takes  $\mathcal{B}_k := \{1, 2, \dots, k\}$ . We will consider  $\mathcal{B}_k \subset \{1, 2, \dots, k\}$  (or something a bit different)
- ▶ The method generates a sequence of trial points  $\{x_k\} \subset X$  by solving a QP:

$$x_{k+1} := \arg \min \left\{ f_k^M(x) + \frac{1}{2t_k} \|x - \hat{x}_k\|^2 : x \in X \right\}.$$

## SOLVING THE QP SUBPROBLEM

The QP

$$\min \left\{ f_k^M(x) + \frac{1}{2t_k} \|x - \hat{x}_k\|^2 : x \in X \right\}$$

can be rewritten as

$$\begin{cases} \min_{x,r} & r + \frac{1}{2t_k} \|x - \hat{x}_k\|^2 \\ \text{s.a} & f(x_j) + g_j^\top(x - x_j) \leq r, \quad j \in \mathcal{B}_k \\ & x \in X, r \in \mathfrak{R}. \end{cases}$$

We can apply specialized softwares.

## PROXIMAL BUNDLE METHOD

$$x_{k+1} := \arg \min \left\{ f_k^M(x) + \frac{1}{2t_k} \|x - \hat{x}_k\|^2 : x \in X \right\}.$$

A rule decide when to update the stability center  $\hat{x}_k$ . Such rule depends on the predicted decrease by the model  $f_k^M$

$$v_k = f(\hat{x}_k) - f_k^M(x_{k+1})$$

and a constant  $\kappa \in (0, 1)$ :

- Serious step: if  $f(x_{k+1}) \leq f(\hat{x}_k) - \kappa v_k$ , then

$$\hat{x}_{k+1} \leftarrow x_{k+1}$$

- Null step: if  $f(x_{k+1}) > f(\hat{x}_k) - \kappa v_k$ , then

$$\hat{x}_{k+1} \leftarrow \hat{x}_k$$

The serious-step sequence  $\{\hat{x}_k\}$  is a subsequence of  $\{x_k\}$

## NEXT ITERATE

## LEMMA

Suppose that  $X$  is a polyhedron or  $\text{ri}(X) \neq \emptyset$ . Then

$$x_{k+1} = \hat{x}_k - t_k \hat{g}_k \quad \text{com} \quad \hat{g}_k = p_f^k + p_X^k,$$

where  $p_f^k \in \partial f_k^M(x_{k+1})$  and  $p_X^k \in \partial i_X(x_{k+1})$ .  
( $i_X$  is the indicator function of  $X$ .)

Furthermore, the affine function

$$f_k^L(x) := f_k^M(x_{k+1}) + \langle \hat{g}_k, x - x_{k+1} \rangle$$

is a lower approximation for the model  $f_k^M$ :

$$f_k^L(x) \leq f_k^M(x) \quad \forall x \in X.$$

## OPTIMALITY MEASURE

## PROPOSITION

Let the predicted decrease and aggregate linearization error defined by

$$v_k := f(\hat{x}_k) - f_k^M(x_{k+1}) \quad \text{and} \quad \hat{e}_k := f(\hat{x}_k) - f_k^L(\hat{x}_k).$$

Then,

$$\hat{e}_k \geq 0, \quad \hat{e}_k + t_k \|\hat{g}_k\|^2 = v_k \geq 0 \quad \text{for all } k.$$

Furthermore

$$f(\hat{x}_k) \leq f(x) + \hat{e}_k + \|\hat{g}_k\| \|\hat{x}_k - x\| \quad \text{for all } x \in X \text{ and } k.$$

If  $(\hat{e}_k, \hat{g}_k) = 0$ , then  $\hat{x}_k$  is solution to the problem

## ALGORITHM: PROXIMAL BUNDLE METHOD

$$f_k^M(x) = \max_{j \in \mathcal{B}_k} \{f(x_j) + g_j^\top(x - x_j)\}, \quad x_{k+1} = \arg \min \left\{ f_k^M(x) + \frac{1}{2t_k} \|x - \hat{x}_k\|^2 : x \in X \right\}$$

- Step 0.** Choose  $\kappa \in (0, 1)$ ,  $t_1 \geq t_{\min} > 0$ ,  $x_1 \in X$  and tolerance  $\text{tol} > 0$ . Call the oracle to compute  $(f(x_1), g_1)$ . Define  $\hat{x}_1 \leftarrow x_1$ ,  $k \leftarrow 1$ ,  $\mathcal{B}_1 \leftarrow \{1\}$ ,
- Step 1.** Solve the QP to obtain  $x_{k+1}$ . Define  $\hat{g}_k \leftarrow (\hat{x}_k - x_{k+1})/t_k$ ,  $v_k \leftarrow f(\hat{x}_k) - \check{f}_k(x_{k+1})$ , and  $\hat{e}_k \leftarrow v_k - t_k \|\hat{g}_k\|^2$
- Step 2.** If  $\hat{e}_k \leq \text{tol}$  and  $\|\hat{g}_k\| \leq \text{tol}$ , stop:  $\hat{x}_k$  is an approximate solution
- Step 3.** Call the oracle to obtain  $(f(x_{k+1}), g_{k+1})$   
 Serious step. **If**  $f(x_{k+1}) \leq f(\hat{x}_k) - \kappa v_k$ , **then**  $\hat{x}_{k+1} \leftarrow x_{k+1}$  and choose  $t_{k+1} \geq t_k$   
 Null step. **Otherwise**, define  $\hat{x}_{k+1} \leftarrow \hat{x}_k$  and choose  $t_{k+1} \in [t_{\min}, t_k]$
- Step 4.** Choose  $\mathcal{B}_{k+1} \supset \{k+1, k^a\}$   
 Set  $k \leftarrow k+1$  and go back to Step 1.

## SOME COMMENTS

- ▶ Only 2 linearizations are required:  $f_k^L$  and  $f_k^{L^a}$ , i.e.,

$$\mathcal{B}_{k+1} = \{k+1, k^a\} \quad \text{suffices!}$$

- ▶ the prox-parameter  $t_k$  is non-increasing along null steps
- ▶ a simple heuristic to update the prox-parameter is the following
  - ▶ compute  $t_{\text{aux}} := t_k \left( 1 + \frac{(g_{k+1} - g_k)^\top (x_{k+1} - x_k)}{\|g_{k+1} - g_k\|^2} \right)$
  - ▶ if null step:  $t_{k+1} \leftarrow \min\{t_k, \max\{t_{\text{aux}}, t_k/2, t_{\min}\}\}$
  - ▶ if serious step:  $t_{k+1} \leftarrow \max\{t_k, \min\{t_{\text{aux}}, 10t_k\}\}$
- ▶ it is advisable to consider different tolerances for the measures  $\hat{e}_k$  and  $\hat{g}_k$
- ▶ the sequence  $\{f(\hat{x}_k)\}$  is non-increasing
- ▶ any accumulation point of  $\{\hat{x}_k\}$  is a solution to the problem