



# A variant of a projected variable metric method for bound constrained optimization problems

J. Frederic Bonnans

## ► To cite this version:

J. Frederic Bonnans. A variant of a projected variable metric method for bound constrained optimization problems. [Research Report] RR-0242, INRIA. 1983. inria-00076316

**HAL Id: inria-00076316**

**<https://inria.hal.science/inria-00076316>**

Submitted on 24 May 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



CENTRE DE ROCQUENCOURT

Institut National  
de Recherche  
en Informatique  
et en Automatique

Domaine de Voluceau  
Rocquencourt  
B.P. 105  
78153 Le Chesnay Cedex  
France  
Tél. (3) 954 90 20

Rapports de Recherche

N° 242

**A VARIANT OF A PROJECTED  
VARIABLE METRIC METHOD  
FOR BOUND CONSTRAINED  
OPTIMIZATION PROBLEMS**

Joseph Frédéric BONNANS

Octobre 1983

**A VARIANT OF A PROJECTED VARIABLE METRIC METHOD**  
**FOR BOUND CONSTRAINED OPTIMIZATION PROBLEMS**

Joseph Frédéric BONNANS

INRIA

Domaine de Voluceau

BP 105 - Rocquencourt

78153 LE CHESNAY CEDEX (France)

**RESUME :** Ce rapport étudie une méthode à métrique variable de résolution des problèmes d'optimisation sous contraintes de borne. La définition de la métrique utilise des idées de D.P. Bertsekas, tandis que la règle de recherche linéaire est une extension d'une règle due à P. Wolfe pour les problèmes sans contraintes. Un théorème de convergence est obtenu. Si l'information du second ordre est disponible, on en déduit une méthode superlinéairement convergente, même en l'absence des conditions de complémentarité stricte. Un résultat partiel de convergence est obtenu si les méthodes de quasi-Newton sont employées dans la définition de la métrique. La résolution numérique d'un problème de contrôle par cette méthode est présentée.

**ABSTRACT :** This paper studies a variable metric method for bound constrained optimization problems. The definition of the metric uses ideas of D.P. Bertsekas, and the line search rule is an extension of one studied by P. Wolfe for unconstrained problems. A global convergence theorem is obtained. If the second order information is available, a superlinearly convergent method is deduced, even without the strict complementary conditions. A partial convergence result is obtained if quasi-Newton methods are used to define the metric. Numerical resolution of a control problem by this method is presented.



## I - INTRODUCTION

We are concerned with the problem

$$(1.1) \quad \begin{cases} \text{Min } f(x), \\ x_i \geq 0, i = 1, \dots, n; \end{cases}$$

where  $x = (x_1, \dots, x_n)$  is in  $\mathbb{R}^n$  and  $f$  is continuously differentiable from  $\mathbb{R}^n$  into  $\mathbb{R}$ . This type of problem appears, for instance, when using the augmented lagrangian method for dualizing all constraints except bound constraints. It also appears when dualizing a problem with inequality constraints. Among many existing algorithms to solve (1.1) are

- projected gradient methods,
- methods consisting in a sequence of subproblems of type (1.1), in which some components of  $x$  for which the constraint is locally active are set to zero,
- recursive quadratic programming (RQP) (see for instance S.P. Han [5], M.J.D. Powell [9]) which consists here in the resolution at each iteration of a bound-constrained quadratic problem to get a descent direction of some exact penalty function.

In a recent article, D.P. Bertsekas [2] generalized the projected gradient method [1], i.e.

$$x^{k+1} = (x^k - t^k g^k)^+,$$

where  $(x)^+$  is the vector of  $\mathbb{R}^n$  whose  $i$ th. component is the positive part of  $x_i$ , to the variable metric method

$$x^{k+1} = (x^k - t^k M^k g^k)^+,$$

$M^k$  being a definite positive matrix whose some non diagonal elements are set to zero in order to get a decrease of the criterion for some  $t^k > 0$ . This

allows to design methods having a better convergence rate than the projected gradient method, but without the need to solve at each iteration a quadratic programming problem. D.P. Bertsekas [2] analyses the convergence of such methods associated to an Armijo-like linear search. He shows that if the condition number of  $M^k$  remains bounded and if  $x^k$  converges to  $x^*$ , the first-order necessary conditions hold at  $x^*$ . If, in addition, the strict complementarity conditions hold, the binding set of  $x^k$  is, for  $k$  great enough, identical to the binding set of  $x^*$  so that the problem reduces to an unconstrained problem. This allows to prove the superlinear convergence when  $M^k$  is deduced from the Hessian of  $f$  by some way.

There are two main differences between [2] and our study. First, we suppose that the linear search is subject to some natural extension of the Wolfe conditions [6], [12]. We also include the case of non strict complementarity conditions at the optimum (except in §5).

The organization of this paper is as follows. In §2 we state the algorithm and give a global convergence theorem. In §3 we restrict the class of admissible matrices  $M^k$ . This allows to prove that any isolated local minimum of the problem is a local attractor of the sequences generated by the algorithm. In §4 the matrix  $M^k$  is obtained in a simple way from the Hessian of  $f$ . We show that the method is then superlinearly convergent. In §5 we get the matrix  $M^k$  from  $B^k$ , approximation of the Hessian of  $f$  obtained with the BFGS algorithm. Here we need to assume that the strict complementarity conditions hold and that the condition number of  $B^k$  remains bounded. Numerical resolution of a bang-bang control problem with this method is exposed in §6.

## II - A CLASS OF PROJECTED VARIABLE METRIC METHODS

We consider problem (1.1) and say that  $\bar{x} \in \mathbb{R}^n$  is a Kuhn-Tucker point if the following necessary optimality conditions hold :

$$(2.1) \quad \begin{cases} \bar{x}_i \geq 0 \\ \bar{x}_i = 0 \implies \frac{\partial f}{\partial x_i}(\bar{x}) \geq 0 \\ \bar{x}_i > 0 \implies \frac{\partial f}{\partial x_i}(\bar{x}) = 0 \end{cases} \quad \forall i = 1, \dots, n.$$

Denote  $g(x) = \frac{\partial f}{\partial x}$ , the gradient of  $f$ . Let  $M$  be a  $n \times n$  definite positive matrix whose elements are denoted  $m_{ij}$ . Define

$$(2.2) \quad x(t) = (x - tMg(x))^+, \quad t \in \mathbb{R}^+.$$

We search conditions on  $M$  insuring that

$$(2.3) \quad \{x \text{ is not a Kuhn-Tucker point}\} \implies \{\exists t > 0 ; f(x(t)) < f(x)\}.$$

This may not happen in general, for a counter-example see [2]. We note

$$\mathbb{R}^{n+} = \{x \in \mathbb{R}^n ; x_i \geq 0, i = 1, \dots, n\},$$

and from now on we suppose :

$$(2.4) \quad f \text{ is continuously differentiable on } \mathbb{R}^{n+}.$$

Define

$$h(t) = f(x(t)), \quad t \in \mathbb{R}^+,$$

$$I = \{1, 2, \dots, n\},$$

$$I_B(x) = \{i \in I ; x_i = 0 \text{ and } g_i(x) > 0\},$$

$$I_L(x) = I - I_B(x).$$

Note that with (2.4),  $h$  has everywhere a right derivative, everywhere continuous except for, at most,  $n$  values of  $t$  we denote it  $h'(t)$ .

Proposition 2.1. (D.P. Bertsekas [2]). If (2.4) holds and

$$(2.5) \quad m_{ij} = 0 \text{ if } \{i \neq j \text{ and } i \text{ or } j \text{ is in } I_B(x)\},$$

then (2.3) holds and :

$$(2.6) \quad h'(0) \leq - \sum_{i,j \in I_L(x)} m_{ij} g_i g_j. \quad \square$$

Proof

The proof is given for the convenience of the reader. Define

$$I_M(x) = \{i \in I, x_i = 0 \text{ and } -(Mg)_i < 0\}.$$

Then

$$h'(0) = - \sum_{i \in I - I_M(x)} (Mg)_i g_i,$$

or equivalently

$$h'(0) = - g^t Mg + \sum_{i \in I_M(x)} (Mg)_i g_i.$$

If  $i \in I_M(x)$  and  $g_i \leq 0$ ,  $(Mg)_i g_i \leq 0$  hence

$$h'(0) \leq - g^t Mg + \sum_{i \in I_M(x)} (Mg)_i g_i,$$

with

$$\hat{I}_M(x) = \{i \in I ; x_i = 0, g_i > 0, -(Mg)_i < 0\}.$$

We now use hypothesis (2.5), which implies that  $\hat{I}_M(x) = I_B(x)$  and also

$$g^t M g = \sum_{i,j \in I_L(x)} m_{ij} g_i g_j + \sum_{i \in I_B(x)} m_{ii} g_i^2,$$

and

$$(Mg)_i = m_{ii} g_i, \quad \forall i \in \hat{I}_M(x) \equiv I_B(x).$$

This implies (2.6). As  $M$  is definite positive and some  $g_i$ ,  $i \in I_L(x)$ , is not nul if  $x$  is not a Kuhn-Tucker point, we deduce that  $h'(0) < 0$ . This proves the proposition.  $\square$

The preceding proposition gives a mean to define a descent path. To get an algorithm we must combine it with some strategy for the linear search. Alternatively to the Armijo-type method of [2] we propose the following Wolfe type method :

Definition 2.1. Let  $C_1, C_2$  be two constants such that

$$(2.7) \quad 0 < C_1 < C_2 < 1.$$

Then  $t$  is said admissible if the two following relations hold :

$$(2.8) \quad h(t) \leq h(0) + C_1 h'(0)t,$$

$$(2.9) \quad h'(t) \geq C_2 h'(0). \quad \square$$

These relations are identical to those of [12] if no bound constraint is active.

Proposition 2.2. We suppose that (2.4), (2.5) hold and that  $f$  has a finite minorant. Then if  $x$  is not a Kuhn-Tucker point, there exist four constants  $a, a', b', b$  such that

$$(i) \quad 0 < a \leq a' < b' \leq b < +\infty,$$



(ii) any  $t$  such that (2.8) and (2.9) hold is in  $[a, b]$ ,

(iii) for any  $t$  in  $[a', b'[,$  (2.8) and (2.9) hold.  $\square$

Proof

Supposing  $x$  is not a Kuhn-Tucker point, put

$$b = \sup \{t \geq 0 ; (2.8) \text{ holds}\}.$$

The continuity of  $h'(t)$  at  $t = 0$  implies the strict positivity of  $b$ . From the existence of a minorant of  $f$  we deduce that  $b$  is finite. Put

$$a = \inf \{t \geq 0 ; (2.9) \text{ holds}\}.$$

The fact that  $f$  has a finite minorant implies that  $a$  is not  $+\infty$ . Because  $h(t)$  is continuous at  $t = 0$ ,  $a$  is strictly positive. In addition ,

$$h(a) = h(0) + \int_0^a h'(t) dt,$$

and by the definition of  $a$ :

$$h(a) \leq h(0) + C_2 h'(0) a,$$

or equivalently

$$h(a) \leq h(0) + C_1 h'(0) a + (C_2 - C_1) h'(0) a.$$

Proposition 2.1. implies that  $h'(0) < 0$ . Using (2.7), we get

$$(2.10) \quad h(a) < h(0) + C_1 h'(0) a$$

so that  $a < b$  ; (ii) is a consequence of the definition of  $a$  and  $b$ . Put

$$b'' = \sup \{\tau \geq 0 ; (2.8) \text{ holds } \forall t \in [0, \tau]\}.$$

Necessarily  $a < b'' \leq b$  and

$$h(b'') = h(0) + C_1 h'(0) b''.$$

Subtracting (2.10) from it, we get

$$h(b'') - h(a) > C_1 h'(0) (b'' - a).$$

As  $h'$  is piecewise continuous, this implies the existence of  $a', b'$  such that  $a \leq a' < b' \leq b''$  and that

$$h'(t) \geq C_1 h'(0), \quad \forall t \in [a', b'],$$

and a fortiori with (2.7), as  $h'(0) < 0$  :

$$h'(t) \geq C_2 h'(0), \quad \forall t \in [a', b'].$$

As  $b' \leq b''$ , (2.8) also holds on  $[a', b']$  by the definition of  $b''$ . This proves (i) and (iii).  $\square$

For the sequel we need the following notations. Being given sequences  $\{x^k\}$  in  $\mathbb{R}^{n^+}$  and  $\{M^k\}$  in  $L(\mathbb{R}^n, \mathbb{R}^n)$ , every  $M^k$  being symmetric, positive definite, we denote :

$$I_B^k = I_B(x^k),$$

$$I_L^k = I_L(x^k),$$

$$n_L^k = \text{card } (I_L^k),$$

$$f^k = f(x^k),$$

$$g^k = g(x^k),$$

$g_L^k = \text{restriction of } g^k \text{ to indices of } I_L^k$  ;  $g_L^k$  is a vector of dimension  $n_L^k$ ,

$$x^k(t) = (x^k - t M^k g^k)^+, \quad t \in \mathbb{R}^+,$$

$$h^k(t) = f(x^k(t)),$$

$$||\cdot||_m \text{ } L^\infty \text{ norm on } \mathbb{R}^m, m \in \mathbb{N}, \text{ i.e. } ||z||_m = \sup \{|z_i|, i = 1, \dots, m\},$$

$$||\cdot||_{E,m} \text{ euclidean norm on } \mathbb{R}^m, m \in \mathbb{N}, \text{ i.e. } ||z||_{E,m} = \left( \sum_{i=1}^m (z_i)^2 \right)^{1/2}.$$

We denote  $m_{ij}^k$  the elements of  $M^k$  and use the following hypothesis on  $\{x^k\}, \{M^k\}$  :

$$(2.11) \quad m_{ij}^k = 0 \text{ if } \{i \neq j \text{ and } i \text{ or } j \text{ is in } I_B^k\}.$$

$$(2.12) \quad \left\{ \begin{array}{l} \text{There exists two constants } \alpha, \beta \text{ such that } 0 < \alpha < \beta < +\infty, \text{ and} \\ \text{such that all eigenvalues of } M^k \text{ remain in } [\alpha, \beta]. \end{array} \right.$$

We now state an algorithm of resolution of (1.1) :

#### Alg1

1° Choose  $x^0 \in \mathbb{R}^{n+}$  and  $(C_1, C_2)$  such that (2.7) holds. Set  $k = 0$ .

2° If  $x^k$  is a Kuhn-Tucker point, stop. Otherwise choose  $t^k$  such that

$$(2.13) \quad h^k(t^k) \leq h^k(0) + C_1 t^k h'^k(0),$$

$$(2.14) \quad h'^k(t^k) \geq C_2 h'^k(0).$$

Set

$$x^{k+1} = x^k(t^k),$$

$$k = k+1.$$

Go to 2°.  $\square$

Then a first result is :

Proposition 2.3. We suppose that (2.4) holds. Let  $\{x^k\}$  be a sequence generated by Alg1 such that (2.11), (2.12) hold. Then if  $\{x^k\}$  converges to some  $x^*$  :

- (i)  $x^*$  is a Kuhn-Tucker point,
- (ii) if  $i \in I$  is such that  $x_i^* = 0$  and  $g_i^* > 0$ , there exists  $k_0 \in \mathbb{N}$  such that  $k > k_0 \Rightarrow x_i^k = 0$ .  $\square$

Proof

(i) Because of hypothesis (2.4), (2.11) we can apply proposition 2.1, so that relation (2.6) at step  $k$  holds, and can be written here

$$h^k(0) \leq - \sum_{i,j \in I_L^k} m_{ij}^k g_i^k g_j^k ;$$

and so, with (2.12) :

$$(2.15) \quad h^k(0) \leq - \alpha \|g_L^k\|_{E, n_L^k}^2 .$$

If a subsequence of  $h^k(0)$  tends towards 0, we deduce easily from (2.4), (2.15) that  $g_L(x^*) = 0$ , i.e.  $x^*$  is a Kuhn-Tucker point. Suppose now that

$$\lim_{k \rightarrow +\infty} \sup h^k(0) = - \gamma < 0 .$$

From (2.14) we deduce that

$$(2.16) \quad \lim_{k \rightarrow +\infty} \inf [h^k(t^k) - h^k(0)] \geq (1-C_2)\gamma .$$

Define

$$I_{M^k}^k(t) = \{i \in I ; x_i^k(t) = 0 \text{ and } -(M^k g^k)_i < 0\} .$$

Then it is easy to see that

$$(2.17) \quad h'^k(t) = - (g^k)^t M^k g(x^k(t)) + \sum_{i \in I_{M^k}^k(t)} (M^k g^k)_i g_i(x^k(t)).$$

As  $\|x^{k+1} - x^k\|_n \rightarrow 0$ ,  $\|g^{k+1} - g^k\|_n \rightarrow 0$  so that the variations of  $h'$  induced by the first term of the right hand side of (2.17) tend to zero. The same result holds for contributions to the second term of indices  $i$  in  $I_{M^k}^k(t^k) \cap I_{M^k}^k(0)$ . Then, with (2.16), and as  $I_{M^k}^k(0) \subset I_{M^k}^k(t^k)$  :

$$(2.18) \quad \lim_{k \rightarrow +\infty} \sum_{i \in I_{M^k}^k(t^k) - I_{M^k}^k(0)} (M^k g^k)_i g_i^{k+1} \geq (1 - C_2)\gamma.$$

Let  $i$  be in  $I$ . Then :

- if  $x_i^* > 0$ ,  $i \notin I_{M^k}^k(t^k)$  for  $k$  great enough.
- if  $x_i^* = 0$  and  $g_i(x^*) = 0$  the contribution of indice  $i$  to the limit in (2.18) is null.
- if  $x_i^* = 0$  and  $g_i^*(x) < 0$ ,  $g_i^{k+1}$  will be negative for  $k$  great enough ; because of the definition of  $I_{M^k}^k(t)$ ,  $(M^k g^k)_i g_i^{k+1}$  will be negative for  $k$  great enough.
- if  $x_i^* = 0$  and  $g_i^* > 0$ ,  $g_i^{k+1}$  will be strictly positive for  $k > k_0$ , so that if  $i \in I_{M^k}^k(t^k)$  for  $k > k_0$ ,  $x_i^{k+1} = 0$ , and because of (2.11)  $x_i^{k'} = 0$  for any  $k' > k$ . Consequently  $i$  is not in  $I_{M^k}^k(t^k) - I_{M^k}^k(0)$  for  $k$  great enough.

As we considered all possible cases for  $i$ , we see that there is a contradiction in (2.18). So  $h'^k \rightarrow 0$  and that proves (i).

(ii) Let  $i$  be such that  $x_i^* = 0$  and  $g_i(x^*) > 0$ . Then  $g_i^k > 0$  for  $k > k_0$ . If  $x_i^k = 0$  for some  $k > k_0$ , then  $x_i^{k'} = 0$  for  $k' > k$  because of (2.11). If  $x_i^k > 0$ ,  $\forall k > k_0$ , then the limit of  $h'^k(0)$  is not zero, which is in contradiction with (i).  $\square$

If the strict complementarity conditions hold :

$$(2.19) \quad x_i^* = 0 \Rightarrow g_i(x^*) > 0, \forall i \in I,$$

we deduce from proposition 2.3 the

Corollary 2.1. Under hypothesis of proposition 2.3., if hypothesis (2.19) holds, there exists  $k_0 \in \mathbb{N}$  such that

$$\begin{aligned} x_i^* > 0 &\Rightarrow x_i^k > 0 \\ x_i^* = 0 &\Rightarrow x_i^k = 0 \end{aligned} \quad \forall i \in I, \forall k > k_0. \quad \square$$

In this case, after a finite number of iterations, Alg1 reduces to an algorithm for unconstrained minimization : a variable metric method associated to the classical Wolfe conditions for the linear search. The same result holds for the Armijo like linear search used in [2].

We now state a global convergence result for Alg1.

Theorem 2.1. We suppose that (2.4) holds and that  $f$  has a finite minorant. Let  $x^k$  be a sequence generated by Alg1 such that (2.11), (2.12) hold. We suppose also that

$$(2.20) \quad \sup \{ \|g(x)\|_n ; x \in \mathbb{R}^{n+} \text{ and } f(x) < f(x^0) \} < +\infty.$$

Then

$$\liminf_{k \rightarrow +\infty} \|g_L^k\|_{n_L^k} = 0. \quad \square$$

Proof

We remind that the norms  $\|\cdot\|_{n_L^k}$  and  $\|\cdot\|_{E, n_L^k}$  are equivalent. Suppose that

$$\liminf_{k \rightarrow +\infty} \|g_L^k\|_{E, n_L^k} = \gamma > 0.$$

Then, as proposition 2.1 holds, by (2.6) :

$$\liminf_{k \rightarrow +\infty} h^k(0) = -\alpha\gamma^2.$$

Relation (2.13) can be written as

$$-C_1 t^k h^k(0) \leq f^k - f^{k+1}$$

and so, for some  $k_0 \in \mathbb{N}$  :

$$\frac{C_1}{2} \alpha\gamma^2 t^k \leq f^k - f^{k+1}, \quad \forall k > k_0.$$

Denoting  $\bar{f}$  a minorant of  $f$  and summing on  $k$ , we get

$$(2.21) \quad \sum_{k=k_0}^{+\infty} t^k \leq \frac{2}{C_1 \alpha\gamma^2} (f^{k_0} - \bar{f}) < +\infty.$$

On the other hand, with (2.12)

$$\begin{aligned} \|x^{k+1} - x^k\|_n &\leq t^k \|M^k g^k\|_n, \\ &\leq t^k \beta \|g^k\|_{E, n}. \end{aligned}$$

Using (2.20) and (2.21) we deduce that

$$\sum_{k=k_0}^{+\infty} \|x^{k+1} - x^k\|_n < +\infty,$$

so that  $\{x^k\}$  converges. Then we have a contradiction with the conclusion of proposition 2.3.  $\square$

We now introduce additional constraints on the matrices  $M^k$  in order to improve the properties of the algorithm.

### III - AN ADDITIONAL CONDITION ON MATRICES $M^k$

Suppose that Alg1 computes some point  $x^k$  such that for some  $i$ ,  $x_i^k > 0$  is small and  $g_i(x^k) > 0$  is not too small. Then it is likely that  $i$  is in the set of binding constraints. To take this information into account we define two function from  $\mathbb{R}^{n+}$  into  $\mathbb{R}^{n+}$  called  $\epsilon(x)$  and  $\bar{\epsilon}(x)$  and we impose

$$(3.1) \quad m_{ij}^k = m_{ji}^k = 0 \text{ if } x_i^k \leq \epsilon_i(x^k), g_i^k > \bar{\epsilon}_i(x^k) \text{ and } i \neq j.$$

Note that relation (2.11) is a particular case of (3.1) with  $\epsilon_i^k = \bar{\epsilon}_i^k = 0$ . In the sequel we suppose (2.11) and (3.11) both satisfied though it might be possible to relax (2.11) for  $\bar{\epsilon}$  small enough. If  $x^k$  converges towards a Kuhn-Tucker point  $x^*$  it is desirable that  $\epsilon_i(x^k)$  and  $\bar{\epsilon}_i(x^k)$  converge to zero, so that, if the strict complementarity conditions hold, there is no restriction on  $M_{ij}^k$  for  $i, j$  such that  $x_i^* > 0$  and  $x_j^* > 0$ , in order to extend any variable metric method for unconstrained problems to (1.1). We give examples of such functions. For  $\nu > -1$ , we define

$$\phi_{iv} : \mathbb{R}^{n+} \rightarrow \mathbb{R}$$

by

$$\phi_i(x) = |x_i - (x_i - |g_i(x)|^\nu g_i(x))^+|, \quad i = 1, \dots, n.$$

For  $\mu, \nu > -1$ , we consider the following choices for  $\epsilon$  and  $\bar{\epsilon}$  :

$$(3.2) \quad \begin{cases} \epsilon_i(x) = \phi_{i\mu}(x), \\ \bar{\epsilon}_i(x) = \sum_{j=1}^n \phi_{j\nu}(x) \end{cases} \quad i = 1, \dots, n.$$

$$(3.3) \quad \begin{cases} \epsilon_i(x) = \sum_{j=1}^n \phi_{j\mu}(x) \\ \bar{\epsilon}_i(x) = \sum_{j=1}^n \phi_{j\nu}(x). \end{cases} \quad i = 1, \dots, n.$$

Note that in [2] are introduced some restrictions on  $M^k$  corresponding to the choice



$$(3.4) \quad \left\{ \begin{array}{l} \epsilon_i^k = \sum_{j=1}^n \phi_{jo}(x) \\ \bar{\epsilon}_i^k = 0 \end{array} \right\} \quad i = 1, \dots, n.$$

We state a first result giving conditions on  $\epsilon(x)$  and  $\bar{\epsilon}(x)$  for which a local minimum of  $f$  on  $\mathbb{R}^{n+}$  is a local attractor of the sequences computed by Alg1. For this purpose we do the following hypothesis on the linear search :

$$(3.5) \quad \left\{ \begin{array}{l} t^k \text{ is the last term of the finite sequence } t_p^k \text{ defined as follows :} \\ t_1^k = 1, \\ \text{we stop if (2.13), (2.14) hold with } t^k = t_p^k, \\ \text{if } h^k(t_p^k) \geq h^k, \text{ then } t_p^k < t_{p'}^k \text{ if } p' > p. \\ t_{p+1}^k \leq \delta t_p^k \text{ with } \delta \text{ given, } 1 < \delta < +\infty. \end{array} \right.$$

We also need the following hypothesis :

$$(3.6) \quad g(x) \text{ is locally lipschitzian.}$$

Note that (3.5) is natural and holds in practice. Then :

Theorem 3.1. We suppose that (2.4), (3.6) hold and that  $x^*$  is a strict minimum of  $f$  on some neighbourhood  $\mathcal{Q}(x^*)$  of  $x^*$  in  $\mathbb{R}^{n+}$ , such that  $x^*$  is the only Kuhn-Tucker point on  $\mathcal{Q}(x^*)$ . We suppose also that

$$(3.7) \quad \left\{ \begin{array}{l} \forall i \in I_B(x^*), \exists a_1 > 0 \text{ such that } \|x - x^*\|_n < a_1 \text{ implies :} \\ x_i \leq \epsilon_i(x) \\ \text{and} \\ g_i(x) > \bar{\epsilon}_i(x). \end{array} \right.$$

Let  $\{x^k\}$  be a sequence generated by Alg1 such that (2.11), (2.12), (3.1), (3.5) hold.

Then there exists  $a_2 > 0$  such that  $\|x^0 - x^*\|_n < a_2$  implies  $x^k \rightarrow x^*$ .  $\square$

Remark 3.1. If  $\epsilon$ ,  $\bar{\epsilon}$  are choosed by (3.2), (3.3) or (3.4), hypothesis (3.7) holds.  $\square$

To simplify the proof of the theorem we suppose that

$$I_B(x^*) = \{i \in I ; x_i^* = 0 \text{ and } g_i(x^*) > 0\}$$

corresponds to the first  $p$  indices. For any vector in  $\mathbb{R}^n$  we denote  $\tilde{y}$  the vector composed of its first  $p$  components and  $\tilde{y}$  the vector composed of the last  $q = n-p$  components.

Proof of Theorem 3.1.

Let us first prove that

$$(3.8) \quad \begin{cases} \text{If } \mathcal{U}(x^*) \text{ is closed, } \forall a > 0, \text{ there exists } f_a > f(x^*) \text{ such that} \\ x \in \mathcal{U}(x^*) \text{ and } f(x) < f_a \implies |x - x^*| < a. \end{cases}$$

If that was not true, for some  $a > 0$ ,  $\forall \tilde{f} > f(x^*)$ , there would exist  $x \in \mathcal{U}(x^*)$  such that  $f(x) < \tilde{f}$  and  $|x - x^*| > a$ . Set  $a$  and take a sequence of values of  $\tilde{f}$  ahving  $f(x^*)$  for limit : we obtain then a point  $\tilde{x}$  in  $\mathcal{U}(x^*)$ , different from  $x^*$ , with  $f(\tilde{x}) \leq f(x^*)$ , which is in contradiction with the hypothesis. Hence (3.8) holds.

As  $\mathcal{U}(x^*)$  contains some closed ball of center  $x^*$  and positive radius, we can suppose that  $\mathcal{U}(x^*)$  is closed.

Because of (3.1), (3.7), there exists  $a_3 > 0$  such that if  $|x^k - x^*| < a_3$  :

$$(3.9) \quad M_{ij}^k = 0 \text{ if } i \neq j \text{ and } i \text{ or } j \text{ is in } I_B(x^*).$$

Then we can write the relation between  $x^k$  and  $x^k(t)$  as

$$(3.10) \quad \begin{cases} \tilde{x}^{k+1}(t) = (\tilde{x}^k - tD^k \tilde{g}^k)^+, \\ \tilde{\tilde{x}}^{k+1}(t) = (\tilde{\tilde{x}}^k - t\tilde{M}^k \tilde{\tilde{g}}^k)^+, \end{cases}$$

with  $D^k$  diagonal and  $\tilde{M}^k$  a symmetric definite positive matrix whose eigenvalues are in  $[\alpha, \beta]$ . Obviously

$$(3.11) \quad |x_i^k(t) - x_i^*| \leq |x_i^k - x_i^*|, \quad \forall i \in I_B(x^*).$$

Note that  $\tilde{g}^k$  is lipschitzian and null at  $x^*$ , so that, denoting  $n_L^*$  the cardinal of  $I_L(x^*)$  :

$$||\tilde{g}^k||_{E, n_L^*} \leq L ||x^k - x^*||_n,$$

for some  $L > 0$ , and with (2.12), (3.9), (3.11) :

$$\begin{aligned} ||x^k(t) - x^*||_n &\leq ||x^k(t) - x^k + x^k - x^*||_n \\ &\leq t ||M^k g^k||_n + ||x^k - x^*||_n. \end{aligned}$$

As  $||M^k g^k||_n \leq \sqrt{n} ||M^k g^k||_{E, n} \leq \sqrt{n\beta} ||g^k||_{E, n}$ , we get

$$(3.12) \quad ||x^k(t) - x^*||_n \leq (\sqrt{n\beta} Lt + 1) ||x^k - x^*||_n.$$

Let  $d > 0$  be such that

$$\bar{B}(x^*, d) \equiv \{x \in \mathbb{R}^{n^+} ; ||x - x^*||_n \leq d\} \subset \mathcal{U}(x^*)$$

and  $a_4$  such that ( $\delta$  is the constant in (3.5)) :

$$(3.13) \quad \begin{cases} 0 < a_4 \leq a_3, \\ a_4 [\sqrt{n\beta} L + 2\delta + 1] \leq d. \end{cases}$$

We may suppose that,  $f_{a_4}$  being given by (3.8) :

$$(3.14) \quad ||x^k - x^*||_n < a_4,$$

$$(3.15) \quad f(x^k) < f_{a_4}.$$

We perform the linear search, starting with  $t = 1$ . Because of (3.12), (3.14) :

$$||x^k(1) - x^*||_n \leq a_4(\sqrt{n\beta L} + 1).$$

With (3.13) this implies that  $x^k(1) \in \bar{B}(x^*, d)$ , hence  $x^k(1) \in \mathcal{U}(x^*)$ . Suppose that  $t_p^k$  is for some  $p$  in  $\bar{B}(x^*, d)$ . If  $f(x^k(t_p^k)) \leq f(x^k)$ , then by (3.15)

$$(3.16) \quad ||x^k(t_p^k) - x^*||_n \leq a_4.$$

Then the maximum value of  $t_{p+1}^k$  is  $\delta t_p^k = \bar{t}$ . But for any  $t \leq \bar{t}$

$$\begin{aligned} ||x^k(t) - x^*||_n &\leq ||x^k(t) - x^k||_n + ||x^k - x^*||_n, \\ &\leq \delta ||x^k(t_p^k) - x^k||_n + ||x^k - x^*||_n, \\ &\leq \delta ||x^k(t_p^k) - x^*||_n + (\delta+1) ||x^k - x^*||_n \end{aligned}$$

and with (3.14), (3.16)

$$||x^k(t) - x^*||_n \leq a_4(2\delta+1),$$

so that, with (3.13),  $x^k(t) \in \bar{B}(x^*, d)$ ,  $\forall t \leq t_{p+1}^k$ . Hence the sequence  $x(t_p^k)$  remains in  $B(x^*, d)$ , so that

$$x^{k+1} \in B(x^*, d) \subset \mathcal{U}(x^*).$$

In addition,  $f(x^{k+1}) < f(x^k)$  so that, by (3.15) :

$$f(x^{k+1}) \leq f_{a_4},$$

and with (3.8) this implies that

$$||x^{k+1} - x^*||_n \leq a_4.$$

So the hypothesis made on  $x^k$  to insure that  $x^{k+1}$  is in  $\mathcal{U}(x)$  also hold for  $x^{k+1}$  : this proves that the all sequence  $\{x^k\}$  remains in  $\mathcal{U}(x^*)$ . As  $x^*$  is the only Kuhn-Tucker point on  $\mathcal{U}(x^*)$ , by theorem (2.1),  $\{x^k\}$  converges to  $x^*$ .  $\square$

We now specialize to a case where elements of  $M^k$  are simply deduced of elements of the hessian of  $f$ .

#### IV - A PROJECTED NEWTON METHOD

In this section we suppose that

(4.1)  $f$  is twice continuously differentiable on  $\mathbb{R}^{n+}$  and its hessian is locally lipschitzian.

We note  $H(x)$  the hessian of  $f$  at point  $x$ . We use  $H(x^k)$  in the design of  $M^k$  as in [2]:  $M^k$  is the solution of :

$$(4.2) \quad \begin{cases} (M^k)_{ij}^{-1} = 0 \text{ if (2.11) or (3.1) imply that } m_{ij}^k \text{ must be null,} \\ (M^k)_{ij}^{-1} = (H(x^k))_{ij} \text{ otherwise.} \end{cases}$$

If the inverse of  $M^k$  defined by (4.2) is really invertible,  $M^k$  is well defined and (2.11), (3.1) hold.

We suppose that  $x^0$  is in the neighbourhood of some Kuhn-Tucker point  $x^*$  checking the sufficient conditions for optimality

$$(4.3) \quad \begin{cases} z^t H(x^*) z > 0, \\ \forall z \in \mathbb{R}^n ; z_i g_i(x^*) = 0, i = 1, \dots, n. \end{cases}$$

It is not difficult to see that if  $x^k$  is sufficiently close to  $x^*$ , relations (2.11), (3.1), (4.2) define in a unique way a symmetric definite positive matrix  $M^k$  such that (2.12) holds. Then Theorem 3.1. gives sufficient conditions of convergence of  $x^k$  towards  $x^*$  if  $x^0$  is close enough from  $x^*$ . As  $H(x)$  is not everywhere positive definite, some corrections must be brought to the method to get a well-defined and globally convergent algorithm. This situation is similar to the unconstrained case. We now focus on the rate of convergence. If the conditions of strict complementarity hold, after a finite number of iterations

the algorithm reduces to Newton's method for an unconstrained problem. Then it is easy to obtain a superlinear convergence rate. The following theorem makes no hypothesis of strict complementarity.

Theorem 4.1. Let  $\{x^k\}$  be generated by Alg1 with  $M^k$  defined by (4.2). We suppose that (4.1) holds and that  $\{x\}$  has a limit  $x^*$  such that (2.1), (4.3) hold.

We suppose that

$$\left. \begin{array}{l} (4.4i) \quad \frac{||\epsilon(x^k)||_n}{||\tilde{g}(x^k)||_q} \rightarrow 0 \\ (4.4ii) \quad ||\bar{\epsilon}(x^k)||_n \rightarrow 0 \end{array} \right\} \text{ when } k \rightarrow +\infty.$$

Then, if the linear search checks (3.5), there exists  $C_1^0 > 0$  such that, if  $C_1 < C_1^0$ ,  $x^k \rightarrow x^*$  superlinearly.  $\square$

Remark 4.1. Suppose that  $\epsilon(x)$ ,  $\bar{\epsilon}(x)$  are chosen by (3.2) or (3.3) with  $\mu > 0$  and  $\nu > -1$ . Then (3.7) holds so that  $x^*$  is a local attractor and, after a finite number of iterations,  $x_i = 0$  for  $i \in I_B$ . Then, as  $\mu > 0$ , it is not difficult to see that (4.4) holds. If  $\epsilon(x)$ ,  $\bar{\epsilon}(x)$  are chosen by (3.4), however, (4.4) does not hold.  $\square$

#### Proof of Theorem 4.1.

We first check that the choice  $t^k = 1$  in the linear search insures that  $x^k \rightarrow x^*$  superlinearly. For  $k$  great enough, because of corollary 2.1. and of the way we choose  $M^k$ , the indices  $i$  for which  $x_i = 0$  and  $g_i^* > 0$  have no action on the algorithm; so, to simplify the proof, we suppose there is no such indice, i.e.  $I_B(x^*) = \emptyset$ . Let  $k$  be fixed. We denote

$$I_P = \{i \in I ; x_i^k \leq \varepsilon_i(x^k) \text{ and } g_i(x^k) > \bar{\varepsilon}_i(x^k)\},$$

$$I_Q = I - I_P,$$

$$n_P = \text{card}(I_P), n_Q = \text{card}(I_Q).$$

For  $z \in \mathbb{R}^n$  we denote  $z_P, z_Q$  the restriction of  $z$  to indices of  $I_P$  and  $I_Q$ . Then the iteration with  $t^k = 1$  can be written as

$$(4.5) \quad \hat{x}_P^k = x_P^k - D^k g_P^k,$$

$$(4.6) \quad \hat{x}_Q^{k+1} = x_Q^k - \hat{M}^k g_Q^k;$$

and

$$(4.7) \quad x_P^{k+1} = (\hat{x}_P^k)^+,$$

$$(4.8) \quad x_Q^{k+1} = (\hat{x}_Q^k)^+.$$

Here  $D$  (resp.  $\hat{M}$ ) is a diagonal (resp. symmetric) definite positive matrix whose elements are obtained in a direct way from  $M^k$ . Denote

$$I^\# = \{i \in I ; x_i^* > 0\}.$$

Because of (4.4), when  $k \rightarrow \infty$ ,  $\varepsilon_i^k \rightarrow 0, \forall i \in I$ , hence for some  $k_0, k > k_0$  implies

$$x_i^k > \varepsilon_i^k, \forall i \in I^\#.$$

Consequently, for  $k > k_0, I^\# \cap I_P = \emptyset$  i.e.

$$x_i^* = 0, \forall i \in I_P.$$

With (4.5), (4.7) and noticing that  $g_i(x^k) \geq 0$  if  $i \in I_P$ , we get

$$(4.9) \quad |x_i^{k+1} - x_i^*| \leq |x_i^k - x_i^*|, \forall i \in I_P.$$



As  $I_B(x^*) = \emptyset$ ,  $\tilde{g}(x^k) = g(x^k)$  and by (4.4), for any  $\xi > 0$  there exists  $k > 0$  such that if  $k > k_1$

$$|x_i^k - x_i^*| \leq \xi \|g^k\|_n, \forall i \in I_P.$$

Now,  $g(x^*) \equiv 0$  and (4.1) implies that  $g$  is lipschitzian with some constant  $L$ , so that

$$(4.10) \quad |x_i^k - x_i^*| \leq \xi L \|x^k - x^*\|_n, \forall i \in I_P.$$

With (4.9) this implies

$$(4.11) \quad |x_i^{k+1} - x_i^*| \leq \xi L \|x^k - x^*\|_n, \forall i \in I_P.$$

We now consider (4.6), (4.8). Let  $q$  be the cardinal of  $I_Q$  and, by reindexing,  $x^k = (x_P^k, x_Q^k)$ . Define

$$\begin{aligned} F : \mathbb{R}^q &\rightarrow \mathbb{R} \\ x_Q &\rightarrow f(x_P^k, x_Q). \end{aligned}$$

Then (4.6) is the Newton iteration to solve the equation

$$(4.12) \quad \frac{\partial F}{\partial x_Q} = \frac{\partial f}{\partial x_Q}(x_P^k, x_Q) \equiv 0.$$

As (4.1) (4.3) holds and as  $\frac{\partial f}{\partial x}(x_P^*, x_Q^*) \equiv 0$ , we can apply the implicit function theorem. It states that in the neighbourhood of  $x_P^*$ ,  $\hat{x}_Q$  solution of (4.12) can be considered as lipschitzian function of  $x_P$ , i.e.

$$\|\hat{x}_Q - x_Q^*\|_{n_Q} \leq a_1 \|x_P^k - x_P^*\|_{n_P}$$

for some  $a_1 > 0$ . Because of (4.10) :

$$(4.13) \quad \|\hat{x}_Q - x_Q^*\|_{n_Q} \leq a_1 L \xi \|x^k - x^*\|_n.$$

We now evaluate  $||x_Q^{k+1} - \hat{x}_Q||_{n_Q}$ . We have

$$\frac{\partial f}{\partial x_Q}(x_P^k, x_Q^{k+1}) = \frac{\partial f}{\partial x_Q}(x_P^k, x_Q^k) + \frac{\partial^2 f}{\partial x_Q^2}(x_P^k, x_Q^k)(x_Q^{k+1} - x_Q^k) + \varepsilon(x_Q^{k+1} - x_Q^k)$$

and (4.1) implies that for some  $a_2 > 0$  :

$$||\varepsilon(x_Q^{k+1} - x_Q^k)||_{n_Q} \leq a_2 ||x_Q^{k+1} - x_Q^k||_{n_Q}^2.$$

As we did a Newton step to solve  $\frac{\partial f}{\partial x_Q}(x_P^k, x_Q^k) = 0$ , we get

$$\begin{aligned} ||\frac{\partial f}{\partial x_Q}(x_P^k, x_Q^{k+1})||_{n_Q} &\leq a_2 ||x_Q^{k+1} - x_Q^k||_{n_Q}^2 \\ &\leq a_3 ||g_Q^k||_{n_Q}^2 \\ &\leq a_4 ||x^k - x^*||_n^2, \end{aligned}$$

for some  $a_3, a_4 > 0$ , independant of  $x_P^k$  in some neighbourhood of  $x^*$ , and this implies

$$||x_Q^{k+1} - \hat{x}_Q||_{n_Q} \leq a_5 ||x^k - x^*||_n^2.$$

This with (4.11) implies the superlinear convergence of  $\{x^k\}$  towards  $\{x\}$ .

We now check that the choice  $t^{k=1}$  is compatible with conditions (2.13) (2.14), (3.5) on the linear search. Condition (3.5) says that the first trial of  $t^k$  is 1, so we have to check (2.13), (2.14). Remind that we cancelled the components for wich  $g_i(x^*) > 0$ , so that  $g(x)$  is null at  $x^*$  and because of (4.1) :

$$||g(x)||_n \leq L ||x - x^*||_n$$

in a neighbourhood of  $x^*$ , for some  $L > 0$ . As the step  $t=1$  gives a superlinear convergence, we easily deduce that condition (2.14) holds for any  $C_2$  in  $]0, 1[$  if  $k$  is great enough. We now want to check (2.13). In a neighbourhood of  $x^*$  :

$$f(x) \leq f(x^*) + a_5 ||x - x^*||_n^2.$$

For any  $\xi > 0$ , if  $k > k_4$  (and if  $t^k = 1$ ) :

$$||x^{k+1} - x^*||_n \leq \xi ||x^k - x^*||_n,$$

and so :

$$(4.14) \quad f(x^{k+1}) \leq f(x^*) + a_5(\xi)^2 ||x^k - x^*||_n^2.$$

We also have for  $k > k_6$  :

$$f(x^k) \geq f(x) + \frac{1}{2}(x^k - x^*)^T H(x^*) (x^k - x^*) - \xi ||x^k - x^*||_n^2,$$

and also

$$||g(x^k) - H(x^*)(x^k - x^*)|| \leq \xi ||x^k - x^*||,$$

so that for  $k > k_7$

$$f(x^k) \geq f(x^*) + \frac{1}{2}(g^k)^T H(x^*)^{-1} g^k - \xi ||x^k - x^*||_n^2,$$

and with (4.14), for  $k > k_8$  :

$$f(x^k) - f(x^{k+1}) \geq \frac{1}{2}(g^k)^T H(x^*)^{-1} g^k - \xi ||x^k - x^*||_n^2.$$

Then (2.14) certainly holds if

$$(4.15) \quad \frac{1}{2} (g^k)^T H(x^*)^{-1} g^k - \xi ||x^k - x^*||_n^2 \geq C_1 (g^k)^T M g^k.$$

Denote  $\lambda_m(B)$ ,  $\lambda_M(B)$  the smallest and largest eigenvalue of a  $n \times n$  matrix  $B$ . The definition of  $M^k$  implies that

$$\lambda_m(H(x^k)) \leq \lambda_m(M^k) \leq \lambda_M(M^k) \leq \lambda_M(H(x^k)).$$

As  $\lambda_m(H(x^k)) \rightarrow \lambda_m(H(x^*))$  and  $\lambda_M(H(x^k)) \rightarrow \lambda_M(H(x^*))$ , we have, for any  $\epsilon > 0$ , if  $k > k_8$  :

$$(g^k)^t M^k g^k \geq (\lambda_m(H(x)) - \epsilon) \|g^k\|_{E,n}^2,$$

so that (2.14) holds for  $k > k_9$  if

$$(4.16) \quad C_1 < \frac{1}{2} \frac{\lambda_m(H(x))}{\lambda_M(H(x))}. \quad \square$$

Remark 4.2. In the unconstrained case,  $M^k = H(x^k)^{-1}$  so that (2.15) holds for  $k$  great enough if  $C_1 < \frac{1}{2}$ . Condition (4.16) is far more severe.  $\square$

Remark 4.3. If, instead of (4.4i) we suppose that  $\|\epsilon(x^k)\|_n \leq a_1 \|\tilde{g}(x^k)\|_q^2$  for some  $a_1 > 0$ , we easily deduce from the proof of Theorem 4.1. that  $x^k \rightarrow x$  quadratically. This holds with choice (3.2) or (3.3) if  $\mu \geq 1$ .  $\square$

Remark 4.4. Theorem 4.1. can be easily extended to some modifications of Newton's method, for instance if in (4.2)  $H(x^k)$  is not the true hessian of  $f$  but tends towards the hessian of  $H(x^*)$  when  $x^k \rightarrow x^*$ .  $\square$

## V - A PROJECTED QUASI NEWTON METHOD

We now consider the case when matrices  $B^k$ , obtained by the BFGS formula [4] are used to design  $M^k$ , in an analogous way to (4.2). We first remind the BFGS formula. A sequence  $\{x^k\}$  in  $\mathbb{R}^n$  be given, we note

$$\begin{aligned} s^k &= x^{k+1} - x^k, \\ y^k &= g^{k+1} - g^k. \end{aligned}$$

A symmetric definite positive matrix  $B^0$  being given, matrices  $B^k$  are defined in a recurrent way by

$$(5.1) \quad B^{k+1} = B^k + \frac{y^k (y^k)^t}{(y^k)^t s^k} - \frac{B^k s^k (B^k s^k)^t}{(s^k)^t B^k s^k}.$$

These matrices can be viewed as an approximation of the hessian of  $f$ . We note  $\sigma^k = (y^k)^t s^k$ . If  $B^k$  is definite positive, so is  $B^{k+1}$  if and only if  $\sigma^k > 0$ ; so to preserve the positive definiteness of  $\{B^k\}$  if  $\sigma^k \leq 0$ , we generalize (5.1) as in [9] :

$$(5.2) \quad \left\{ \begin{aligned} B^{k+1} &= B^k + \frac{z^k (z^k)^t}{(z^k)^t s^k} - \frac{B^k s^k (B^k s^k)^t}{(s^k)^t B^k s^k} \\ \text{with} \\ \theta^k &= \frac{0.8 (s^k)^t B^k s^k}{(s^k)^t B^k s^k - (y^k)^t s^k}, \\ z^k &= \begin{cases} \theta^k y^k + (1 - \theta^k) B^k s^k & \text{if } \sigma^k \leq 0.2 (s^k)^t B^k s^k, \\ y^k & \text{otherwise.} \end{cases} \end{aligned} \right.$$

Now we define  $M^k$  by

$$(5.3) \quad \left\{ \begin{aligned} (M^k)^{-1}_{ij} &= 0 \text{ if } (M^k)_{ij} \text{ is imposed to be null by (2.11) or (3.1),} \\ (M^k)^{-1}_{ij} &= B^k_{ij} \text{ otherwise.} \end{aligned} \right.$$

Then if no constraint is active, Alg1 reduces to the algorithm studied by M.J.D. Powell [8], who obtained a global convergence result if  $f$  is convex and proved that the superlinear convergence occurred if  $x^k$  had some limit point  $x^*$  such that  $H(x^*)$  is definite positive. None of these results seem easy to extend. However, we have the following result :

Theorem 5.1. Let  $x^k$  be generated by Alg1 with  $M^k$  defined by (5.3). We suppose that the condition number of  $\{B^k\}$  remains bounded. Then

$$\lim_{k \rightarrow \infty} \inf ||g_L^k||_{n_L^k} = 0.$$

If in addition  $\{x^k\}$  has some limit point  $x^*$  such that (4.1), (4.3) and the strict complementarity conditions (2.20) hold, then  $x^k \rightarrow x^*$  superlinearly.  $\square$

#### Proof

The first part of the theorem concerning the global convergence is a direct application of theorem 2.1. Now, if (2.20) hold, because of corollary 2.1, for  $k > k_0$  the problem reduces to the algorithm of [8] for unconstrained problems ; in [8], superlinear convergence under our hypothesis was proved.  $\square$

Remark 5.1. Problem (1.1) can be solved by a recursive quadratic programming approach associated to a Quasi-Newton formula [5], [9]. For this kind of algorithm, it seems that no result concerning the convergence, stronger than those of theorem 5.1, have been yet obtained (see [10], [3]).  $\square$

## VI - NUMERICAL RESULTS

We applied the algorithm exposed in §5 to a discretised version of the following control problem. Let  $T > 0$  be given and denote

$$\begin{cases} \Omega = ]0,1[ \\ Q = \Omega \times ]0,T[. \end{cases}$$

The state equation is

$$(6.1) \quad \begin{cases} \frac{\partial y}{\partial t} - \Delta y = 0 \text{ in } Q, \\ \frac{\partial y}{\partial n}(t,1) = u(t), \\ \frac{\partial y}{\partial n}(t,0) = 0, \\ y(0,x) = 0 \quad \forall x \in ]0,1[. \end{cases}$$

The criterion to be minimized is

$$J(u) = \int_{\Omega} (y(u,T,x) - y_d(x))^2 dx ,$$

with

$$y_d(x) = \begin{cases} 0 & \text{if } 0 < x < 1/2, \\ 1 & \text{if } 1/2 < x < 1. \end{cases}$$

Note that for  $u$  given in  $L^2(0,T)$ , (6.1) has a unique solution with trace at  $t=T$  in  $L^2(\Omega)$  (see J.L. Lions [7] for the mathematical aspects) so that  $J(u)$  makes sense. We include some bounds constraints  $-0.5 \leq u \leq 2$ , so that the problem is

$$\begin{cases} \text{Min } J(u), \\ -0.5 \leq u(t) \leq 2, \quad \forall t \in [0,T]. \end{cases}$$

The problem is discretized in time using the usual implicit Euler scheme, then in space with a finite-difference method equivalent to P1 finite elements. We denote

$$\begin{cases} n_x & \text{number of space steps,} \\ n & \text{number of time steps.} \end{cases}$$

In our test  $n_x$  is set to 40 and  $n$  varies from 50 to 200. Note that the discretized problem is not well conditioned : the hessian has the eigenvalue 0 with an order of multiplicity of at least  $n - n_x$ .

The solution of the discretized problem has a bang-bang structure : except for a small number of time steps, the value of the optimal control is equal to one bound. In practice, we computed solutions with 0 or 1 non binding variable.

The parameters  $\epsilon$  is given by (3.3) with  $\mu = 1$  and  $\bar{\epsilon}$  is set to zero. We denote

$$\begin{cases} \text{nit} & \text{number of iterations in optimization,} \\ \text{nf} & \text{number of calls of function and gradient,} \\ ||g||_r & L^\infty \text{ norm of residual gradient,} \\ \text{CT} & \text{computing time in seconds.} \end{cases}$$

All tests were run in HB-68 Multics system of INRIA.

The results are stated in table 1.

These results show that  $\text{nf}$  does not increase too much with respect to  $n$ . However as computations are made on matrices of dimension  $n^2$ , the total amount of computation is between a linear and a quadratic function of  $n$ .



n	nit	nf	$  g  _r$	CT	CT/n	CT/n <sup>2</sup>
50	41	73	1.5 e-10	80	1.6	3.2 e-2
100	60	119	.5 e-12	279	2.8	2.8 e-2
150	61	109	1. e-12	504	3.36	2.24 e-2
200	64	114	1.5 e-13	788	3.94	1.97 e-2

## VII - CONCLUSION

This article studies the projected variable metric method of D.P. Bertsekas [2] for bound constrained problems associated to a new line search which reduces in the unconstrained case to one of those studied in P. Wolfe [12]. Some new formulae of the parameters used to define the metric are obtained. With them, we show that a local minimum is a local attractor point of sequences computed by the algorithm and that, if the second-order information is available, the superlinear convergence occurs even if the strict complementarity conditions do not hold.

If the variable metric is obtained through the BFGS quasi-Newton formula, a partial result of convergence is stated as a simple consequence of a theorem of M.J.D. Powell [8].

This article does not consider extensions of the conjugate gradient method. However, it is obvious that the memoryless quasi-Newton algorithms as those of D.F. Shanno [11] have a direct extension for bound constrained problems, which should be competitive for large scale problems.

### Acknowledgments

The author thanks D. Gabay, C. Lemarechal and E. Panier for helpful comments and discussions.

## REFERENCES

- [1] BERTSEKAS D.P. (1976). On the Goldstein-Levitin-Poljak gradient projection method. IEEE Trans. AC-21, pp. 174-184.
- [2] BERTSEKAS D.P. (1982). Projected Newton methods for optimization problems with simple constraints. SIAM J. on Control and Opt., 20, pp. 221-246.
- [3] BOGGS P.T. - TOLLE J.W. - PYNG WANG (1982). On the local convergence of Quasi-Newton methods for constrained optimization. SIAM J. on Control and Opt., 20, pp. 161-171.
- [4] DENNIS J.E. - MORE J.J. (1977). Quasi-Newton methods, motivation and theory. SIAM Review 19. pp. 46-89.
- [5] HAN S.P. (1977). A globally convergent method for non linear programming. J.O.T.A. 22. pp. 297-309.
- [6] LEMARECHAL C. (1981). A view of line-searches in optimization and optimal control. Auslender, Oetti, Stoer eds. Lectures notes in Control and Information Sciences n° 30, Springer Verlag. pp. 59-78.
- [7] LIONS J.L. (1968). Contrôle optimal de systèmes gouvernés par des équations aux dérivées partielles. Dunod, Gauthier-Villars, Paris.
- [8] POWELL M.J.D. (1976). Some global convergence properties of a variable metric algorithm for minimization without exact line search, in : Nonlinear Programming, SIAM-AMS Proc. 9 Providence, R.I.
- [9] POWELL M.J.D. (1978). A fast algorithm for nonlinearly constrained optimization calculations, in Numerical Analysis, Lecture Notes in Mathematics n° 630, Springer-Verlag.

- [10] POWELL M.J.D. (1978). The convergence of variable metric methods for nonlinearly constrained optimization calculations, in Nonlinear Programming 2, O.L. Mangasarian, R.R. Meyer, S.M. Robinson eds. Academic Press, New York.
- [11] SHANNO D.F. (1978). Conjugate gradient methods with inexact searches. Maths. Op. Res. 3, pp. 244-256.
- [12] WOLFE P. (1969). Convergence conditions for ascent methods (I), SIAM Review 11, pp. 226-235.
- [13] WOLFE P. (1971). Convergence conditions for ascent methods (II), SIAM Review 13, pp. 185-188.

7)

8)

9)

10)

11)

12)