

E. POLAK

G. RIBIERE

**Note sur la convergence de méthodes de
directions conjuguées**

Revue française d'informatique et de recherche opérationnelle. Série rouge, tome 3, n° R1 (1969), p. 35-43

http://www.numdam.org/item?id=M2AN_1969__3_1_35_0

© AFCET, 1969, tous droits réservés.

L'accès aux archives de la revue « Revue française d'informatique et de recherche opérationnelle. Série rouge » implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/legal.php>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

Article numérisé dans le cadre du programme
Numérisation de documents anciens mathématiques
<http://www.numdam.org/>

NOTE SUR LA CONVERGENCE DE METHODES DE DIRECTIONS CONJUGUEES

par E. POLAK ⁽¹⁾ et G. RIBIERE ⁽²⁾

Résumé. — Ce papier utilise un théorème général qui donne des conditions suffisantes de convergence d'une classe d'algorithmes à direction de déplacement, pour construire un algorithme convergent de directions conjuguées, destiné à la minimisation non contrainte de fonctions réelles dans \mathbb{R}^n . Cet algorithme est dérivé, par une modification simple, de la méthode de Fletcher-Reeves. On donne quelques résultats numériques pour illustrer le comportement de ce nouvel algorithme.

On montre également que la convergence de la méthode de Newton à pas variable, aussi bien que celle de la méthode de la plus grande pente, peut être obtenue à l'aide de ce théorème de convergence.

I. INTRODUCTION

On va s'attacher à démontrer la convergence de plusieurs algorithmes destinés à trouver le minimum d'une fonction. Parmi ceux-ci on retrouve la méthode de plus grande pente et la méthode de Newton à pas variable.

Tout d'abord on donnera un théorème de convergence assez général, puis on l'appliquera à divers algorithmes. On démontrera en particulier la convergence d'une nouvelle méthode de gradient conjugué qui est une modification de la méthode de Fletcher-Reeves (1).

2. THEOREME GENERAL DE CONVERGENCE

On définit tout d'abord un algorithme assez général destiné à minimiser une fonction continue $f(x)$, $x \in \mathbb{R}^n$. On part de x_0 arbitraire et on forme successivement :

$$(2.1) \quad x_{i+1} = a(x_i) \quad ; \quad f(x_{i+1}) < f(x_i),$$

(le procédé s'arrête à un point x_k si $f(x_{k+1}) \geq f(x_k)$).

(1) Department of electrical engineering and computer sciences. University of California Berkeley, California.

(2) C.N.R.S. Institut Blaise-Pascal, Service de Développement, Paris.

2.1. Définition : — On dit que x est un point « désirable » du problème de minimisation si

$$(2.2) \quad f(a(x)) \geq f(x).$$

2.2. Lemme (E. Polak [2]).

Si, pour tout x non « désirable », il existe des constantes dépendant de x , $\varepsilon(x) > 0$ et $\delta(x) > 0$, telles que

$$(2.3) \quad \|x - x_i\| \leq \varepsilon(x) \Rightarrow f(x_i) - f[a(x_i)] \geq \delta(x),$$

alors toute sous-suite convergente de la suite définie par l'algorithme (2.1) a pour point d'accumulation un point « désirable ».

Démonstration :

Voir l'article de E. Polak « On the convergence of optimization algorithms » dans ce même numéro.

Considérons maintenant une classe d'algorithmes pour minimiser f , supposée continûment différentiable ($f \in C^1(\mathbf{R})$), qui sont caractérisés par les propriétés suivantes :

On part de x_0 arbitraire, puis on calcule les itérés successifs par

$$(2.4) \quad x_{i+1} = x_i + \mu_i p_i$$

où p_i , qui ne dépend pas uniquement de x_i , est la direction de déplacement astreinte à la condition

$$(2.5) \quad -f'(x_i) \cdot p_i \geq \alpha \|f'(x_i)\| \|p_i\|$$

où $\alpha > 0$ est indépendant de x_i et μ_i est calculé de telle sorte que p_i soit orthogonal au gradient de f en x_{i+1} , plus précisément

$$(2.6) \quad \mu_i = \inf \{ \mu \mid \mu \geq 0, \quad f'(x_i + \mu p_i) \cdot p_i = 0 \}.$$

On peut dire aussi que x_{i+1} réalise le premier minimum de f sur la demi-droite d'origine x_i et de direction p_i . Bien entendu $f(x_{i+1}) \leq f(x_i)$.

2.3. Définition : — On dit que x est stationnaire si

$$(2.7) \quad f'(x) = 0.$$

2.4. Théorème. — Toute sous-suite convergente de la suite définie par un algorithme du type (2.4), (2.5) et (2.6) a pour point d'accumulation un point stationnaire.

Démonstration

On peut toujours supposer que les directions p sont normalisées de telle sorte que $\|p\| = 1$ et on définit :

$$S = \{ p \mid \|p\| = 1 \}$$

Soit x_0 un point non stationnaire. Alors

$$(2.8) \quad \|f'(x_0)\| = \nu > 0.$$

Du fait de la continuité de la différentielle de f , il existe une boule $B_\varepsilon(x_0)$ de centre x_0 et rayon ε telle que

$$(2.9) \quad x \in B_\varepsilon(x_0) \Rightarrow \|f'(x)\| \geq \nu/2.$$

On définit

$$g(x, p) = -f'(x) \cdot p.$$

Alors d'après (2.5) et (2.9) on a

$$(2.10) \quad x \in B_\varepsilon(x_0) \Rightarrow g(x, p) \geq \alpha \frac{\nu}{2} \|p\| = \beta$$

L'uniforme continuité de $g(x, p)$ dans le produit $B_\varepsilon(x_0) \times S$ entraîne qu'il existe λ_m tel que

$$(2.11) \quad \lambda \in [0, \lambda_m] \Rightarrow |g(x, p) - g(x + \lambda p, p)| \leq \beta/2.$$

Utilisant le développement de Taylor on a

$$f(x + \lambda_m p) = f(x) - \lambda_m g(x + \lambda p, p)$$

où

$$\lambda \in [0, \lambda_m].$$

D'après (2.10) et (2.11) on a aussi

$$x \in B_\varepsilon(x_0) \Rightarrow g(x + \lambda p, p) \geq \beta/2$$

d'où

$$f(x) - f(x + \lambda_m p) \geq \lambda_m \beta/2$$

mais d'après la définition (2.6)

$$f(x + \mu p) \leq f(x + \lambda_m p)$$

d'où

$$(2.12) \quad f(x) - f(x + \mu p) \geq \lambda_m \frac{\alpha}{4} \|f'(x_0)\| > 0.$$

Avec cela, on va montrer tout d'abord que x « désirable » est équivalent à x stationnaire dans le cas d'un algorithme du type (2.4), (2.5) et (2.6).

En effet, supposons que x_0 soit « désirable »; alors si x_0 n'était pas stationnaire on pourrait appliquer (2.12) avec $x = x_0$ d'où

$$0 \geq f(x_0) - f(x_0 + \mu p) \geq \lambda_m \frac{\alpha}{4} \|f'(x_0)\| > 0$$

ce qui est absurde, donc x_0 est stationnaire.

Maintenant supposons que x_0 soit stationnaire; alors d'après (2.6) $\mu = 0$ et $x_0 + \mu p \equiv x_0$ d'où $f(x_0 + \mu p) = f(x_0)$ et x_0 est « désirable ».

Pour finir la démonstration du théorème il suffit maintenant d'utiliser le lemme 2.2. avec $\delta = \lambda_m \frac{\alpha}{4} \|f'(x_0)\|$.

3. APPLICATIONS

3.1. Méthodes quasi-Newton

Ce type de méthodes, dont fait partie celle de Davidon modifiée par Fletcher-Powell [3] est définie ainsi:

On dispose d'une suite de matrices carrées d'ordre n , $H_0, H_1, \dots, H_i, \dots$ toutes symétriques définies positives. La direction de déplacement est donnée, en posant $r_i = -[f'(x_i)]^T$, par

$$(3.1) \quad p_i = H_i r_i$$

et on a

$$(3.2) \quad x_{i+1} = x_i + \mu_i H_i r_i.$$

Supposant toujours $f \in C^1(\mathbf{R}^n)$ on va chercher dans quelles conditions l'hypothèse (2.5) est satisfaite, de façon à pouvoir appliquer le théorème 2.4. On a:

$$(3.3) \quad \frac{-f'(x_i) \cdot p_i}{\|f'(x_i)\| \|p_i\|} = \frac{r_i^T p_i}{\|r_i\| \|p_i\|} = \frac{r_i^T H_i r_i}{\|r_i\| \|H_i r_i\|}$$

Appelant $\beta_1^{(i)}$ et $\beta_n^{(i)}$ respectivement la plus grande et la plus petite valeur propre de H_i on a la minoration

$$(3.4) \quad \frac{r_i^T H_i r_i}{\|r_i\| \|H_i r_i\|} \geq \frac{\beta_n^{(i)}}{\beta_1^{(i)}},$$

où le rapport $\frac{\beta_n^{(i)}}{\beta_1^{(i)}}$ n'est autre que l'inverse du conditionnement de H_i . Par conséquent l'hypothèse (2.5) sera vérifiée si ce conditionnement est borné pour tout i . Il ne semble pas encore qu'on ait montré que cela était vrai de la méthode de Fletcher-Powell.

Par contre on peut considérer un cas particulier de méthode quasi-Newton. On doit, pour cela, faire les hypothèses:

$$(3.5) \quad \text{i) } f \in C^2(\mathbf{R}^n)$$

ii) la matrice J des dérivées partielles secondes de f vérifie

$$(3.6) \quad \lambda_n \|y\|^2 \leq y^T J(x) y \leq \lambda_1 \|y\|^2 \quad \forall x, y \in \mathbf{R}^n$$

$$\lambda_n > 0, \quad \lambda_1 > 0$$

ce qui entraîne que f est au moins strictement convexe.

On choisit pour matrices

$$(3.7) \quad H_i = J^{-1}(x_i).$$

Avec (3.1) et (3.2) la méthode ainsi définie diffère de la méthode bien connue de Newton du fait que dans (3.2) μ_i est généralement différent de 1.

Les inégalités (3.6) indiquent que le conditionnement des matrices H_i est borné et on peut appliquer le théorème 2.4.

Topkis et Veinott [4] ont obtenu des résultats analogues par un procédé différent.

3.2. Méthode du gradient conjugué

A partir de maintenant on suppose que f vérifie (3.5) et (3.6).

On va montrer que l'algorithme de Daniel [5], dont il a déjà démontré la convergence, vérifie (2.5). Cet algorithme n'est autre qu'une méthode de directions conjuguées [6] lorsque f est quadratique.

On part de x_0 arbitraire, $p_0 = r_0 = -[f'(x_0)]^T$ et on définit les itérés :

$$(3.8) \quad x_{i+1} = x_i + \mu_i p_i$$

$$(3.9) \quad p_{i+1} = r_{i+1} + \gamma_i p_i; \quad r_{i+1} = -[f'(x_{i+1})]^T$$

$$(3.10) \quad \gamma_i = -\frac{r_{i+1}^T J(x_{i+1}) p_i}{p_i^T J(x_{i+1}) p_i}$$

On voit facilement qu'en utilisant (2.6)

$$(3.11) \quad p_{i+1}^T r_{i+1} = r_{i+1}^T r_{i+1} + \gamma_i p_i^T r_{i+1} = \|r_{i+1}\|^2.$$

En outre d'après (3.9)

$$\|p_{i+1}\| \leq \|r_{i+1}\| + |\gamma_i| \|p_i\|$$

et d'après (3.10) et (3.6)

$$|\gamma_i| \leq \frac{\|r_{i+1}\| \|p_i\| \lambda_1}{\|p_i\|^2 \lambda_n}$$

d'où

$$\|p_{i+1}\| \leq \|r_{i+1}\| \left(1 + \frac{\lambda_1}{\lambda_n}\right)$$

et finalement

$$\frac{p_{i+1}^T r_{i+1}}{\|p_{i+1}\| \|r_{i+1}\|} \geq \frac{1}{1 + \lambda_1/\lambda_n} \quad \text{C.Q.F.D.}$$

3.3. Une modification de la méthode de Fletcher-Reeves

Dans l'algorithme de Daniel, le calcul de γ_i serait prohibitif puisqu'il faudrait connaître la matrice des dérivées secondes en chacun des itérés. C'est

pourquoi Fletcher-Reeves avaient proposé une expression plus simple; soit

$$(3.12) \quad \gamma_i = \frac{\|r_{i+1}\|^2}{\|r_i\|^2}$$

On va montrer tout d'abord que les expressions (3.10) et (3.12) sont identiques lorsque f est quadratique (1).

En effet, on a le développement de Taylor

$$(3.13) \quad -r_{i+1} = -r_i + \mu_i J p_i$$

où cette fois J est constante. Multipliant à gauche par p_{i-1} on a :

$$(3.14) \quad -p_{i-1}^T r_{i+1} = -p_{i-1}^T r_i + \mu_i p_{i-1}^T J p_i$$

Tout d'abord $p_{i-1}^T r_i$ est nul d'après (2.9).

Ensuite prenant (3.9) et (3.10) on a

$$(3.15) \quad p_i^T J p_{i+1} = p_i^T J r_{i+1} + \gamma_i p_i^T J p_i,$$

dont le second membre est nul d'après la définition (3.10) de γ_i . En remplaçant i par $i-1$ dans (3.15) on a également $p_{i-1}^T J p_i = 0$, ce qui, porté dans (3.14) entraîne que $p_{i-1}^T r_{i+1} = 0$. Le vecteur r_{i+1} orthogonal à p_{i-1} et p_i sera aussi orthogonal à r_i d'après (3.9) où on remplace i par $i-1$. Pour finir, utilisant (3.13) et (3.10) on a

$$(3.16) \quad \gamma_i = -\frac{r_{i+1}^T J p_i}{p_i^T J p_i} = -\frac{r_{i+1}^T (r_i - r_{i+1})}{p_i^T (r_i - r_{i+1})}$$

ce qui, du fait que $r_{i+1}^T r_i = 0$ ainsi que $p_i^T r_{i+1} = 0$, donne (3.12).

Revenons maintenant au cas où f n'est plus quadratique. On va écrire (3.12) d'une autre façon. Le développement de Taylor

$$(3.17) \quad -r_{i+1} = -r_i + \mu_i J_i p_i,$$

où

$$J_i = \int_0^1 J(x_i + \theta \mu_i p_i) d\theta;$$

donne

$$-p_i^T r_{i+1} = -p_i^T r_i + \mu_i p_i^T J_i p_i$$

et

$$\mu_i = \frac{p_i^T r_i}{p_i^T J_i p_i}$$

ou bien, d'après (3.11)

$$(3.18) \quad \mu_i = \frac{\|r_i\|^2}{p_i^T J_i p_i}$$

(1) Ceci est bien connu. Par exemple on peut en voir aussi la démonstration dans DURAND, *Résolution des équations algébriques*, t. II, Masson éd. (1962).

D'autre part en utilisant (3.17) on a

$$-r_{i+1}^T r_{i+1} = -r_{i+1}^T r_i + \mu_i r_{i+1}^T J_i' p_i$$

d'où la nouvelle forme de (3.12)

$$\gamma_i = \frac{r_{i+1}^T r_i - \mu_i r_{i+1}^T J_i' p_i}{\|r_i\|^2}$$

qui, en utilisant (3.18), donne finalement

$$(3.19) \quad \gamma_i = \frac{r_{i+1}^T r_i}{\|r_i\|^2} - \frac{r_{i+1}^T J_i' p_i}{p_i^T J_i' p_i}$$

On voit que (3.19) diffère assez peu de l'expression (3.10) utilisée par Daniel. Le deuxième terme du second membre n'en diffère que par la matrice J dont cette fois on prend la moyenne sur le segment $[x_i, x_{i+1}]$ alors que Daniel la calcule en x_{i+1} . En outre, comme on l'a démontré au début de ce paragraphe, le terme $r_{i+1}^T r_i$ est nul lorsque la fonction f est quadratique.

Par conséquent si, au lieu de définir γ_i comme Fletcher-Reeves, on choisissait

$$(3.20) \quad \gamma_i = \frac{\|r_{i+1}\|^2 - r_{i+1}^T r_i}{\|r_i\|^2} = - \frac{r_{i+1}^T J_i' p_i}{p_i^T J_i' p_i}$$

on pourrait facilement, comme on l'a fait de l'algorithme de Daniel, vérifier que (2.5) est satisfait.

Pour conclure remarquons que le minimum d'une fonction f quadratique pourra être obtenu en n itérations exactement avec le nouvel algorithme proposé [avec le choix (3.20) de γ_i] puisque dans ce cas $r_{i+1}^T r_i = 0$.

3.4. Convergence globale des algorithmes de directions conjuguées

En fait avec les hypothèses (3.5) et (3.6) qu'on a dû introduire pour parler des méthodes de directions conjuguées, on peut obtenir un résultat plus fort que le théorème (2.4).

En effet, on a le développement de Taylor

$$f(x) = f(x_0) + f'(x_0) \cdot (x - x_0) + \frac{1}{2} (x - x_0)^T J(x_0 + \theta(x - x_0))(x - x_0)$$

où

$$0 \leq \theta \leq 1.$$

d'où la minoration, en utilisant (3.6)

$$f(x) \geq f(x_0) - \|f'(x_0)\| \|x - x_0\| + \frac{1}{2} \lambda_n \|x - x_0\|^2.$$

Cela entraîne que l'ensemble des x tels que $f(x) \leq f(x_0)$ est borné et fermé dans \mathbf{R}^n donc compact. On peut donc effectivement extraire de la suite x_i de l'algorithme une sous-suite convergente. D'après le théorème 2.4, au point d'accumulation x^* on a $f'(x^*) = 0$ et puisqu'il n'existe qu'un seul point ayant cette propriété, toute la suite x_i converge en fait vers x^* le minimum unique de f (ce résultat est aussi donné par Daniel).

4. ESSAIS NUMERIQUES

On a comparé numériquement la méthode de Fletcher-Reeves et la méthode modifiée que nous proposons.

Il convient tout d'abord de faire une remarque. Fletcher et Reeves dans leur papier insistent sur la nécessité de briser la suite des directions conjuguées p_i à intervalles réguliers afin d'assurer la convergence de l'algorithme. Cela consiste à ce moment à choisir pour direction de déplacement $p_i = r_i$, donc à repartir suivant le gradient. Ils proposent, comme optimum, de briser toutes les $n + 1$ itérations. Effectivement en procédant ainsi on obtient la convergence dans la plupart des cas. Toutefois, lorsque f est convexe, on a pu observer que cela n'était pas nécessaire et c'est bien ce que montre la démonstration de convergence de notre méthode modifiée. En outre, couper la suite des directions conjuguées trop tôt peut être désastreux pour la convergence si la matrice J est très mal conditionnée au voisinage de la solution. Dans ce cas il est nécessaire de poursuivre les itérations beaucoup plus loin comme on peut le voir dans un article de Ginsburg [7] où f est quadratique.

Lorsque f n'est plus convexe, il est pratiquement nécessaire de briser la suite des directions conjuguées avec les deux méthodes. Toutefois on a vu des cas où, en utilisant la méthode modifiée, on obtenait une convergence raisonnable sans briser la suite.

En particulier avec le test proposé par A. R. Colville (New York Scientific Center I.B.M.):

minimiser

$$\begin{aligned} f(x) = & 100(\xi_2 - \xi_1^2)^2 + (1 - \xi_1)^2 + 90(\xi_4 - \xi_3^2)^2 \\ & + (1 - \xi_3)^2 + 10 \cdot 1[(\xi_2 - 1)^2 + (\xi_4 - 1)^2] \\ & + 19 \cdot 8(\xi_2 - 1)(\xi_4 - 1), \end{aligned}$$

où

$$x = \{ \xi_1, \xi_2, \xi_3, \xi_4 \}^T,$$

le minimum $f(x^*) = 0$ est obtenu pour

$$\xi_1^* = \xi_2^* = \xi_3^* = \xi_4^* = 1.$$

Le point de départ choisi était

$$\xi_1 = \xi_3 = -3; \quad \xi_2 = \xi_4 = -1.$$

Avec les deux méthodes, en utilisant une recherche très précise du minimum d'une fonction sur une droite, on a obtenu les résultats suivants :

MÉTHODE	FLETCHER-REEVES brisé à $n + 2$	FLETCHER-REEVES sans brisure	FLETCHER-REEVES modifié brisé à $3n$	FLETCHER-REEVES modifié sans brisure
Nombres d'itéra- tions	32	> 400	42	112

Dans tous les cas la valeur finale du coût était inférieure à 10^{-10} .

5. CONCLUSION

En somme, on voit que, en utilisant un résultat général de convergence, comme le lemme 2.2., on peut non seulement démontrer la convergence de nombreux algorithmes bien connus mais aussi voir comment modifier certains algorithmes heuristiques pour obtenir des méthodes convergentes.

ACKNOWLEDGMENT

This research was supported in part by the national Aeronotics and Space Administration under Grant No Ns G 354 (supp. 5).

REFERENCES

- [1] R. FLETCHER et C. M. REEVES, « Function minimization by conjugate gradients », *The Computer Journal*, pp. 149-154 (1964).
- [2] E. POLAK, « On Primal and Dual Methods for solving discrete optimal control Problems », Second Int. Conf. on Computing methods in optimization problems, San Remo, Italy, Sept. 1968.
- [3] R. FLETCHER et M. J. D. POWELL, « A rapidly convergent descent method for minimization », *The Computer Journal*, vol. 6, p. 163 (1963).
- [4] D. M. TOPKIS et A. F. VEINOTT (Jr), « On the convergence of some feasible direction algorithms for non-linear programming », *Siam. J. on Control*, vol. 5, n° 2, p. 268 (1967).
- [5] J. W. DANIEL, « The conjugate gradient method for linear and non linear operator equations », *Siam J. Num. Anal.*, vol. 4, n° 1 (1967).
- [6] M. R. HESTENES et E. STIEFEL, « Methods of conjugate gradients for solving linear systems », *J. Res. N.B.S.*, vol. 49, p. 409 (1952).
- [7] T. GINSBURG, « The conjugate gradient method », *Numerische Mathematik Band 5*, Heft 2, p. 191 (1963).