

# IFT 3515

## Fonctions à plusieurs variables

### Optimisation sans contraintes

### Méthode de Steihaug-Toint

Fabian Bastin  
DIRO  
Université de Montréal

## Motivations

Considérons le problème

$$\min f(x)$$

où  $f \in C^2$ .

Nous avons vu que l'algorithme de région de confiance permet de résoudre ce problème et tout point limite est un point critique au premier ordre si la minimisation approximative du modèle donne une fraction suffisante de la réduction atteinte au point de Cauchy.

Cependant, la convergence est lente si seul le point de Cauchy est utilisé.

## Motivations

Typiquement, un modèle quadratique est utilisé pour le sous-problème de région de confiance. La méthode du gradient conjugué est efficace pour minimiser approximativement une fonction quadratique strictement convexe. Peut-on l'adapter si le modèle n'est plus strictement convexe?

Si le modèle n'est pas convexe, la région de confiance devrait prévenir la génération d'une séquence de points allant à l'infini. Par contre, si le modèle est strictement convexe et atteint son minimum à l'intérieur de la région de confiance, celle-ci ne devrait pas interférer avec la méthode du gradient conjugué.

## Sous-problème

Nous considérons le sous-problème

$$\min_s q(s) := g^T s + \frac{1}{2} s^T H s$$

sous la contrainte  $\|s\|_M \leq \Delta$ .

Le sous-problème est défini en prenant

$$g = \nabla f(x_k)$$

$$H = \nabla^2 f(x_k)$$

$$M = R^T R \quad (\text{préconditionneur})$$

## Préconditionnement et changement d'échelle

La condition

$$\|s\|_M \leq \Delta$$

peut se réécrire

$$\langle s, Ms \rangle \leq \Delta^2$$

ou encore

$$\langle Rs, Rs \rangle \leq \Delta^2.$$

Le préconditionnement revient à effectuer un changement d'échelle dans les variables.

## Application du gradient conjugué

Supposons que nous appliquions l'algorithme du gradient conjugué pour minimiser  $q$ . Plusieurs cas de figure peuvent se produire.

Si à chaque itération du gradient conjugué  $\langle d_k, Hd_k \rangle > 0$  et tous les itérés  $s_k$  restent à l'intérieur de la région de confiance, nous obtenons un problème d'optimisation convexe sans contrainte.

Il peut toutefois arriver qu'à une certaine itération  $k$ ,  $\langle d_k, Hd_k \rangle \leq 0$ . Dans ce cas, le problème de recherche linéaire

$$\min_{\alpha} q(s_k + \alpha d_k)$$

donnerait  $\alpha_k = +\infty$ . Mais si nous ajoutons la contrainte de région de confiance, alors  $\alpha_k$  sera choisi de manière à atteindre la frontière de la région de confiance. Dans ce cas, nous cherchons la racine positive du problème

$$\|s_k + \alpha_k d_k\|_M \leq \Delta_k.$$

## Application du gradient conjugué

La troisième possibilité est que  $s_k$  soit en dehors de la région de confiance à l'itération  $k$ .

Il est concevable que les itérés suivants puissent revenir dans la région de confiance. Toutefois, ce n'est pas le cas.

### Théorème

*Supposons que nous appliquions l'algorithme du gradient conjugué préconditionné pour minimiser  $q(s)$ , en partant de  $s_0 = 0$ , et que  $\langle d_i, Hd_i \rangle > 0$  pour  $0 \leq i \leq k$ . Les itérés  $s_j$  satisfont alors les inégalités*

$$\|s_j\|_M \leq \|s_{j+1}\|_M,$$

*pour  $0 \leq j \leq k - 1$ .*

## Preuve

Montrons tout d'abors que

$$\langle d_j, M d_i \rangle = \frac{g_j^T v_j}{g_i^T v_i} \langle d_i, M d_i \rangle > 0.$$

pour tout  $0 \leq i \leq j \leq k$ . C'est trivial pour  $i = j$ . Supposons que cela soit aussi vrai pour  $i \leq l$ . De l'algorithme du gradient conjugué, nous avons

$$d_{l+1} = -v_{l+1} + \frac{g_{l+1}^T v_{l+1}}{g_l^T v_l} d_l.$$

Dès lors,

$$\langle d_{l+1}, M d_i \rangle = -\langle v_{l+1}, M d_i \rangle + \frac{g_{l+1}^T v_{l+1}}{g_l^T v_l} \langle d_l, M d_i \rangle$$

## Preuve

Par hypothèse de récurrence

$$\begin{aligned}\langle d_{l+1}, Md_i \rangle &= -\langle v_{l+1}, Md_i \rangle + \frac{g_{l+1}^T v_{l+1}}{g_l^T v_l} \frac{g_l^T v_l}{g_i^T v_i} \langle d_i, Md_i \rangle \\ &= -\langle v_{l+1}, Md_i \rangle + \frac{g_{l+1}^T v_{l+1}}{g_i^T v_i} \langle d_i, Md_i \rangle\end{aligned}$$

Or, la propriété de conjugaison indique que

$$\langle v_{l+1}, Md_i \rangle = 0$$

et donc

$$\langle d_{l+1}, Md_i \rangle = \frac{g_{l+1}^T v_{l+1}}{g_i^T v_i} \langle d_i, Md_i \rangle$$

De plus, comme  $M$  est définie positive, pour  $0 \leq i \leq k$ ,

$$g_i^T v_i = g_i^T M g_i > 0.$$

## Preuve

Par conséquent, nous avons bien

$$\langle d_{l+1}, M d_i \rangle > 0.$$

D'autre part, l'algorithme du gradient conjugué donne

$$s_j = s_0 + \sum_{i=0}^{j-1} \alpha_i d_i = \sum_{i=0}^{j-1} \alpha_i d_i.$$

Dès lors

$$\begin{aligned}\langle s_j, M d_j \rangle &= \left\langle \sum_{i=0}^{j-1} \alpha_i d_i, M d_j \right\rangle \\ &= \sum_{i=0}^{j-1} \alpha_i \langle d_i, M d_j \rangle > 0\end{aligned}$$

## Preuve

Par conséquent,

$$\begin{aligned}\|s_{j+1}\|_M^2 &= \langle s_{j+1}, Ms_{j+1} \rangle \\&= \langle s_j + \alpha_j d_j, M(s_j + \alpha_j d_j) \rangle \\&= \langle s_j, Ms_j \rangle + 2\alpha_j \langle s_j, Md_j \rangle + \alpha_j^2 \langle d_j, Md_j \rangle \\&> \langle s_j, Ms_j \rangle = \|s_j\|_M^2.\end{aligned}$$

## Discussion

Dès lors, aussi longtemps qu'une courbure positive, i.e.  $H$  est définie positif, est rencontrée dans la méthode du gradient conjugué préconditionné, la norme-M des itérés est strictement croissante, pourvu que nous partions de  $s_0 = 0$ . Ainsi

$$\|s_j\|_M \leq \|\arg \min_{s \in \mathbb{R}^n} q(s)\|_M,$$

où l'inégalité est stricte excepté à l'itération finale.

En particulier, quand  $H$  est définie positif, et que  $s^*$  se trouve hors de la région de confiance, la solution du sous-problème de région de confiance doit se trouver sur la frontière de la région de confiance. Dès lors, une stratégie adéquate est de revenir sur la frontière lorsqu'un itéré du gradient conjugué arrive en dehors de la région de confiance.

## GC tronqué de Steihaug-Toint

Étant donné  $x_0$ , posons  $g_0 = \nabla f(x_0)$ ,  $v_0 = M^{-1}g_0$  et  $d_0 = -v_0$ . Pour  $k = 0, 1, 2, \dots$ , jusqu'à convergence, faire

- Set  $\kappa_k = \langle d_k, Hd_k \rangle$ .
- Si  $\kappa_k \leq 0$ , calculer  $\sigma_k$  comme la racine positive de  $\|s_k + \sigma_k d_k\|_M = \Delta$ , et poser

$$s_{k+1} = s_k + \sigma_k d_k.$$

Arrêt.

- Sinon, calculer

$$\alpha_k = \frac{\langle g_k, v_k \rangle}{\kappa_k}$$

- Si  $\|s_k + \alpha_k d_k\|_M \geq \Delta$ , calculer  $\sigma_k$  comme la racine positive de  $\|s_k + \sigma_k d_k\|_M = \Delta$ , et poser

$$s_{k+1} = s_k + \sigma_k d_k.$$

Arrêt

- Sinon

## GC tronqué de Steihaug-Toint

Calculer

$$s_{k+1} = s_k + \alpha_k d_k$$

$$g_{k+1} = g_k + \alpha_k H d_k$$

$$v_{k+1} = M^{-1} g_{k+1}$$

$$\beta_k = \frac{\langle g_{k+1}, v_{k+1} \rangle}{\langle g_k, v_k \rangle}$$

$$d_{k+1} = -v_{k+1} + \beta_k d_k$$

## Convergence

Aussi longtemps que le nombre de conditionnement de  $M$  reste borné sur la séquence de sous-problèmes approximativement résolus par l'algorithme de région de confiance, alors n'importe quel itéré généré par l'algorithme de gradient conjugué tronqué est suffisant pour assurer la convergence vers un point critique au premier ordre.

Ceci résulte du fait que le premier itéré généré par l'algorithme de gradient conjugué tronqué est le point de Cauchy pour le modèle, et n'importe quel itéré généré par la suite donne une valeur plus petite pour le modèle.

## Considérations pratiques

À chaque étape de la méthode de Steihaug-Toint, nous devons calculer  $\|s_k + \alpha_k d_k\|_M$ .

Ce n'est pas un problème si  $M$  est disponible, mais ce peut l'être si tout ce qui est disponible est une procédure qui renvoie  $M^{-1}v$ , pour un  $v$  donné, de sorte que  $M$  est non disponible.

Heureusement, ce n'est pas un inconvénient majeur comme il est possible de calculer  $\|s_k + \alpha_k d_k\|_M$  à partir de l'information disponible..

Observons tout d'abord que

$$\|s_k + \alpha_k d_k\|_M^2 = \|s_k\|_M^2 + 2\alpha_k \langle s_k, M d_k \rangle + \alpha^2 \|d_k\|_M^2.$$

## Considérations pratiques

La racine carrée positive de  $\|s_k + \sigma_k d_k\|_M = \Delta$  s'obtient en développant l'égalité

$$\|s_k + \sigma_k d_k\|_M^2 = \Delta^2$$

Nous pouvons la réécrire comme

$$\begin{aligned} 0 &= \|s_k + \sigma_k d_k\|_M^2 - \Delta^2 \\ &= \|s_k\|_M^2 - \Delta^2 + 2\sigma_k \langle s_k, M d_k \rangle + \sigma^2 \|d_k\|_M^2 \end{aligned}$$

Les racines carrées sont

$$\sigma_k = \pm \frac{\langle s_k, M d_k \rangle - \sqrt{\langle s_k, M d_k \rangle^2 - \|d_k\|_M^2 (\|s_k\|_M^2 - \Delta^2)}}{\|d_k\|_M^2}$$

et la racine positive est donnée par

$$\sigma_k = \frac{-\langle s_k, M d_k \rangle + \sqrt{\langle s_k, M d_k \rangle^2 - \|d_k\|_M^2 (\|s_k\|_M^2 - \Delta^2)}}{\|d_k\|_M^2}$$

## Considérations pratiques

Dès lors, nous pouvons calculer  $\|s_{k+1}\|_M^2$  à partir de  $\|s_k\|_M^2$  aussi longtemps que nous connaissons déjà  $\langle s_k, M d_k \rangle$  et  $\|d_k\|_M^2$ .

D'une part, en exploitant la propriété de conjugaison,

$$\begin{aligned}\|d_k\|_M^2 &= \langle d_k, M d_k \rangle \\&= \langle -v_k + \beta_{k-1} d_{k-1}, M(-v_k + \beta_{k-1} d_{k-1}) \rangle \\&= \langle -v_k, -M v_k \rangle + \langle \beta_{k-1} d_{k-1}, M \beta_{k-1} d_{k-1} \rangle \\&= \langle v_k, M M^{-1} g_k \rangle + \beta_{k-1}^2 \langle d_{k-1}, M d_{k-1} \rangle \\&= g_k^T v_k + \beta_{k-1}^2 \|d_{k-1}\|_M^2\end{aligned}$$

## Considérations pratiques

D'autre part, comme

$$s_k = \sum_{i=0}^{k-1} \alpha_i d_i,$$

nous avons

$$\begin{aligned}\langle s_k, M d_k \rangle &= \langle s_{k-1} + \alpha_{k-1} d_{k-1}, M(-v_k + \beta_{k-1} d_{k-1}) \rangle \\&= -\langle s_{k-1}, M v_k \rangle + \beta_{k-1} \langle s_{k-1}, M d_{k-1} \rangle \\&\quad - \alpha_{k-1} \langle d_{k-1}, M v_k \rangle + \alpha_{k-1} \beta_{k-1} \langle d_{k-1}, M d_{k-1} \rangle \\&= \beta_{k-1} (\langle s_{k-1}, M d_{k-1} \rangle + \alpha_{k-1} \|d_{k-1}\|_M^2).\end{aligned}$$

## Terminaison

Il n'est pas nécessaire de minimiser le modèle avec une grande précision. On pourra s'arrêter lorsqu'une réduction suffisante du gradient est atteinte

$$\|g_k\| \leq \|g_0\| \min\{\chi, \|g_0\|^\theta\}$$

avec  $\chi < 1$ ,  $\theta \geq 0$  ou  $k > k_{\max} \geq 0$ .

Des valeurs usuelles sont  $\chi = 0.1$ ,  $\theta = 0.5$  et  $k_{\max} = n$ .