

Data Mining

TP8 - Optimisation sans contraintes

Rémi MUSSARD - Thomas ROBERT

Le but du TP est d'étudier différentes méthodes permettant de converger vers le minimum d'une fonction convexe de plusieurs variables.

Pour ces méthodes, on utilise le gradient ou la matrice Hessienne afin de déterminer la direction de descente permettant de faire décroître la fonction de coût. On avance dans cette direction d'une valeur définie par un pas qui peut être fixe ou variable. Enfin, on itère ce processus tant que le gradient est supérieur à un seuil fixé à 10^{-3} .

1 Ce qu'on veut trouver

Afin de calculer le gradient et la Hessienne, on calcule les dérivées partielles d'ordre 1 et 2 par rapport à θ_1 et θ_2 . On obtient :

$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial \theta_1} = -2(1 - \theta_1) - 400\theta_1(\theta_2 - \theta_1^2)^2 \\ \frac{\partial J}{\partial \theta_2} = 200(\theta_2 - \theta_1^2)^2 \end{bmatrix}$$
$$H = \begin{bmatrix} \frac{\partial^2 J}{\partial \theta_1^2} = 2 - 400(\theta_2 - \theta_1^2)^2 + 800\theta_1^2 & \frac{\partial^2 J}{\partial \theta_1 \partial \theta_2} = -400\theta_1 \\ \frac{\partial^2 J}{\partial \theta_1 \partial \theta_2} = -400\theta_1 & \frac{\partial^2 J}{\partial \theta_2^2} = 200 \end{bmatrix}$$

On peut également calculer analytiquement l'extremum de cette fonction, il suffit d'annuler les dérivées, ce qui permet d'obtenir un extremum en $\theta^* = (1, 1)^\top$.

Le calcul de la Hessienne permet de voir qu'il s'agit bien d'un minimum puisqu'elle est définie positive :

$$H(\theta^*) = \begin{bmatrix} 802 & -400 \\ -400 & 200 \end{bmatrix}$$

2 Préparation

```
1 clear all
2 close all
3 clc
4
5 % parametres du probleme
6 a = [1; 3];
7 b = [1; -3];
8 c = [-1; 0];
9
10 % Create a grid of x and y points
11 n = 100;
12 [X, Y] = meshgrid(linspace(-1.25, 2.5, n), linspace(-1.75, 2, n));
13 ptx = reshape(X, n*n, 1);
14 pty = reshape(Y, n*n, 1);
15 pt = [ptx pty];
```

```

16
17 % Define the function J = f(\theta)
18 Jmat = (1-pt(:,1)).^2 + 100*(pt(:, 2) - pt(:, 1).^2).^2;
19
20 % solution initiale
21 theta0 = [-1; 0];

```

3 Fonction init_fig

Afin d'afficher une figure avec les lignes de niveaux, on utilisera la fonction init_fig.

```

1 function init_fig(theta0, Jmat, n, X, Y)
2
3 % Create the surface plot using contour command
4 figure;
5 set(gcf, 'PaperUnits', 'points');
6 set(gcf, 'PaperPosition', [0 0 750*1.7 500*1.7]);
7 set(gcf, 'Position', [0 0 750*1.7 500*1.7]);
8 contour(X, Y, reshape(Jmat, n, n), [40:-5:0 1 0], 'linewidth', 1.5);
9 colorbar
10 axis tight
11
12 % trace de theta0
13 hold on
14 h = plot(theta0(1,:), theta0(2,:), 'ro');
15 set(h, 'MarkerSize', 8, 'markerfacecolor', 'r');
16 text(theta0(1), theta0(2)+0.025, '\theta_0', 'fontsize', 15)

```

4 Première version avec α constant

On réalise un premier code avec un pas α constant. Ce pas est défini par des essais successifs pour trouver un pas adapté au problème, c'est à dire convergeant à une vitesse satisfaisante : ni trop rapide, ni trop lent.

Notons qu'en réalité, le pas d'avancement à chaque itération n'est pas constant car on ne normalise pas le gradient. Ainsi, le pas d'avancement réel dépend de la valeur du gradient, et est donc d'autant plus grand que la fonction augmente rapidement selon la direction de descente.

On remarque que la fonction présente une zone où le coût ne varie quasiment pas, donnant un gradient de valeur faible, et nous faisant donc faire de très petits pas vers le minimum, rendant la convergence très lente.

```

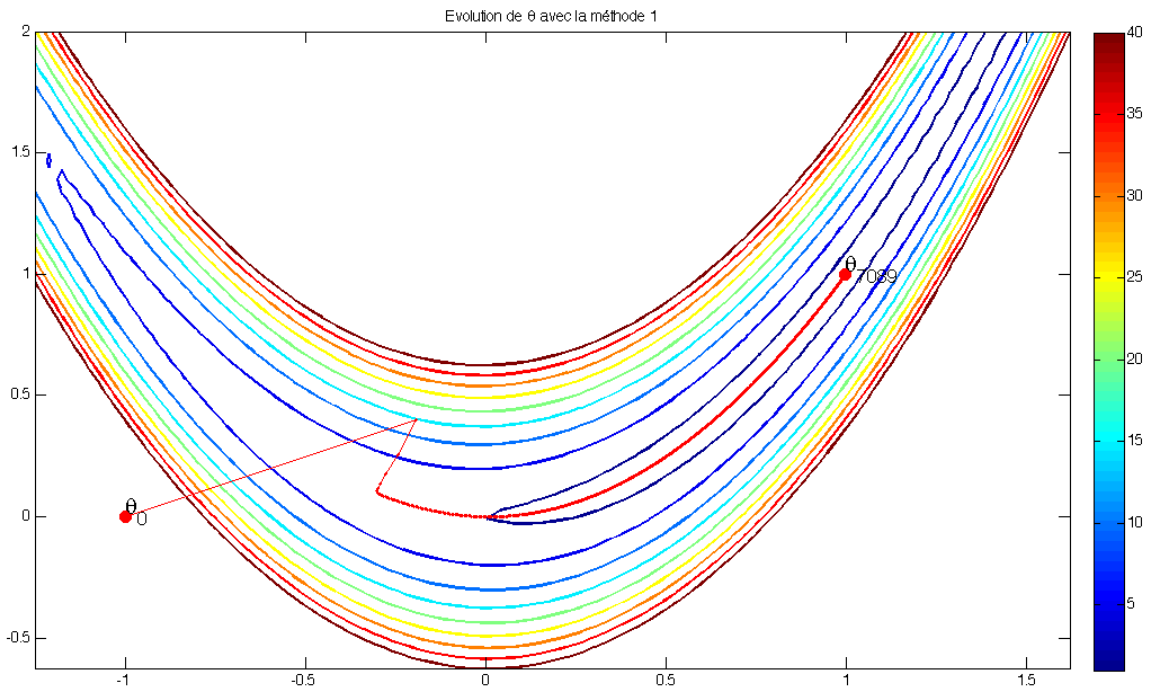
1 % pas alpha fixé
2 alpha = 0.002;
3
4 % variables pour la boucle et initialisation
5 Jlist = [];
6 grad = mongradient(theta0);
7 theta_old = theta0;
8 theta = theta0;
9 i = 1;
10
11 % initialiser la figure
12 close all;
13 init_fig(theta0, Jmat, n, X, Y);
14
15 % tant qu'on a pas convergé, on itère
16 while norm(grad) > 1e-3 && i < 15000
17
18 % calculs du nouveau theta
19 grad = mongradient(theta);
20 direction = -grad;
21 theta_old = theta;
22 theta = theta + alpha * direction;
23
24 % trace du theta courant

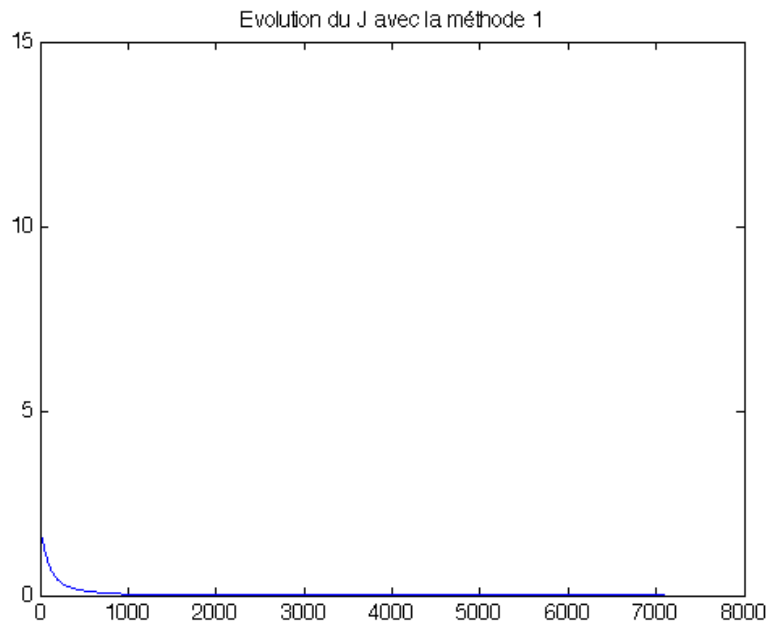
```

```

25     h = plot([theta_old(1) theta(1)], [theta_old(2) theta(2)], '-ro');
26     set(h, 'MarkerSize', 2, 'markerfacecolor', 'r');
27
28     % calcul de J
29     Jlist = [Jlist moncritere(theta)];
30     i = i + 1;
31 end
32
33 % theta final
34 h = plot(theta(1,:), theta(2,:), 'ro');
35 set(h, 'MarkerSize', 8, 'markerfacecolor', 'r');
36 text(theta(1,1), theta(2,1)+0.025, ['\theta_{', int2str(i-1), '}'], 'fontsize', 15)
37 title('Evolution de \theta avec la méthode 1');
38
39 % affichage évolution J
40 figure;
41 plot(Jlist);
42 title('Evolution du J avec la méthode 1');

```





5 Deuxième version avec α variable

Cette fois, on décide d'appliquer une règle permettant de faire varier α afin de converger plus vite.

A chaque itération, si le α que l'on a nous permet de converger, on le multiplie par 1,15 afin d'essayer de converger plus vite à l'itération suivante.

Sinon, si le α que l'on avait fait augmenter le coût, alors on annule les calculs réalisés et on divise α par 2.

On voit que la méthode met toujours longtemps à converger car la direction de descente ne suit pas vraiment une direction optimale qui serait au centre de la ligne de niveau, c'est à dire la direction curviligne qui suivrait le « creux » de la fonction.

A la place, on oscille autour de cette direction optimale en suivant des directions qui ne sont pas optimales (souvent environ à 45° entre la direction optimale et sa perpendiculaire).

Le fait que cette direction ne sont pas optimale fait qu'il est impossible d'avancer beaucoup dans cette direction sans faire réaugmenter le coût J . Ainsi, le pas ne peut pas vraiment augmenter.

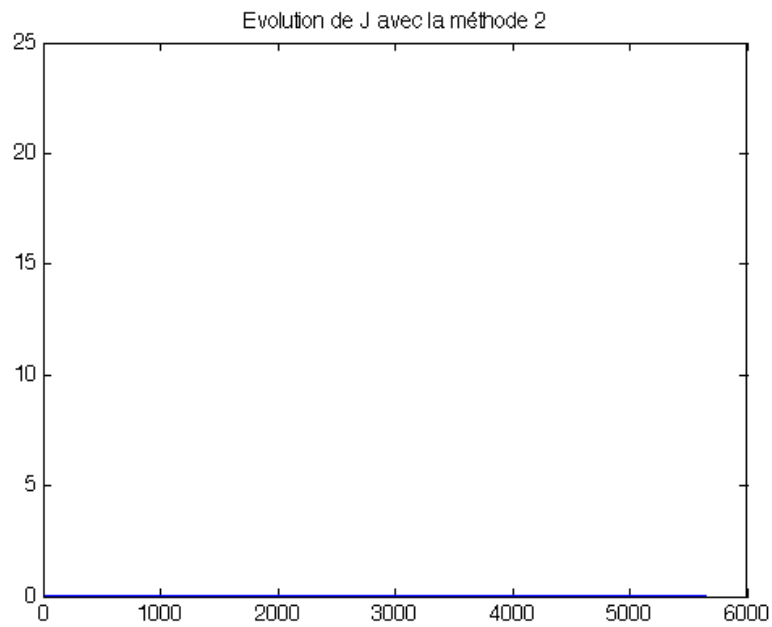
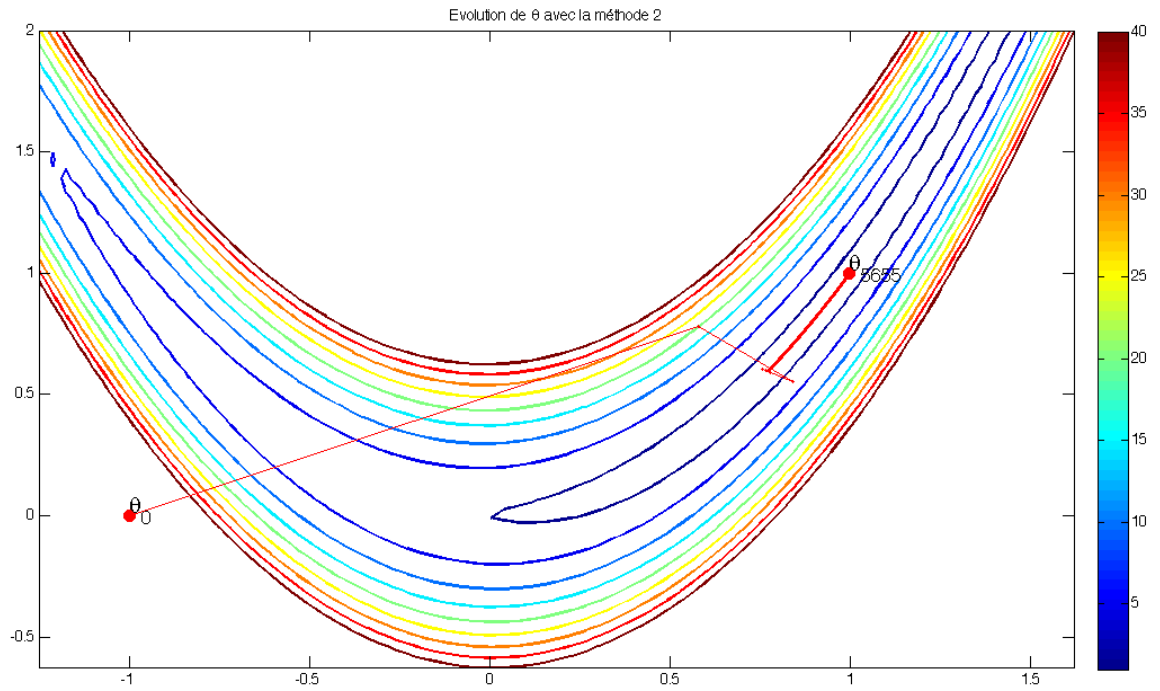
Ce problème de direction vient du fait que le gradient nous donne une approximation de premier ordre. On essaye descend donc forcément en ligne droite ce qui n'est pas optimal.

```

1 % initialisation des variables
2 alpha = 1;
3 grad = mongradient(theta0);
4 Jlist = [];
5 J = moncritere(theta0);
6 Jprec = J+1e-3;
7 theta = theta0;
8 theta_old = theta0;
9 i = 1;
10 j = 1;
11
12 % initialiser la figure
13 close all;
14 init_fig(theta0, Jmat, n, X, Y);
15
16 % tant qu'on a pas convergé, on itère

```

```
17 while norm(grad) > 1e-3 && i < 15000 && j < 10000
18
19 % calculs
20 grad = mongradient(theta); % calcul du gradient
21 direction = -grad; % direction de descente
22 theta_new = theta + alpha * direction; % MAJ de theta
23 J_new = moncritere(theta_new); % calcul de J
24
25 % si on améliore J avec le calcul réalisé
26 if(J - J_new > 0)
27
28     % augmentation de alpha pour l'itération suivante
29     alpha = alpha*1.3;
30
31     % enregistrement de ce qui a été fait
32     Jprec = J;
33     J = J_new;
34     Jlist = [Jlist J];
35     theta_old = theta;
36     theta = theta_new;
37
38     % affichage theta
39     h = plot([theta_old(1) theta(1)], [theta_old(2) theta(2)], '-ro');
40     set(h, 'MarkerSize', 2, 'markerfacecolor', 'r');
41
42     j = j + 1;
43
44     % sinon si on a augmenté J, on enregistre rien et on diminue alpha
45     else
46         alpha = alpha/2;
47     end
48
49     i = i + 1;
50 end
51
52 % theta final
53 h = plot(theta(1,:), theta(2,:), 'ro');
54 set(h, 'MarkerSize', 8, 'markerfacecolor', 'r');
55 text(theta(1,1), theta(2,1)+0.025, ['\theta_{', int2str(j-1), '}'], 'fontsize', 15)
56 title('Evolution de \theta avec la méthode 2');
57
58 % affichage évolution J
59 figure;
60 plot(Jlist);
61 title('Evolution de J avec la méthode 2');
```



6 Troisième version avec la matrice Hessienne

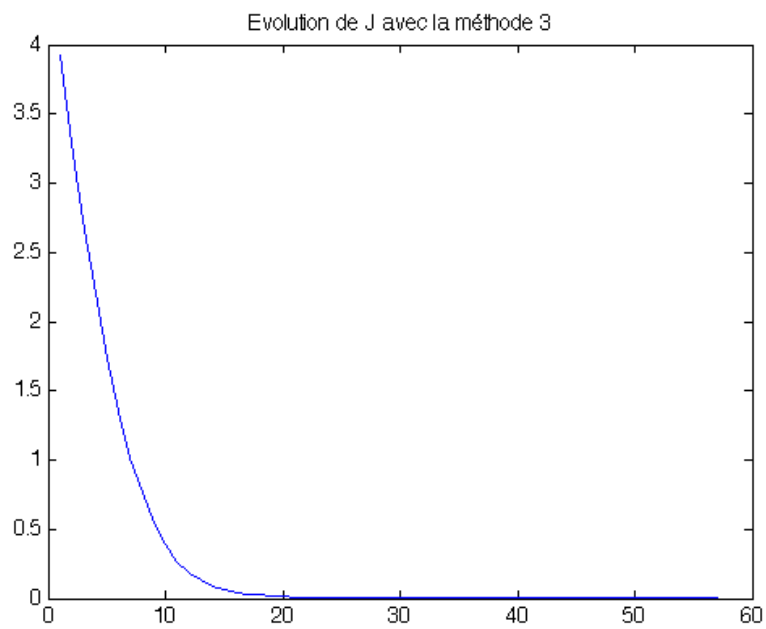
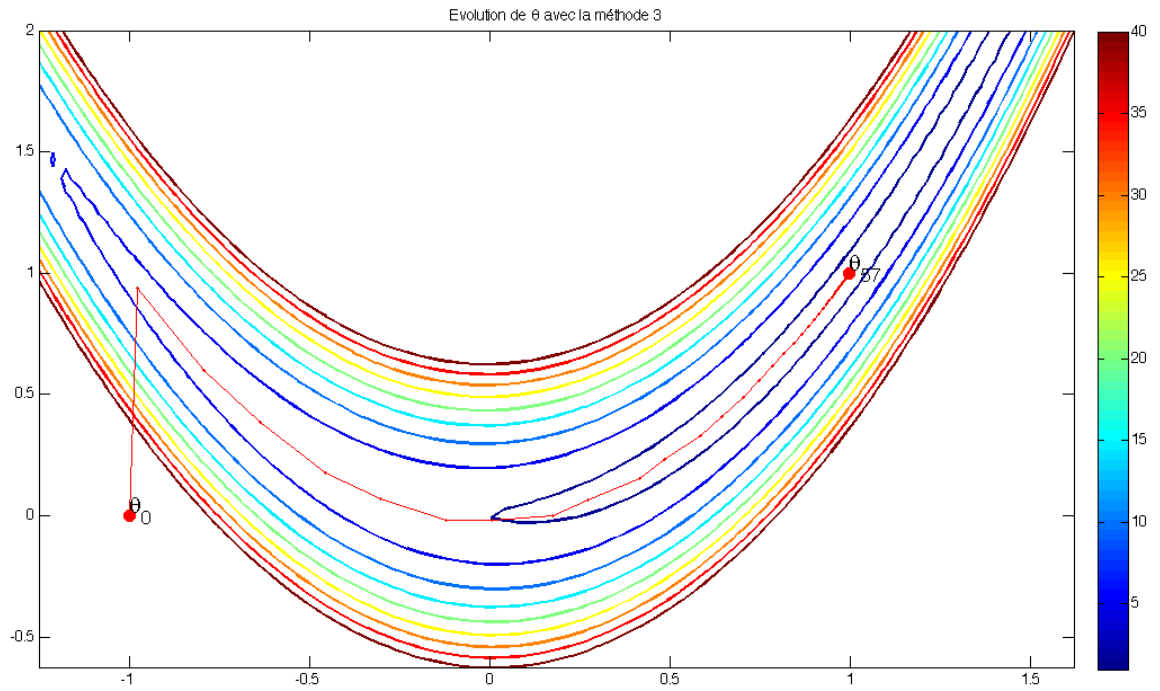
Cette fois, on calcule la matrice Hessienne afin d'augmenter la vitesse de convergence. On calcule donc (via la matrice Hessienne) les dérivées secondes de la fonction étudiée, permettant de converger beaucoup plus rapidement.

Le problème est que notre matrice Hessienne n'est pas définie positive, ou du moins pas toujours bien conditionnée. On utilise donc la méthode de *quasi-Newton* en utilisant au lieu de la matrice H une matrice B proche de H mais définie positive. Dans notre cas, on se contente de reconditionner la matrice H en avant $\lambda = 3$ sur la diagonale.

On remarque qu'avec cette méthode, on converge beaucoup plus rapidement qu'avec les méthodes précédentes.

```

1 % initialisation des variables
2 lambda = 3;
3 Jlist = [];
4 J = moncritere(theta0);
5 Jprec = J+1e-3;
6 grad = mongradient(theta0);
7 theta = theta0;
8 i = 1;
9
10 % initialiser la figure
11 close all;
12 init_fig(theta0, Jmat, n, X, Y);
13
14 % tant qu'on a pas convergé, on itère
15 while norm(grad) > 1e-3 && i < 3000
16 %for i=1:2
17     % calculs
18     grad = mongradient(theta);           % calcul du gradient
19     H = monHessien(theta);               % calcul de la matrice Hessienne
20     B = H + lambda * eye(size(H));      % calcul d'une matrice proche de H def pos
21     direction = -B\grad;                 % direction de descente
22     theta_old = theta;                   % sauvegarde ancien theta
23     theta = theta + direction;           % MAJ theta
24
25     % calcul de J
26     Jprec = J;
27     J = moncritere(theta);
28     Jlist = [Jlist J];
29
30     % trace du theta courant
31     h = plot([theta_old(1) theta(1)], [theta_old(2) theta(2)], '-ro');
32     set(h, 'MarkerSize', 2, 'markerfacecolor', 'r');
33
34     % inc i
35     i = i + 1;
36 end
37
38 % theta final
39 h = plot(theta(1,:), theta(2,:), 'ro');
40 set(h, 'MarkerSize', 8, 'markerfacecolor', 'r');
41 text(theta(1,1), theta(2,1)+0.025, ['\theta_{', int2str(i-1), '}'], 'fontsize', 15)
42 title('Evolution de \theta avec la méthode 3');
43
44 % affichage évolution J
45 figure;
46 plot(Jlist);
47 title('Evolution de J avec la méthode 3');
```



7 Conclusion

Après avoir épuré le code afin de ne conserver que le code permettant de calculer le θ final pour chacune des méthodes (suppression des affichages, de la sauvegarde des J, etc.), on mesure combien de temps met chacune des méthodes pour converger.

On trouve les résultats suivants :

- Méthode α fixé : 0.29 secondes
- Méthode α variable : 0.48 secondes
- Méthode quasi-Newton : 0.023 secondes

On voit que dans ce cas, la méthode quasi-Newton est particulièrement efficace comparé aux autres. Ce phénomène aurait probablement tendance à s'atténuer avec d'autres fonctions, ou si nous étions dans un espace de plus grande dimension que \mathbb{R}^2 , en effet, la méthode de Newton a une complexité en $O(n^3)$ rendant l'algorithme très gourmand et peu performant en grande dimension.