

# Projet Optimisation A4

Léon Thomir et Hugo Yeremian

March 2023

## 1 Partie Théorique

On considère la fonction  $f$  définie sur  $R^2$  par

$$f(x, y) = x^4 + y^4 - 2(x - y)^2 \quad (1)$$

1. Montrer qu'il existe  $(\alpha, \beta) \in R_+^2$  (et les déterminer) tels que  $f(x, y) \geq \alpha\|(x, y)\|^2 - \beta$  pour tous  $(x, y) \in R^2$ , où la notation  $\|\cdot\|$  désigne la norme euclidienne sur  $R^2$ .

En déduire que le problème  $\mathcal{P}$  donné par

$$\inf_{(x, y) \in R^2} f(x, y) \quad (2)$$

possède au moins une solution.

Comme fonction polynomiale,  $f \in C^\infty$ . On a donc  $\forall (x, y) \in R^2$

$$\begin{aligned} f(x, y) &= x^4 + y^4 - 2x^2 - 2y^2 + 4xy \\ &\geq x^4 + y^4 - 4x^2 - 4y^2 \end{aligned} \quad (3)$$

car  $\forall (x, y) \in R^2, (x + y)^2 \geq 0 \Leftrightarrow xy \geq -\frac{x^2 + y^2}{2}$ .

On rappelle que  $\forall (K, z) \in R^2, (K + z)^2 K^4 + z^4 - 2zK^2 \geq 0$ , d'où

$$f(x, y) \geq 2(2z - 4)x^2 + 2(z - 4)y^2 - 2z^4 \quad (4)$$

En prenant  $z = 4$ , on a

$$\begin{aligned} f(x, y) &\geq 4(x^2 + y^2) - 512 \\ &\geq 4\|(x, y)\|^2 - 512 \end{aligned} \quad (5)$$

donc avec  $\alpha = 4$  et  $\beta = 512$ , on a  $\forall (x, y) \in R^2, f(x, y) \geq \alpha\|(x, y)\|^2 - \beta$ .

En utilisant  $\alpha\|(x, y)\|^2 - \beta \xrightarrow{\|(x, y)\|^2 \rightarrow +\infty} +\infty$ , et  $f$  coercive sur  $R^2$  (car  $R^2$  est fermée et de dimension finie), le problème  $\mathcal{P}$  possède au moins une solution.

2. La fonction  $f$  est-elle convexe sur  $R^2$  ?

Comme fonction polynomiale,  $f \in C^2$  sur  $R^2$ , on peut donc étudier les dérivées de  $f$ .  $\forall (x, y) \in R^2$ ,

$$\frac{\partial f(x, y)}{\partial x} = 4x^3 - 4x + 4y = 4x^3 + 4(y - x)$$

$$\frac{\partial f(x, y)}{\partial y} = 4y^3 + 4(y - x)$$

$$\frac{\partial^2 f(x, y)}{\partial x^2} = 12x^2 - 4$$

$$\frac{\partial^2 f(x, y)}{\partial y^2} = 12y^2 + 4$$

$$\frac{\partial^2 f(x, y)}{\partial x \partial y} = 4$$

On définit la Hessienne  $H$  de  $f$ :

$$H(x, y) = \begin{bmatrix} 12x^2 - 4 & 4 \\ 4 & 12y^2 + 4 \end{bmatrix} \Rightarrow H(0, 0) = 4 \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix} \quad (6)$$

d'où  $|H(0,0)| = 0$ , donc la Hessienne n'est pas semi-définie positive en tout point. Ainsi,  $f$  n'est pas convexe en tout point. **3.** Déterminer les points critiques de  $f$  et préciser leur nature. Résoudre alors le problème dans  $\mathcal{P}$ .  
On a  $\forall (x,y) \in R^2$ ,

$$\begin{aligned}\nabla f(x,y) = \begin{bmatrix} 0 \\ 0 \end{bmatrix} &\Leftrightarrow \begin{cases} x^3 - (x-y) = 0 \\ y^3 + (x-y) = 0 \end{cases} \\ &\Leftrightarrow \begin{cases} x^3 + y^3 = 0 \\ y^3 + (x-y) = 0 \end{cases} \\ &\Leftrightarrow \begin{cases} x = -y \\ x^3 - 2x = 0 \end{cases} \\ &\Leftrightarrow \begin{cases} x = -y \\ x(x^2 - 2) = 0 \end{cases}\end{aligned}\tag{7}$$

Donc les points critiques sont  $\{(-\sqrt{2}, \sqrt{2}); (\sqrt{2}, -\sqrt{2}); (0,0)\}$

- Pour le point  $(-\sqrt{2}, \sqrt{2})$   
On a  $H(-\sqrt{2}, \sqrt{2}) = \begin{bmatrix} 20 & 4 \\ 4 & 20 \end{bmatrix}$ , de trace 40 et de déterminant 384. Donc la Hessienne a deux valeurs propres strictement positives et  $(-\sqrt{2}, \sqrt{2})$  est un minimum local avec  $f(-\sqrt{2}, \sqrt{2}) = -8$
- Pour le point  $(\sqrt{2}, -\sqrt{2})$   
On a  $H(\sqrt{2}, -\sqrt{2}) = \begin{bmatrix} 20 & 4 \\ 4 & 20 \end{bmatrix}$ , de trace 40 et de déterminant 384. Donc la Hessienne a deux valeurs propres strictement positives et  $(\sqrt{2}, -\sqrt{2})$  est un minimum local avec  $f(\sqrt{2}, -\sqrt{2}) = -8$
- Pour le point  $(0,0)$   
On a  $H(0,0) = \begin{bmatrix} -4 & 4 \\ 4 & -4 \end{bmatrix}$  donc on ne peut pas conclure avec la matrice Hessienne

Prenons  $\epsilon \in R$  avec  $|\epsilon| < 2$

$$\begin{aligned}f(\epsilon, -\epsilon) &= 2\epsilon^4 - 8\epsilon^2 \\ &= -2\epsilon^2(4 - \epsilon^2)\end{aligned}\tag{8}$$

avec  $4 - \epsilon^2 > 0$  donc  $f(\epsilon, -\epsilon) < 0$  et  $f(\epsilon, \epsilon) = 2\epsilon^4 \geq 0$

Ainsi, on peut conclure que  $(0,0)$  est un point selle de  $f$ . On sait que  $\mathcal{P}$  possède au moins une solution, et donc

$$\inf_{(x,y) \in R^2} f(x,y) = f(\sqrt{2}, -\sqrt{2}) = f(-\sqrt{2}, \sqrt{2}) = -8\tag{9}$$

## 2 Partie Théorique

Nous allons maintenant nous intéresser au problème de minimisation de la fonction de Rosenbrock définie sur  $R^2$  par

$$f(x,y) = (x-1)^2 + 100(x^2 - y)^2\tag{10}$$

### 2.1 Etude théorique

(a) Trouver les points critiques de  $f$ .

Comme fonction polynomiale,  $f \in \mathcal{C}^\infty$  sur  $R^2$ . On détermine les dérivées partielles :

$$\frac{\partial f(x,y)}{\partial x} = 2x - 2 + 400(x^2 - y)x \text{ and } \frac{\partial f(x,y)}{\partial y} = 200(-x^2 + y)$$

On a donc

$$\begin{aligned}\nabla f(x,y) = \begin{bmatrix} 0 \\ 0 \end{bmatrix} &\Leftrightarrow \begin{cases} 2x - 2 + 400(x^2 - y)x = 0 \\ 200(-x^2 + y) = 0 \end{cases} \\ &\Leftrightarrow \begin{cases} x = 1 \\ y = 1 \end{cases}\end{aligned}\tag{11}$$

Donc  $f$  admet un unique point critique au point  $(1, 1)$

(b) Démontrer que  $f$  admet un minimum global qu'elle atteint en  $(1, 1)$ .

On détermine les dérivées secondes de  $f$  sur  $R^2$  avec  $f \in C^\infty$ .  $\forall (x, y) \in R^2$  :

$$\frac{\partial^2 f(x, y)}{\partial x^2} = 2 + 400(3x^2 - y) ; \frac{\partial^2 f(x, y)}{\partial y^2} = 200 \text{ and } \frac{\partial^2 f(x, y)}{\partial x \partial y} = -400x$$

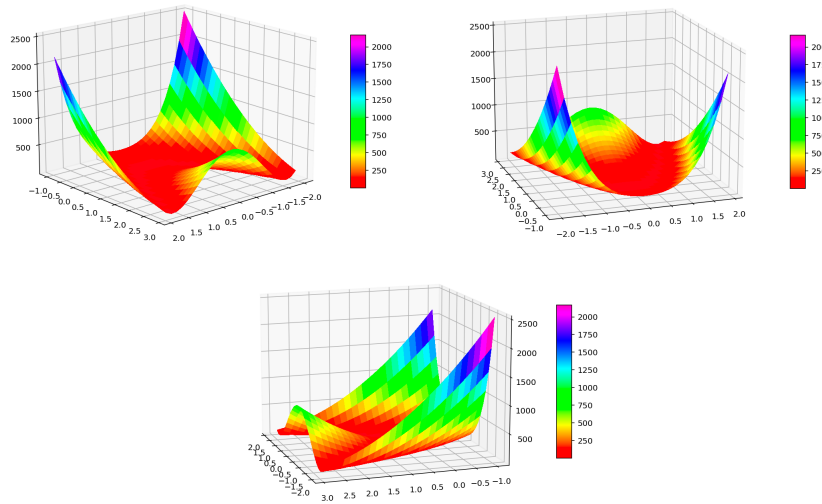
On a donc  $H(x, y) = \begin{bmatrix} 2 + 400(3x^2 - y) & -400x \\ -400x & 200 \end{bmatrix}$  et  $H(1, 1) = \begin{bmatrix} 802 & -400 \\ -400 & 200 \end{bmatrix}$  La trace vaut 1002 et le déterminant vaut 400. Ainsi, la Hessienne a deux valeurs propres strictement positives au point  $(1, 1)$  donc  $f$  admet un minimum local en  $(1, 1)$  avec  $f(1, 1) = 0$ . De plus,  $f$  est positive sur  $R^2$  donc  $f$  admet un minimum global en  $(1, 1)$ .

(c) Déterminer le Jacobien de la fonction  $f$  ainsi que sa Hessienne.

D'après les questions précédentes,

$$J_f = \begin{bmatrix} 2(x-1) + 400(x^2 - y)x \\ -200(x^2 - y) \end{bmatrix} \text{ et } H(x, y) = \begin{bmatrix} 2 + 400(3x^2 - y) & -400x \\ -400x & 200 \end{bmatrix}$$

## 2.2 Plot



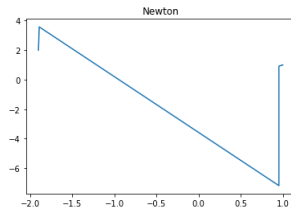
## 3 Partie Numérique : Methodes de descente

```

1 def NewtonR(x, y, epsilon):
2     x = [x]
3     y = [y]
4     start_time = time.time()
5     Matr = -np.linalg.inv(hessrosen(x[0], y[0])).dot(diffrosen(x[0], y[0])) + np.array([x[0],
6     y[0]])
7     x.append(Matr[0])
8     y.append(Matr[1])
9     n = 0
10    nmax = 10000
11    while (np.sqrt((x[n+1]-x[n])**2+(y[n+1]-y[n])**2) > epsilon and n < nmax):
12        Matr = -np.linalg.inv(hessrosen(x[n+1], y[n+1])).dot(diffrosen(x[n+1], y[n+1])) + np.
13        array([x[n+1], y[n+1]])
14        x.append(Matr[0])
15        y.append(Matr[1])
16        n += 1
17    rep = [x[-1], y[-1]]
18    X = np.transpose([x, y])
19    temps = time.time() - start_time
20    print(temps)
21    print(n)
22    return rep, X, temps

```

Listing 1: Méthode Newton



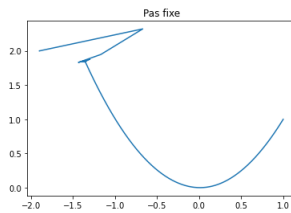
### 3.1 Méthode de descente à pas fixe

```

1 def gradfixe(x, y, eps, rho):
2     start_time = time.time()
3     k = 0
4     epsk = 2*eps
5     nmax = 100000
6     x_list = [x]
7     y_list = [y]
8     while (epsk > eps) and (k < nmax):
9         wk = - np.array(diffrosen(x_list[k], y_list[k]))
10        Xnew = np.array([x_list[k], y_list[k]]) + rho * wk
11        x_list.append(Xnew[0])
12        y_list.append(Xnew[1])
13        epsk = np.sqrt((x_list[k+1]-x_list[k])**2 + (y_list[k+1]-y_list[k])**2)
14        k += 1
15    rep = np.array([x_list[-1], y_list[-1]])
16    X = np.transpose(np.array([x_list, y_list]))
17    temps = time.time() - start_time
18    print(k)
19    print(epsk)
20    return rep, X, temps

```

Listing 2: Méthode de descente à pas fixe



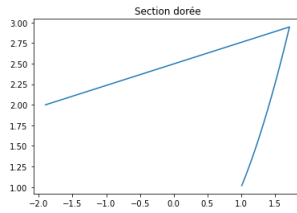
### 3.2 Section Dorée

```

1 def gradopt_section_doree(x, y, eps):
2     start_time = time.time()
3     k = 0
4     epsk = 2 * eps
5     nmax = 100000
6     x_list = [x]
7     y_list = [y]
8     while epsk > eps and k < nmax:
9         wk = - np.array(diffrosen(x_list[k], y_list[k]))
10        rhok = section_doree(x_list[k], y_list[k], wk, 0, 1, 50)
11        Xnew = np.array([x_list[k], y_list[k]]) + rhok * wk
12        x_list.append(Xnew[0])
13        y_list.append(Xnew[1])
14        epsk = np.sqrt((x_list[k + 1] - x_list[k]) ** 2 + (y_list[k + 1] - y_list[k]) ** 2)
15        k += 1
16    rep = np.array([x_list[-1], y_list[-1]])
17    X = np.transpose(np.array([x_list, y_list]))
18    temps = time.time() - start_time
19    print(k)
20    print(epsk)
21    return rep, X, temps

```

Listing 3: Grad Opt Section Dorée



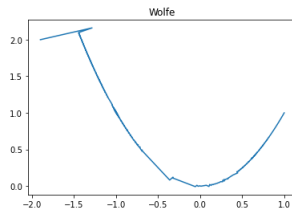
### 3.3 GradOpt Méthode Wolfe

```

1 def gradopt_m thode_wolfe(x, y, eps, rho):
2     start_time = time.time()
3     k = 0
4     epsk = 2 * eps
5     nmax = 100000
6     x_list = [x]
7     y_list = [y]
8     rhok = rho
9     while (epsk > eps) and (k < nmax):
10         wk = -np.array(diffrosen(x_list[k], y_list[k]))
11         rhok = m thode_wolfe(x_list[k], y_list[k], rhok, wk, 0.1, 0.9)
12         Xnew = np.array([x_list[k], y_list[k]]) + rhok * wk
13         x_list.append(Xnew[0])
14         y_list.append(Xnew[1])
15         epsk = np.sqrt((x_list[k+1] - x_list[k])**2 + (y_list[k+1] - y_list[k])**2)
16         k += 1
17     rep = np.array([x_list[-1], y_list[-1]])
18     X = np.transpose(np.array([x_list, y_list]))
19     temps = time.time() - start_time
20     print(k)
21     print(epsk)
22     return rep, X, temps

```

Listing 4: Grad Opt Méthode Wolfe



### 3.4 DFP

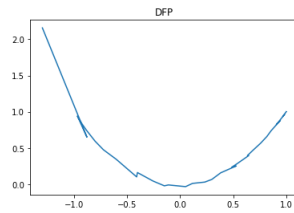
```

1 def dfp(x, y, eps, rho):
2     tic = time.time()
3     k = 0
4     rhok = rho
5     Bk = np.eye(2)
6     epsk = 2 * eps
7     nmax = 100000
8     x=[x]
9     y=[y]
10    while epsk > eps and k < nmax:
11        wk = -np.array(diffrosen(x[k], y[k]))
12        rhok = m thode_wolfe(x[k], y[k], rhok, wk, 0.01, 0.99)
13        Mat1 = rhok * Bk @ wk + np.array([x[k], y[k]])
14        x.append( Mat1[0])
15        y.append( Mat1[1])
16        sk = np.array([x[k+1], y[k+1]]) - np.array([x[k], y[k]])
17        yk = diffrosen(x[k+1], y[k+1]) - diffrosen(x[k], y[k])
18        Bk = Bk + np.outer(sk, sk.T) / (sk.T @ yk) - (Bk @ yk @ yk.T @ Bk.T) / (yk.T @ Bk @ yk)
19    )
20    epsk = np.sqrt((x[k+1] - x[k]) ** 2 + (y[k+1] - y[k]) ** 2)
21    k += 1
22    rep = np.array([x[-1], y[-1]])
23    X = np.vstack((x, y)).T
24    print(k)
25    temps = time.time() - tic

```

```
25 return rep, X, temps
```

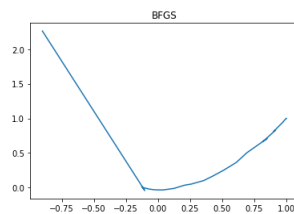
Listing 5: DFP



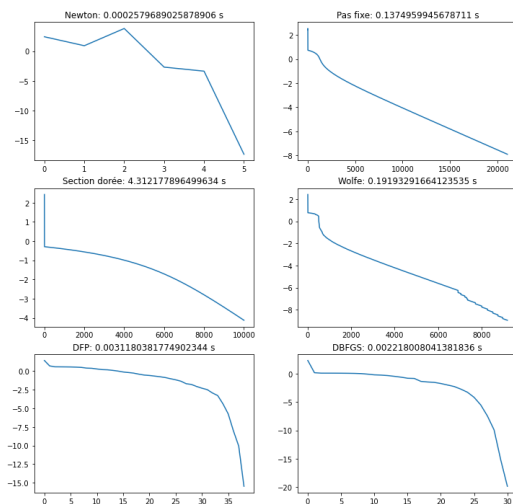
### 3.5 BFGS

```
1 def bfgs(x, y, eps, rho):
2     tic = time.time()
3     k = 0
4     rhok = rho
5     Bk = np.eye(2)
6     epsk = 2 * eps
7     nmax = 100000
8     x=[x]
9     y=[y]
10    while epsk > eps and k < nmax:
11        wk = -np.array(diffrosen(x[k], y[k]))
12        rhok = mthode_wolfe(x[k], y[k], rhok, wk, 0.01, 0.99)
13        Mat1 = rhok * Bk @ wk + np.array([x[k], y[k]])
14        x.append(Mat1[0])
15        y.append(Mat1[1])
16        sk = np.array([x[k+1], y[k+1]]) - np.array([x[k], y[k]])
17        yk = diffrosen(x[k+1], y[k+1]) - diffrosen(x[k], y[k])
18        Bk = Bk + np.outer(yk, yk.T) / (yk.T @ sk) - (Bk @ sk @ sk.T @ Bk.T) / (sk.T @ Bk @ sk)
19    epsk = np.sqrt((x[k+1] - x[k]) ** 2 + (y[k+1] - y[k]) ** 2)
20    k += 1
21    rep = np.array([x[-1], y[-1]])
22    X = np.vstack((x, y)).T
23    print(k)
24    temps = time.time() - tic
25    return rep, X, temps
```

Listing 6: BFGS



### 3.6 Recapitulatif des Méthodes



Les temps de calcul sont inclus à chaque début de graphique, et nous précisons les nombres d'itérations de chaque méthode :

$$\begin{pmatrix} \textit{Newton} & 6 \\ \textit{Pas fixe} & 21056 \\ \textit{section dorée} & 9996 \\ \textit{Wolfe} & 9196 \\ \textit{DFP} & 38 \\ \textit{BFGS} & 30 \end{pmatrix}$$

**Gradient à pas fixe:** Cette méthode utilise un pas fixe pour toutes les itérations, ce qui signifie que la mise à jour des variables est effectuée en utilisant un pas constant. Le temps d'itération et le temps de calcul sont généralement faibles, mais la convergence peut être lente, en particulier pour des fonctions complexes ou mal conditionnées.

**Pas optimal (section dorée):** La méthode de la section dorée est utilisée pour déterminer le pas optimal à chaque itération. Le temps d'itération et le temps de calcul sont généralement moyens, mais la convergence est généralement rapide. Cette méthode nécessite moins d'itérations pour converger vers la solution optimale que le gradient à pas fixe. Ici nous avons une petite erreur de temps de calcul pour la section dorée, donc c'est plutôt assez lent.

**Pas optimal (Wolfe):** La méthode de Wolfe est une autre approche pour déterminer le pas optimal. Elle est similaire à la section dorée en termes de temps d'itération et de temps de calcul, avec une convergence rapide.

**Quasi-Newton (DFP):** La méthode DFP (Davidon-Fletcher-Powell) est une méthode de Quasi-Newton qui utilise une approximation de la matrice hessienne pour améliorer la convergence. Le temps d'itération est généralement plus long, et le temps de calcul est plus élevé que les autres méthodes. Toutefois, elle présente une convergence très rapide, ce qui la rend souvent préférable pour des problèmes complexes.

**Quasi-Newton (BFGS):** La méthode BFGS (Broyden-Fletcher-Goldfarb-Shanno) est une autre méthode de Quasi-Newton. Elle est similaire à la méthode DFP en termes de temps d'itération et de temps de calcul, mais elle est généralement considérée comme plus efficace et plus robuste. La convergence est également très rapide.

En résumé, nous devrions avoir :

Méthode	Temps d'itération	Temps de calcul	Convergence	Complexité
Gradient à pas fixe	Rapide	Faible	Lente	Faible
Pas optimal (section dorée)	Moyen	Moyen	Rapide	Moyen
Pas optimal (Wolfe)	Moyen	Moyen	Rapide	Moyen
Quasi-Newton (DFP)	Lent	Élevé	Très rapide	Élevé
Quasi-Newton (BFGS)	Lent	Élevé	Très rapide	Élevé

Or nous obtenons des résultats légèrement différents, mais explicables, en effet : les méthodes de Quasi Newton sont sensées être plus lentes, il y a environ 30 itérations c'est très rapide. On peut dire que cela est expliqué par le fait qu'on inverse seulement une matrice  $2 \times 2$ , ainsi ce serait intéressant de tester avec une plus grande dimension.

### 3.7 Complexité

La complexité exacte de ces algorithmes dépend de plusieurs facteurs, tels que la taille du problème, la qualité de l'implémentation et les conditions initiales. Cependant, voici une estimation générale de la complexité en termes d'opérations pour chaque algorithme mentionné:

**Gradient à pas fixe:** La complexité de cet algorithme est généralement de l'ordre de  $O(n^2)$  pour  $n$  variables, car chaque itération nécessite le calcul du gradient et la mise à jour des variables. Le nombre total d'itérations dépend de la fonction et du pas fixe choisi.

**Pas optimal (section dorée):** La complexité de cet algorithme est également de l'ordre de  $O(n^2)$ , en

raison du calcul du gradient et de la recherche unidimensionnelle du pas optimal. Cependant, la convergence est généralement plus rapide, ce qui signifie que moins d'itérations sont nécessaires pour atteindre la solution optimale.

**Pas optimal (Wolfe):** La complexité de cet algorithme est similaire à celle de la section dorée, soit de l'ordre de  $O(n^2)$ . Là encore, la convergence est généralement rapide, nécessitant moins d'itérations que le gradient à pas fixe.

**Quasi-Newton (DFP):** La complexité de la méthode DFP est généralement de l'ordre de  $O(n^3)$  pour  $n$  variables, car elle nécessite le calcul du gradient et la mise à jour de la matrice hessienne inverse. Cependant, la convergence est généralement beaucoup plus rapide que celle des autres méthodes, ce qui peut compenser la complexité accrue.

**Quasi-Newton (BFGS):** La complexité de la méthode BFGS est également de l'ordre de  $O(n^3)$  pour  $n$  variables, en raison du calcul du gradient et de la mise à jour de la matrice hessienne inverse. Tout comme la méthode DFP, la convergence est généralement très rapide, ce qui peut compenser la complexité accrue.

Il est important de noter que ces estimations de complexité sont approximatives et peuvent varier en fonction des conditions spécifiques du problème et de l'implémentation utilisée.

Ici, le problème est à deux dimensions, et manque peut-être de représentativité quid de la dimension de nos méthodes, il serait donc intéressant de comparer ces résultats avec ceux obtenus en dimension 3 ainsi que dans des dimensions supérieures.

### 3.8 Gradient Conjugué

```

1 from scipy.sparse import diags
2 from numpy.linalg import inv
3
4 def GC(A, b, x0, eps):
5     start_time = time.time()
6     k = 0
7     epsk = np.sqrt(np.dot(A @ x0 - b, A @ x0 - b))
8     w0 = -(A @ x0 - b)
9     epsk = 2 * eps
10    nmax = 100
11    xk = x0
12    wk = w0
13    while (epsk > eps and k < nmax):
14        rhok = -np.dot(A @ xk - b, wk) / np.dot(A @ wk, wk)
15        xnew = xk + rhok * wk
16        wnew = -(A @ xnew - b) + ((np.dot(A @ xnew - b, A @ xnew - b)) / (np.dot(A @ xk - b, A
17        @ xk - b))) * wk
18        epsk = np.dot(A @ xnew - b, A @ xnew - b)
19        k += 1
20        wk = wnew
21        xk = xnew
22    temps = time.time() - start_time
23    return temps
24
25 def main_GC():
26     n_list=[10,20,30,50,100]
27     time_list=[]
28     for i in n_list:
29         n = i
30         vec4 = 4 * np.ones(n)
31         vec_m2 = -2 * np.ones(n-1)
32         A = diags([vec_m2, vec4, vec_m2], [-1, 0, 1], shape=(n, n)).toarray()
33         #print(A)
34         b= 1 * np.ones(n)
35         x= 0 * np.ones(n)
36         temps=GC(A, b, x, eps=10e-3)
37         time_list.append(temps)
38     print(time_list)

```

Listing 7: Gradient Conjugué



### 3.8.1 Résultats Gradient Conjugué

$$nlist = \begin{bmatrix} 10 \\ 20 \\ 30 \\ 50 \\ 100 \end{bmatrix} \quad \text{temps correspondants} = \begin{bmatrix} 0.00026702880859375 \\ 0.0003027915954589844 \\ 0.00041103363037109375 \\ 0.0007979869842529297 \\ 0.0018019676208496094 \end{bmatrix}$$

#### Analysons ces résultats :

Comme prévu, le temps d'exécution augmente avec la taille du problème (n). Cela est cohérent avec la complexité théorique de la méthode du gradient conjugué, qui dépend du nombre de variables (n).

Les temps d'exécution restent relativement faibles même pour  $n = 100$ , ce qui suggère que l'algorithme est efficace pour résoudre les problèmes de cette taille. Cela est probablement dû à la structure particulière de notre matrice A (presque diagonales et symétrique).

Les temps d'exécution n'augmentent pas de manière linéaire ou quadratique avec la taille du problème. Cela est la preuve que la convergence est atteinte en moins d'itérations que n. Cela est en adéquation avec le fait que la méthode de gradient conjugué est une méthode exacte qui converge donc en au plus n itérations.

En conclusion, ces résultats montrent que la méthode du gradient conjugué est efficace pour résoudre des problèmes de taille modérée (jusqu'à  $n = 100$ ). Il serait intéressant de modifier la matrice A tant par sa structure que par sa taille.