

# Calcul Numérique TP

YANG Yizhi

8 janvier 2024

## Table des matières

<b>1</b>	<b>Introduction : Problème Poisson 1D</b>	<b>2</b>
<b>2</b>	<b>Test d'environnement</b>	<b>2</b>
<b>3</b>	<b>Méthode directe</b>	<b>3</b>
3.1	Stockage Bande . . . . .	3
3.1.1	Déclaration et allocation en C . . . . .	3
3.1.2	la Constante <code>LAPACK_COL_MAJOR</code> . . . . .	3
3.1.3	la Dimension Principale . . . . .	3
3.2	Explication des Fonctions BLAS/LAPACK . . . . .	3
3.2.1	<code>dgbmv</code> . . . . .	3
3.2.2	<code>dgbtrf</code> . . . . .	3
3.2.3	<code>dgbtrs</code> . . . . .	3
3.2.4	<code>dgbsv</code> . . . . .	3
3.2.5	Calcul de la norme du résidu relatif . . . . .	4
3.2.6	Méthode de validation . . . . .	4
3.2.7	Performance et complexité . . . . .	4
<b>4</b>	<b>Méthodes itératives</b>	<b>5</b>
<b>5</b>	<b>Stockage CSR et CSC</b>	<b>6</b>
<b>6</b>	<b>Conclusion</b>	<b>6</b>

# 1 Introduction : Problème Poisson 1D

On s'intéresse au problème du Poisson à dimension 1. Le problème de Poisson 1D est un système des équations aux dérivées partielles couramment rencontrée en mathématiques appliquées, en physique et en ingénierie. Dans le monde de physique macroscopique, tous les problèmes sont continus dans l'espace réel, et pourtant, dans le monde de l'informatique, la discrétisation est toujours nécessaire.

On applique le maillage uniforme, c'est-à-dire sur le même axe, on prend des points équidistants, et à chaque point on aura une équation de dérivées partielles de second degré par rapport à une condition locale. Pendant les cours de TD précédents, on a écrit le système de ces équations en forme matricielle tridiagonale. Donc résoudre ce problème est équivalent à résoudre l'équation du produit matriciel  $Ax = b$ . Par une simplification de condition des seconds membres, on sait que la solution analytique est sous forme linéaire  $T(x) = T_0 + x(T_1 - T_0)$ .

Dans la suite de ce TP, on va d'abord introduire ce problème et puis proposer des méthodes afin d'approximer la solution, qui sera comparée avec la solution exacte.

## 2 Test d'environnement

En créant une image dans Docker Container, on configure l'environnement du travail pour que les bibliothèques CBLAS et LAPACK soient reconnues lors de la compilation. Un fichier de test a été créé et le résultat se donne :

```
The exponential value is e = 2.718282
The maximum single precision value from values.h is maxfloat = 3.402823e+38
The maximum single precision value from float.h is flt_max = 3.402823e+38
The maximum double precision value from float.h is dbl_max = 1.797693e+308
The epsilon in single precision value from float.h is flt_epsilon = 1.192093e-07
The epsilon in double precision value from float.h is dbl_epsilon = 2.220446e-16
```

```
Test of ATLAS (BLAS/LAPACK) environment
```

```
x[0] = 1.000000, y[0] = 6.000000
x[1] = 2.000000, y[1] = 7.000000
x[2] = 3.000000, y[2] = 8.000000
x[3] = 4.000000, y[3] = 9.000000
x[4] = 5.000000, y[4] = 10.000000
```

```
Test DCOPY y <- x
```

```
y[0] = 1.000000
y[1] = 2.000000
y[2] = 3.000000
y[3] = 4.000000
y[4] = 5.000000
```

## 3 Méthode directe

### 3.1 Stockage Bande

#### 3.1.1 Déclaration et allocation en C

En effet, on doit déclarer et allouer les matrices de manière à ce qu'elles soient compatibles avec les fonctions BLAS/LAPACK, donc de forme General Band ici. C'est une forme efficace pour les matrices multidiagonales, car la plupart de leur éléments sont nuls. Au lieu d'occuper une mémoire de taille  $n * n$ , il suffira de  $lab * n$ , avec  $lab$  le nombre des diagonales.

#### 3.1.2 la Constante LAPACK.COL\_MAJOR

Cette constante signifie que la matrice donnée est stockée colonne par colonne. Afin de convertir les indices usuels, il faut définir la fonction suivante :

```
1 // i, j from 0 to *la -1
2 int indexABCol(int i, int j, int *lab)
3 {
4     return (j + 1) * (*lab - 1) + i - 1;
5 }
```

Listing 1 – General Band Indices

#### 3.1.3 la Dimension Principale

Cela signifie la largeur de la bande stockée, i.e. le nombre total des diagonales (la diagonale principale, les subdiagonales et sousdiagonales).

## 3.2 Explication des Fonctions BLAS/LAPACK

#### 3.2.1 dgbmv

C'est une fonction BLAS qui effectue une multiplication d'une matrice (stockée en bande) par un vecteur. La lettre "d" signifie que les éléments contenus dans les matrices sont de type double.

#### 3.2.2 dgbtrf

C'est une routine LAPACK qui effectue la factorisation LU d'une matrice générale stockée en bande, donc les éléments sont de type double.

#### 3.2.3 dgbtrs

C'est une routine LAPACK qui résout un système d'équations linéaires pour une matrice stockée en bande après une factorisation LU. Ses éléments sont de type double.

#### 3.2.4 dgbsv

C'est similaire de dgbtrs mais avec pivotage partiel, qui est une combinaison de *dgbtrf* et *dgbtrs*.

### 3.2.5 Calcul de la norme du résidu relatif

La norme du résidu relatif s'écrit sous la forme :

$$\|r\|_2 = \frac{\|Ax - b\|_2}{\|b\|_2}$$

On applique d'abord la fonction *cblas\_dgbmv* pour calculer le produit matrice-vecteur  $Ax$ , puis on effectue l'addition de deux vecteurs  $Ax$  et  $b$  par la fonction *cblas\_daxpy*.

On pourra utiliser la fonction *cblas\_dnrm2* pour calculer la norme euclidienne.

### 3.2.6 Méthode de validation

Pour la solution du système  $Ax = b$ , on pourra comparer l'erreur relative de  $x$  par rapport à la solution exacte  $x_{exact}$  :

$$\|err\|_2 = \frac{\|x - e_{exact}\|_2}{\|x_{exact}\|_2}$$

On pourra aussi utiliser le résidu relatif, cela donne un ordre de grandeur  $10^{-16}$  :

The relative forward error is relres = 2.393518e-16

Pour la factorisation LU, on pourra calculer  $\|A - LU\|_2$ , si cette valeur est très petit, on peut dire que la méthode est validée.

### 3.2.7 Performance et complexité

On rajoute une mesure de temps pour les méthodes différentes en utilisant l'entête `time.h` et les fonctions suivantes :

```
1 clock_t start_timer, end_timer;
2 double elapsed_time;
3 start_timer = clock();
4 // code to calculate solution
5 end_timer = clock();
6 elapsed_time = (double)(end_timer - start_timer) / (double)CLOCKS_PER_SEC;
7 printf("Elapsed time for direct method: %lf seconds\n", elapsed_time);
```

Listing 2 – Time Measure

La méthode avec *LU* factorisation *dgbrf* puis *dgbrs* donne :

Elapsed time for direct method: 0.000823 seconds

La méthode direct *dgbsv* donne :

Elapsed time for direct method: 0.000909 seconds

Complexité : *dgbrf* admet une complexité  $O(n^3)$  pour la factorisation LU, *dgbrs* entre  $O(n^2)$  et  $O(n^3)$ ; pour *dgbsv*, il faut  $O(n)$  pour la factorisation LU et  $O(n)$  pour résoudre le système. (avec  $n$  la taille de matrice)

## 4 Méthodes itératives

Pour la méthode alpha-Richardson, on calcule le facteur  $\alpha$  en utilisant les valeurs propres de la matrice du Poisson 1D :

$$\alpha = \frac{2}{\max(\text{eigval}) + \min(\text{eigval})}$$

Pour la méthode Jacobi (respectivement Gauss-Seidel), on calcule d'abord la matrice  $M = D$  (respectivement  $M = D - E$ ), ensuite l'inverse de matrice  $M^{-1}$ . Par itération, on calcule  $x^{k+1} = x^k + M^{-1}r^k$  avec le résidu  $r^k = b - Ax$ .

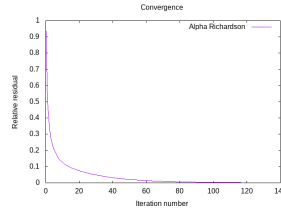


FIGURE 1 – Alpha-Richardson Iteration Convergence

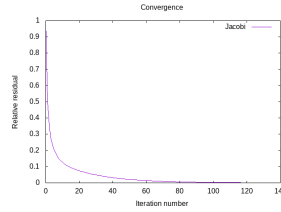


FIGURE 2 – Jacobi Iteration Convergence

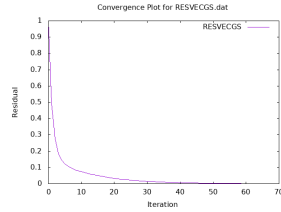


FIGURE 3 – Gauss-Seidel Iteration Convergence

Comme la valeur  $\alpha_{opt} = 0.5$  et  $D^{-1} = 0.5I_n$ , la méthode alpha-Richardson et la méthode Jacobi sont équivalentes :

Optimal alpha for simple Richardson iteration is : 0.500000  
The relative forward error is relres = 5.093853e-03  
Elapsed time for iterative methods: 0.000867 s

Elles partagent aussi la même vitesse de convergence, la méthode de Gauss-Seidel converge deux fois plus vite que celles-ci.

## 5 Stockage CSR et CSC

Pour la méthode CSR (Compressed Sparse Row), on stocke les valeurs significatives, les indices de colonnes et les pointeurs de lignes de la matrice ciblée dans trois tableaux unidimensionnels. Ici un algorithme possible pour la matrice tridiagonale du Poisson 1D :

---

### Algorithm 1 MatriceTridiagonaleVersCSR

---

**Data:** matrice

**Result:** values, columns, row\_ptr

```

// Initialisation des tableaux
1 values ← [] ; // Tableau pour les valeurs non nulles
2 columns ← [] ; // Tableau pour les indices de colonnes
3 row_ptr ← [] ; // Tableau pour les pointeurs de lignes
4 for chaque ligne i de la matrice do
    // Parcourir les éléments de la ligne
5    for chaque élément j de la ligne i do
6        if j est sur la diagonale principale then
7            Ajouter A[i, j] à values Ajouter j à columns
8        else if j est sur la diagonale supérieure ou inférieure then
9            Ajouter A[i, j] à values Ajouter j à columns
    // Ajouter la taille actuelle de values à row_ptr
10 Ajouter la taille actuelle de values à row_ptr
11 return values, columns, row_ptr

```

---

Similairement, pour un stockage CSC (Compressed Sparse Column), il faut just inverser les rôles de lignes et colonnes.

## 6 Conclusion

La convergence plus rapide de Gauss-Seidel par rapport à Jacobi dans le contexte spécifique du problème de Poisson 1D avec une matrice tridiagonale peut être attribuée à la nature séquentielle de la mise à jour, qui peut mieux exploiter la structure de la matrice et conduire à une convergence plus rapide.

Le choix entre le stockage en format général band (GB), Compressed Sparse Row (CSR), et Compressed Sparse Column (CSC) dépend des caractéristiques spécifiques de la matrice et des opérations que l'on a l'intention de réaliser. Le format GB est plus convenable pour les matrices tridiagonales et il donne un accès efficace aux éléments adjacents dans la bande ; les formats CSR

et CSC sont plus convenables pour les matrices creuses, car ils ne nécessitent pas un stockage pour les zéros, mais ils ont un surcoût pour accéder aux éléments dans la même colonne (pour CSR) ou la même ligne (pour CSC).