

TRAVAUX PRATIQUES - MÉTHODES DES DIFFÉRENCES FINIES 2D

TP 6 : MATRICES CREUSES, LAPLACIEN 2D ET PLUS SI AFFINITÉ

Table des matières

<b>1 Les matrices creuses</b>	<b>2</b>
1.1 Stockage . . . . .	2
1.1.1 Occupation mémoire . . . . .	3
1.1.2 Insertion d'éléments . . . . .	3
1.1.3 Construction efficace de matrices creuses . . . . .	4
1.2 Matrices diagonales et tridiagonales . . . . .	4
1.3 Produit de Kronecker de deux matrices (facultatif) . . . . .	5
<b>2 Notations et approximation des dérivées secondes</b>	<b>6</b>
<b>3 Equation de Poisson avec conditions de Dirichlet</b>	<b>6</b>
3.1 Evaluation de fonctions sur la grille discrète . . . . .	7
3.2 A blocs . . . . .	9
3.3 Assemblage de la matrice sans conditions aux limites . . . . .	10
3.3.1 Méthode 1 : utilisation de la bijection . . . . .	10
3.3.2 Méthode 2 : produits de Kronecker (facultatif) . . . . .	10
3.4 Quelques indices . . . . .	11
3.4.1 Indices des bords . . . . .	11
3.4.2 Evaluation de fonctions sur une partie de la grille . . . . .	12
3.5 Assemblage du second membre sans conditions aux limites . . . . .	12
3.6 Prise en compte des conditions aux limites . . . . .	13
3.7 Résolution du système . . . . .	13
3.7.1 Méthode 1 . . . . .	13
3.7.2 Méthode 2 . . . . .	14

En résolvant numériquement des EDPs par des méthodes de type différences finies, éléments finis et/ou volumes finis, on aboutit souvent à la résolution de grands systèmes linéaires creux (i.e. les matrices des systèmes linéaires contiennent un très grand nombre de zéros). Au niveau mémoire (de l'ordinateur), il est alors recommandé (voir nécessaire) d'utiliser un stockage particulier pour ces matrices : CSC, CSR, COO, LIL, DOK, ... Ces différents formats permettent de ne stocker en mémoire que les éléments non nuls des matrices.

Sous Matlab/Octave, le format utilisé est le format CSC (Compressed Sparse Column). Il est très efficace pour les opérations arithmétiques et les produits matrices vecteurs. Nous donnons deux utilisations (parmi d'autres) de la fonction `sparse` de Matlab/Octave :

- `M = sparse(m,n);`

Cette commande crée la matrice creuse nulle `M` de dimension  $m \times n$

- `M = sparse(I,J,K,m,n);`

Cette commande crée la matrice creuse `M` de dimension  $m \times n$  telle que

$$M(Ig(k),Jg(k)) \leftarrow M(Ig(k),Jg(k)) + Kg(k).$$

Les vecteurs `Ig`, `Jg` et `Kg` ont la même longueur. Les éléments nuls de `Kg` ne sont pas pris en compte et les éléments de `Kg` ayant les mêmes indices dans `Ig` et `Jg` sont sommés.

On donne maintenant dans plusieurs langages interprétés l'analogie de ces deux commandes :

- Python (`scipy.sparse` module) :
  - `M=sparse.<format>_matrix((m,n))`
  - `M=sparse.<format>_matrix((Kg, (Ig, Jg)), shape=(m,n))`
où `<format>` est le format de stockage de la matrice creuse (i.e. `csc`, `csr`, `lil`, ...),
- R : format CSC (seulement ?),

- `M=sparse(i=,j=,dims=c(m,n))`
- `M=sparse(i=Ig,j=Jg,x=Kg,dims=c(m,n))`
- Scilab : format row-by-row (?)
  - `M=sparse([],[],[m,n])`
  - `M=sparse([Ig,Jg],Kg,[m,n])`

Dans tous ces langages, de très nombreuses fonctions existent permettant d'utiliser les matrices creuses comme des matrices *classiques*. Il est donc possible d'effectuer des opérations aussi élémentaires que des sommes, produits de matrices et vecteurs mais aussi des opérations plus complexes comme des factorisations (LU, Cholesky, QR, ...), résolution de systèmes linéaires (méthodes directes ou itératives), calcul de valeurs propres et vecteurs propres, ...

Pour les langages de programmation usuels (Fortran, C, C++, ...) des libraires existent permettant de réaliser ces mêmes opérations sur les matrices creuses :

- C/C++ avec les librairies *SuiteSparse* de T. Davis [?] ainsi que *BLAS* (Basic Linear Algebra Subroutines), *LAPACK* (Linear Algebra PACKage) et leurs dérivées.
- CUDA avec la librairie *cuSparse* [?] de Nvidia.

## 1 Les matrices creuses

Nous avons commencé à utiliser les matrices creuses (*sparse matrix*) dans les TPs précédents. Par exemple, pour générer la matrice  $\mathbb{K} \in \mathcal{M}_d(\mathbb{R})$  définie par

$$\mathbb{K} = \begin{pmatrix} -2 & 1 & 0 & \dots & 0 \\ 1 & -2 & 1 & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & 1 \\ 0 & \dots & 0 & 1 & -2 \end{pmatrix} \quad (1.1)$$

nous avons écrit une fonction `Lap1D` pouvant ressembler à ceci :

Listing 1 – Fonction `Lap1D`

```
function K=Lap1D(d)
  K=sparse(d,d);
  K(1,1)=-2;K(1,2)=1;
  for i=2:d-1
    K(i,i)=-2;
    K(i,i-1)=1;
    K(i,i+1)=1;
  end
  K(d,d)=-2;K(d,d-1)=1;
end
```

Toutefois cette manière de construire une matrice creuse est très loin d'être la plus efficace.

Dans un premier temps, une rapide explication du pourquoi vous est proposée et ensuite nous verrons comment améliorer cela en utilisant judicieusement la fonction `sparse`. Il est aussi possible d'utiliser conjointement les fonctions `speye` et `spdiags` pour créer la matrice  $\mathbb{K}$  mais ces commandes ne seront plus utilisables pour assembler les matrices provenant des méthodes d'éléments finis et volumes finis, c'est pourquoi nous nous focaliserons sur la fonction `sparse`.

### 1.1 Stockage

Sous Matlab et Octave, une matrice creuse (ou *sparse matrix*),  $A \in \mathcal{M}_{M,N}(\mathbb{R})$ , est stockée sous le format CSC (Compressed Sparse Column), utilisant trois tableaux :

$$ia(1 : nnz), ja(1 : N + 1) \text{ et } aa(1 : nnz),$$

où  $nnz$  est le nombre d'éléments non nuls de la matrice  $A$ . Ces tableaux (cachés mais contenus dans l'object `sparse` Matlab/Octave) sont définis par

- $aa$  : contient l'ensemble des  $nnz$  éléments non nuls de  $A$  stockés par colonnes.
- $ia$  : contient les numéros de ligne des éléments stockés dans le tableau  $aa$ .
- $ja$  : permet de retrouver les éléments d'une colonne de  $A$ , sachant que le premier élément non nul de la colonne  $k$  de  $A$  est en position  $ja(k)$  dans le tableau  $aa$ . On a  $ja(1) = 1$  et  $ja(N + 1) = nnz + 1$ .

Pour illustrer, voici un exemple simple avec la matrice

$$A = \begin{pmatrix} 1. & 0. & 0. & 6. \\ 0. & 5. & 0. & 4. \\ 0. & 1. & 2. & 0. \end{pmatrix}$$

on a  $M = 3$ ,  $N = 4$ ,  $nnz = 6$  et

$aa$	1.	5.	1.	2.	6.	4.
------	----	----	----	----	----	----

$ia$	1	2	3	3	1	2
------	---	---	---	---	---	---

$ja$	1	2	4	5	7
------	---	---	---	---	---

Le premier élément non nul de la colonne  $k = 3$  de  $A$  est 2, ce nombre est en position 4 dans  $aa$ , donc  $ja(3) = 4$ .

### 1.1.1 Occupation mémoire

Pour illustrer le gain mémoire, nous allons calculer la taille mémoire occupée par la matrice tridiagonale  $\mathbb{K}$  définie en (1.1) dans le cas d'un stockage plein (*full matrix*) ou creux (*sparse matrix*). Cette matrice contient  $d + 2(d - 1) = 3d - 2$  éléments non nuls.

Sous Matlab/Octave le stockage creux de cette matrice va utiliser  $(3d - 2)$  `double` (tableau  $aa$ ),  $(3d - 2)$  `int` (tableau  $ia$ ) et  $(d + 1)$  `int` (tableau  $ja$ ) pour un total de

$$(3d - 2) \text{ double et } (4d - 1) \text{ int.}$$

Or un `double` occupe 8 `Octets` (en anglais *bytes*) et un `int` 4 `Octets`. La place mémoire occupée par la matrice stockée sous forme `sparse` est alors de

$$(3d - 2) \times 8 + (4d - 1) \times 4 = (40d - 20) \text{ Octets.}$$

Par contre, si la matrice était stockée classiquement (`full`), alors la place mémoire serait de

$$d^2 \text{ Octets.}$$

On rappelle que 1 Mo (Méga-octet) correspond à  $10^6$  `Octets`, 1 Go (Giga-octet) correspond à  $10^9$  `Octets` et 1 To (Tera-octet) correspond à  $10^{12}$  `Octets`. On donne maintenant la place mémoire occupée par la matrice suivant sa dimension et son stockage (`full` ou `sparse`) :

$d$	full matrix	sparse matrix
$10^3$	1 Mo	40 Ko
$10^4$	100 Mo	400 Ko
$10^5$	10 Go	4 Mo
$10^6$	1000 Go	40 Mo
$10^7$	100 To	400 Mo

### 1.1.2 Insertion d'éléments

Regardons les opérations à effectuer sur les tableaux  $aa$ ,  $ia$  et  $ja$  si l'on modifie la matrice  $A$  par la commande

$$A(1,2)=8;$$

La matrice devient alors

$$A = \begin{pmatrix} 1. & 8. & 0. & 6. \\ 0. & 5. & 0. & 4. \\ 0. & 1. & 2. & 0. \end{pmatrix}.$$

Dans cette opération un élément nul de  $A$  a été remplacé par une valeur non nulle 8 : il va donc falloir la stocker dans les tableaux sachant qu'aucune place n'est prévue. On suppose que les tableaux sont suffisamment grands (pas de souci mémoire), il faut alors décaler d'une case l'ensemble des valeurs des tableaux  $aa$  et  $ia$  à partir de la 3ème position puis copier la valeur 8 en  $aa(2)$  et le numéro de ligne 1 en  $ia(2)$  :

$aa$	1.	8.	5.	1.	2.	6.	4.
$ia$	1	1	2	3	3	1	2

Pour le tableau  $ja$ , il faut à partir du numéro de colonne 2 plus un, incrémenter de +1 :

$ja$	1	2	5	6	8
------	---	---	---	---	---

La répétition de ces opérations peut devenir très coûteuse en temps CPU lors de l'assemblage de matrices utilisant des méthodes similaires à celle utilisée en Listing 1 (et on ne parle pas ici des problèmes de réallocation dynamique qui peuvent arriver!).

### 1.1.3 Construction efficace de matrices creuses

Pour générer efficacement une matrice creuse, il faut donc éviter d'insérer de manière répétitive des éléments dans celle-ci. Pour cela, on va utiliser l'appel suivant de la fonction `sparse` :

$$M = \text{sparse}(I, J, K, m, n);$$

Cette commande génère une matrice creuse  $m$  par  $n$  telle que  $M(I(k), J(k)) = K(k)$ . Les vecteurs  $I$ ,  $J$  et  $K$  ont la même longueur. Il faut noter que tous les éléments nuls de  $K$  sont ignorés et que tous les éléments de  $K$  ayant les mêmes indices dans  $I$  et  $J$  sont sommés. Cette sommation est très utile lors de l'assemblage de matrices obtenues par une méthode de type éléments finis.

Par exemple, voici deux commandes permettant de générer la matrice  $A$  précédente :

$$A = \text{sparse}([1 \ 1 \ 2 \ 3 \ 3 \ 1 \ 2], [1 \ 2 \ 2 \ 2 \ 3 \ 4 \ 4], [1 \ 8 \ 5 \ 1 \ 2 \ 6 \ 4], 3, 4);$$

$$A = \text{sparse}([1 \ 2 \ 1 \ 1 \ 2 \ 3 \ 3], [4 \ 4 \ 1 \ 2 \ 2 \ 2 \ 3], [6 \ 4 \ 1 \ 8 \ 5 \ 1 \ 2], 3, 4);$$

Il existe aussi la fonction  $A = \text{spalloc}(m, n, nnz)$  qui préalloue une matrice creuse de dimension  $m$  par  $n$  avec au plus  $nnz$  éléments non nuls (*nonzero elements*).

Nous allons maintenant utiliser ceci dans un cadre un peu général : les matrices tridiagonales.

## 1.2 Matrices diagonales et tridiagonales

**Q. 1 (Matlab)** *Ecrire la fonction `spMatDiag` permettant à partir d'un vecteur  $v \in \mathbb{R}^d$  de retourner la matrice diagonale (creuse)  $D \in \mathcal{M}_d(\mathbb{R})$  de diagonale  $v$  (i.e.  $D(i, i) = v(i), \forall i \in \llbracket 1, d \rrbracket$ ) en créant (sans boucle) les tableaux  $I, J$  et  $K$  puis en générant la matrice  $D$  à l'aide de la commande  $D = \text{sparse}(I, J, K, d, d)$ ; ■*

Soit  $A \in \mathcal{M}_d(\mathbb{R})$  la matrice tridiagonale définie par

$$A = \begin{pmatrix} v_1 & w_1 & 0 & \cdots & \cdots & 0 \\ u_1 & v_2 & w_2 & \ddots & & \vdots \\ 0 & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 \\ \vdots & & \ddots & \ddots & v_{d-1} & w_{d-1} \\ 0 & \cdots & \cdots & 0 & u_{d-1} & v_d \end{pmatrix} \quad (1.2)$$

avec  $u \in \mathbb{R}^{d-1}$ ,  $v \in \mathbb{R}^d$  et  $w \in \mathbb{R}^{d-1}$ .

**Q. 2 (Matlab)** *1. Ecrire la fonction `spMatTriDiagSca` permettant à partir des vecteurs  $u, v$  et  $w$  de retourner la matrice creuse  $A$  définie en (1.2). Pour cela on va tout d'abord créer une matrice creuse de  $\mathcal{M}_d(\mathbb{R})$  puis on va la «remplir» à l'aide d'une boucle `for` sur le principe du Listing 1.*

2. Ecrire la fonction `spMatTriDiagVec` permettant à partir des vecteurs  $\mathbf{u}$ ,  $\mathbf{v}$  et  $\mathbf{w}$  de retourner la matrice creuse  $\mathbb{A}$  en créant (sans boucle) les tableaux  $\mathbf{I}$ ,  $\mathbf{J}$  et  $\mathbf{K}$  puis en générant la matrice  $\mathbb{A}$  à l'aide de la commande  $\mathbf{A} = \text{sparse}(\mathbf{I}, \mathbf{J}, \mathbf{K}, \mathbf{d}, \mathbf{d})$ ;
3. Ecrire un programme permettant de mesurer (tic-toc) puis de représenter graphiquement les performances de `spMatTriDiagSca` et `spMatTriDiagVec` en fonction de  $d$ .  
(Attention à la capacité mémoire de la machine !)

### 1.3 Produit de Kronecker de deux matrices (facultatif)

Pour générer efficacement et simplement les matrices creuses obtenus par des méthodes de différences finies en dimension  $\geq 2$  nous utiliserons le produit de Kronecker de deux matrices.

Soient  $\mathbb{A} \in \mathcal{M}_{m,n}(\mathbb{R})$  et  $\mathbb{B} \in \mathcal{M}_{p,q}(\mathbb{R})$ . Le produit tensoriel de Kronecker de  $\mathbb{A}$  par  $\mathbb{B}$ , noté  $\mathbb{A} \otimes \mathbb{B}$ , est la matrice de  $\mathcal{M}_{mp,nq}(\mathbb{R})$  définie avec des blocs de dimension  $p \times q$  par

$$\mathbb{A} \otimes \mathbb{B} = \begin{pmatrix} A_{1,1}\mathbb{B} & \cdots & A_{1,n}\mathbb{B} \\ \vdots & \ddots & \vdots \\ A_{m,1}\mathbb{B} & \cdots & A_{m,n}\mathbb{B} \end{pmatrix} \quad (1.3)$$

Le produit de Kronecker est bilinéaire et associatif : Si les dimensions des matrices  $\mathbb{A}$ ,  $\mathbb{B}$  et  $\mathbb{C}$  sont compatibles on a  $\forall \lambda \in \mathbb{K}$

$$\begin{aligned} \mathbb{A} \otimes (\mathbb{B} + \lambda \cdot \mathbb{C}) &= (\mathbb{A} \otimes \mathbb{B}) + \lambda(\mathbb{A} \otimes \mathbb{C}) \\ (\mathbb{A} + \lambda \cdot \mathbb{B}) \otimes \mathbb{C} &= (\mathbb{A} \otimes \mathbb{C}) + \lambda(\mathbb{B} \otimes \mathbb{C}) \\ \mathbb{A} \otimes (\mathbb{B} \otimes \mathbb{C}) &= (\mathbb{A} \otimes \mathbb{B}) \otimes \mathbb{C} \end{aligned}$$

Par contre, il n'est pas commutatif.

Nous allons écrire une fonction permettant d'effectuer ce produit. Mais auparavant, quelques petits rappels (ou non) des possibilités offertes par le langage Matlab.

Il est facile sous Malab/Octave de modifier une matrice (creuse ou non). Par exemple,

- `A(I,J)=M`

Les tableaux d'indices  $\mathbf{I}$  et  $\mathbf{J}$  sont de dimension respective  $n$  et  $m$ . La matrice  $\mathbf{M}$  est de dimension  $n$ -par- $m$ . La matrice  $\mathbf{A}$  est modifiée de telle sorte que

$$\mathbf{A}(\mathbf{I}(\mathbf{i}), \mathbf{J}(\mathbf{j})) = \mathbf{M}(\mathbf{i}, \mathbf{j}), \quad \forall \mathbf{i} \in \llbracket 1, n \rrbracket, \quad \forall \mathbf{j} \in \llbracket 1, m \rrbracket$$

- `A(I,J)=A(I,J)+M`

Même chose en sommant.

- `A(i,:)=a`

Ici  $\mathbf{a}$  est un scalaire. Tous les éléments de la ligne  $\mathbf{i}$  de la matrice  $\mathbf{A}$  valent  $\mathbf{a}$ .

- `A(:,j)=a`

Ici  $\mathbf{a}$  est un scalaire. Tous les éléments de la colonne  $\mathbf{j}$  de la matrice  $\mathbf{A}$  valent  $\mathbf{a}$ .

- ...

Sous Matlab/Octave, il est possible de récupérer tous les éléments non nuls d'une matrice creuse  $\mathbf{A}$  ainsi que leurs indices de ligne et de colonne à l'aide de la commande

$$\llbracket \mathbf{I}, \mathbf{J}, \mathbf{K} \rrbracket = \text{find}(\mathbf{A});$$

Ici les trois tableaux  $\mathbf{I}$ ,  $\mathbf{J}$  et  $\mathbf{K}$  ont même longueur et on a  $\mathbf{A}(\mathbf{I}(\mathbf{k}), \mathbf{J}(\mathbf{k})) == \mathbf{K}(\mathbf{k})$  pour tout  $\mathbf{k}$  inférieur ou égal à `length(K)`.

**Q. 3 (Matlab)** Ecrire la fonction `spMatKron` permettant à partir deux matrices creuses carrées  $\mathbb{A} \in \mathcal{M}_n(\mathbb{R})$  et  $\mathbb{B} \in \mathcal{M}_m(\mathbb{R})$  de retourner la matrice creuse  $\mathbb{C} = \mathbb{A} \otimes \mathbb{B}$ . Pour celà, on initialisera la matrice retournée  $\mathbf{C}$  à l'aide de la commande `C=sparse(n*m,n*m)` puis on utilisera uniquement une boucle sur les éléments non-nuls de la matrice  $\mathbb{A}$  pour «remplir» la matrice  $\mathbf{C}$ .

On pourra comparer cette fonction à la fonction `kron` de Matlab/Octave.

## 2 Notations et approximation des dérivées secondes

Soient  $\Omega = ]a, b[ \times ]c, d[ \subset \mathbb{R}^2$  et  $\Gamma = \partial\Omega$  la frontière du domaine  $\Omega$ . On note  $\Gamma_N, \Gamma_S, \Gamma_O$  et  $\Gamma_E$  respectivement les frontières nord, sud, ouest et est. on a

$$\Gamma = \Gamma_N \cup \Gamma_S \cup \Gamma_O \cup \Gamma_E.$$

Soit  $v : \Omega \rightarrow \mathbb{R}$  suffisamment régulière.

**Q. 4 (sur feuille)** Montrer que l'on a

$$\frac{\partial^2 v}{\partial x^2}(x, y) = \frac{v(x+h, y) - 2v(x, y) + v(x-h, y)}{h^2} + \mathcal{O}(h^2) \quad (2.1)$$

$$\frac{\partial^2 v}{\partial y^2}(x, y) = \frac{v(x, y+h) - 2v(x, y) + v(x, y-h)}{h^2} + \mathcal{O}(h^2). \quad (2.2)$$

On note  $(x_i)_{i=0}^{N_x}$  et  $(y_j)_{j=0}^{N_y}$  les discrétisations régulières, respectivement, des intervalles  $[a, b]$  et  $[c, d]$  définies par

$$x_i = a + ih_x, \quad \forall i \in \llbracket 0, N_x \rrbracket \quad \text{et} \quad y_j = c + jh_y, \quad \forall j \in \llbracket 0, N_y \rrbracket \quad (2.3)$$

avec  $h_x = (b-a)/N_x$  et  $h_y = (d-c)/N_y$ . On note aussi

$$n_x = N_x + 1, \quad n_y = N_y + 1 \quad \text{et} \quad N = n_x \times n_y \quad (2.4)$$

**Q. 5 (sur feuille)** Montrer que l'on a  $\forall i \in \llbracket 0, N_x \rrbracket, \forall j \in \llbracket 0, N_y \rrbracket$

$$\begin{aligned} \Delta v(x_i, y_j) &\stackrel{\text{def}}{=} \left( \frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2} \right)(x_i, y_j) \\ &= \frac{v(x_{i+1}, y_j) - 2v(x_i, y_j) + v(x_{i-1}, y_j)}{h_x^2} \\ &\quad + \frac{v(x_i, y_{j+1}) - 2v(x_i, y_j) + v(x_i, y_{j-1})}{h_y^2} + \mathcal{O}(h_x^2) + \mathcal{O}(h_y^2) \end{aligned} \quad (2.5)$$

## 3 Equation de Poisson avec conditions de Dirichlet

Soient  $f : \Omega \rightarrow \mathbb{R}$  et  $g : \Gamma \rightarrow \mathbb{R}$  deux fonctions données. On veut résoudre le problème suivant

$$-\Delta u = f, \quad \text{dans } \Omega \quad (3.1)$$

$$u = g, \quad \text{sur } \Gamma \quad (3.2)$$

On utilise par la suite les discrétisations  $(x_i)_{i=0}^{N_x}$  et  $(y_j)_{j=0}^{N_y}$  définies en section 2.

**Q. 6 (sur feuille)** 1. Montrer que ce problème peut s'écrire après discrétisation sous la forme

$$-\frac{U_{i+1,j} - 2U_{i,j} + U_{i-1,j}}{h_x^2} - \frac{U_{i,j+1} - 2U_{i,j} + U_{i,j-1}}{h_y^2} = f(x_i, y_j), \quad \forall (i, j) \in \llbracket 0, N_x \rrbracket \times \llbracket 0, N_y \rrbracket, \quad (3.3)$$

$$U_{0,j} = g(a, y_j), \quad \forall j \in \llbracket 0, N_y \rrbracket, \quad (3.4)$$

$$U_{N_x,j} = g(b, y_j), \quad \forall j \in \llbracket 0, N_y \rrbracket, \quad (3.5)$$

$$U_{i,0} = g(x_i, c), \quad \forall i \in \llbracket 0, N_x \rrbracket, \quad (3.6)$$

$$U_{i,N_y} = g(x_i, d), \quad \forall i \in \llbracket 0, N_x \rrbracket. \quad (3.7)$$

avec  $U_{i,j} \approx u(x_i, y_j)$ .

2. Les inconnues du problème discrétisé sont les  $U_{i,j}$ , pour  $i \in \llbracket 0, N_x \rrbracket$  et  $j \in \llbracket 0, N_y \rrbracket$ . Compter le nombre d'équations distinctes du problème discrétisé (bord compris). Que peut-on en conclure ? ■

En notant  $\beta_x = -\frac{1}{h_x^2}$ ,  $\beta_y = -\frac{1}{h_y^2}$  et  $\mu = 2(\frac{1}{h_x^2} + \frac{1}{h_y^2})$ , les équations (3.3) peuvent aussi s'écrire sous la forme

$$\beta_x(U_{i+1,j} + U_{i-1,j}) + \mu U_{i,j} + \beta_y(U_{i,j+1} + U_{i,j-1}) = f(x_i, y_j), \quad \forall (i, j) \in \llbracket 0, N_x \rrbracket \times \llbracket 0, N_y \rrbracket \quad (3.8)$$

Pour tout  $j \in \llbracket 0, N_y \rrbracket$ , on note  $U_{:,j}$  le vecteur de  $\mathbb{R}^{n_x}$  défini par

$$U_{:,j} = \begin{pmatrix} U_{0,j} \\ \vdots \\ U_{N_x,j} \end{pmatrix}.$$

On note  $\mathbf{V} \in \mathbb{R}^N$  le vecteur bloc

$$\mathbf{V} = \begin{pmatrix} U_{:,0} \\ \hline U_{:,1} \\ \hline \vdots \\ \hline U_{:,N_y} \end{pmatrix}$$

**Q. 7 (sur feuille)** *Explicitiez la bijection  $\mathcal{F} : \llbracket 0, N_x \rrbracket \times \llbracket 0, N_y \rrbracket \longrightarrow \llbracket 1, N \rrbracket$  telle que*

$$\forall (i, j) \in \llbracket 0, N_x \rrbracket \times \llbracket 0, N_y \rrbracket, \quad V_k = U_{i,j}, \quad \text{avec } k = \mathcal{F}(i, j).$$

Dans le cas de la numérotation en  $(i, j) \in \llbracket 0, N_x \rrbracket \times \llbracket 0, N_y \rrbracket$  on parlera de **numérotation 2D** et pour la numérotation en  $k \in \llbracket 1, N \rrbracket$  on parlera de **numérotation globale**.

**Q. 8 (Matlab)** 1. *Ecrire la fonction `k=bijF(i,j,nx)` correspondant à la bijection  $\mathcal{F}$  (**numérotation 2D** vers **numérotation globale**).*

2. *Ecrire la fonction réciproque `[i,j]=bijRecF(k,nx)` correspondant à  $\mathcal{F}^{-1}$  (**numérotation globale** vers **numérotation 2D**). On pourra utiliser la fonction `rem(x,y)` (reste de la division de  $x$  par  $y$ ) et on vérifiera que cette fonction est naturellement vectorisée, c'est à dire que le paramètre  $k$  peut aussi être un tableau. ■*

### 3.1 Evaluation de fonctions sur la grille discrète

Soit  $\mathbf{F} \in \mathbb{R}^N$  le vecteur (bloc) défini par

$$\mathbf{F} = \begin{pmatrix} F_{:,0} \\ \hline F_{:,1} \\ \hline \vdots \\ \hline F_{:,N_y} \end{pmatrix} \quad \text{avec } F_{:,j} = \begin{pmatrix} f(x_0, y_j) \\ \vdots \\ f(x_{N_x}, y_j) \end{pmatrix} \in \mathbb{R}^{N_x+1}, \quad \forall j \in \llbracket 0, N_y \rrbracket.$$

En Algorithme 1, est présenté une version non vectorisée (deux boucles imbriquées) de la construction du vecteur  $\mathbf{F}$ . Dans cet algorithme l'ordre des boucles est primordial pour respecter le choix de la numérotation globale.

---

**Algorithme 1** Algorithme non vectorisé de calcul du vecteur  $\mathbf{F}$

---

```

k ← 1
for j ← 1 to ny do
  for i ← 1 to nx do
    F(k) ← f(x(i), y(j))
    k ← k + 1
  end for
end for

```

---

**Q. 9 (Matlab)** *Ecrire la fonction `EvalFun2DSCa` permettant à partir d'une fonction  $f$  donnée et des discrétisations  $\mathbf{x}$  et  $\mathbf{y}$  de retourner le vecteur  $\mathbf{F}$  associé. On pourra utiliser des boucles `for`. ■*

Pour le cas où la fonction  $f$  est vectorisée sous Matlab/Octave, il est alors possible de calculer le vecteur  $\mathbf{F}$  sans boucle. Pour cela, nous allons construire les deux vecteurs blocs  $\mathbf{X}$  et  $\mathbf{Y}$  de  $\mathbb{R}^n$  définis par

$$\mathbf{X} = \begin{pmatrix} X_{:,0} \\ \vdots \\ X_{:,N_y} \end{pmatrix} \text{ avec } X_{:,j} = \begin{pmatrix} x_0 \\ \vdots \\ x_{N_x} \end{pmatrix} \in \mathbb{R}^{N_x+1}, \quad \forall j \in \llbracket 0, N_y \rrbracket$$

et

$$\mathbf{Y} = \begin{pmatrix} Y_{:,0} \\ \vdots \\ Y_{:,N_y} \end{pmatrix} \text{ avec } Y_{:,j} = \begin{pmatrix} y_j \end{pmatrix} \in \mathbb{R}^{N_x+1}, \quad \forall j \in \llbracket 0, N_y \rrbracket.$$

Avec ces deux vecteurs, le calcul du vecteur  $\mathbf{F}$  pourra être vectorisé (sous la condition que  $f$  le soit) :

$$\mathbf{F} = \mathbf{f}(\mathbf{X}, \mathbf{Y});$$

Il nous faut donc écrire une version vectorisée du calcul des vecteurs  $\mathbf{X}$  et  $\mathbf{Y}$ , mais auparavant on donne plusieurs manières d'écrire le calcul du vecteur  $\mathbf{F}$

---

**Algorithme 2** Calcul de  $\mathbf{F}$  avec construction de  $\mathbf{X}$  et  $\mathbf{Y}$  (version 1)

---

```

k ← 1
for j ← 1 to ny do
  for i ← 1 to nx do
     $\mathbf{X}(k) \leftarrow \mathbf{x}(i)$ 
     $\mathbf{Y}(k) \leftarrow \mathbf{y}(j)$ 
     $\mathbf{F}(k) \leftarrow f(\mathbf{X}(k), \mathbf{Y}(k))$ 
    k ← k + 1
  end for
end for

```

---



---

**Algorithme 3** Calcul de  $\mathbf{F}$  avec construction de  $\mathbf{X}$  et  $\mathbf{Y}$  (version 2)

---

```

1: k ← 1
2: for j ← 1 to ny do
3:   for i ← 1 to nx do
4:      $\mathbf{X}(k) \leftarrow \mathbf{x}(i)$ 
5:      $\mathbf{Y}(k) \leftarrow \mathbf{y}(j)$ 
6:     k ← k + 1
7:   end for
8: end for
9:  $\mathbf{F} \leftarrow f(\mathbf{X}, \mathbf{Y})$ 

```

---

On peut noter que l'ordre des boucles dans l'Algorithme 3 permet d'affirmer que la valeur de  $k$  en ligne 4 est donnée par  $\mathcal{F}(i-1, j-1)$ .

---

**Algorithme 4** Calcul de  $\mathbf{F}$  avec construction de  $\mathbf{X}$  et  $\mathbf{Y}$  (version 3)

---

```

for j ← 1 to ny do
  for i ← 1 to nx do
    k ← BIJF(i-1, j-1, nx)
     $\mathbf{X}(k) \leftarrow \mathbf{x}(i)$ 
     $\mathbf{Y}(k) \leftarrow \mathbf{y}(j)$ 
  end for
end for
 $\mathbf{F} \leftarrow f(\mathbf{X}, \mathbf{Y})$ 

```

---

Grâce à l'application réciproque  $\mathcal{F}^{-1}$  il est alors possible de transformer la double boucle en  $j$  et  $i$  en une seule boucle sur  $k$  : c'est l'objet de l'Algorithme 5. On en déduit alors l'Algorithme 6 totalement vectorisé.



**Algorithme 5** Calcul de  $\mathbf{F}$  avec construction de  $\mathbf{X}$  et  $\mathbf{Y}$  (version 4)

```

for k ← 1 to N do
  [i, j] ← BIJRECFF(k, nx)
  X(k) ← x(i + 1)
  Y(k) ← y(j + 1)
end for
F ← f(X, Y)

```

**Algorithme 6** Calcul de  $\mathbf{F}$  avec construction de  $\mathbf{X}$  et  $\mathbf{Y}$  (version 5 : vectorisée)

```

[I, J] ← BIJRECFF(1 : N, nx) ▷ Deja vectorisee
X ← x(I + 1)
Y ← y(J + 1)
F ← f(X, Y)

```

Un programme Matlab/Octave complet basé sur ce dernier algorithme est proposé en Listing 3.

Listing 2 – Calcul et représentation d’une fonction sur la grille 2D

```

clear all
close all
a=-pi; b=pi; c=0; d=pi;
f=@(x, y) 4*cos(x + y).*sin(x - y);
Nx=133; Ny=222;
x=linspace(a, b, Nx+1);
y=linspace(c, d, Ny+1);
nx=Nx+1; ny=Ny+1; N=nx*ny;

[I, J]=bijRecF(1:N, nx); % Naturellement vectorisee
X=x(I+1); Y=y(J+1); % Attention X et Y vecteurs ←
    ligne
F=f(X, Y); % et donc F aussi!

surf(x, y, reshape(F, nx, ny) ')
shading interp

```

Il faut noter qu’il existe sous Matlab/Octave deux fonctions `meshgrid` et `ndgrid` qui peuvent être utilisées pour le calcul des vecteurs  $\mathbf{X}$  et  $\mathbf{Y}$ . Ces deux fonctions n’utilisent pas la fonction `bijRecF`.

D’autres techniques peuvent aussi être utilisées, par exemple avec la fonction `repmat` et l’opérateur `(:)` (ou la fonction `reshape` ou la commande `substruct('()', {' ':''})`).

**Q. 10 (Matlab)** Ecrire la fonction `EvalFun2DVec` permettant à partir d’une fonction  $f$  et des discrétisations  $\mathbf{x}$  et  $\mathbf{y}$  de retourner le vecteur  $\mathbf{F}$  associé sans utiliser de boucle `for`. ■

### 3.2 A blocs

Chacune des équations du problème discret (3.3)-(3.7) correspond à une discrétisation en un point  $(x_i, y_j)$ . Nous choisissons d’écrire ces équations en utilisant la même numérotation que lors de la construction du vecteur  $\mathbf{V}$  : l’équation écrite au point  $(x_i, y_j)$  sera écrite en ligne  $k = \mathcal{F}(i, j)$  du système.

**Q. 11 (sur feuille)** Etablir que le problème discret (3.3)-(3.7) peut s’écrire sous la forme du système linéaire bloc

$$\left( \begin{array}{c|cccc|c|c}
 \mathbf{E} & \mathbf{0} & \cdots & \cdots & \cdots & \mathbf{0} & \mathbf{0} \\
 \mathbf{M} & \mathbf{D} & \mathbf{M} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\
 \hline
 \mathbf{0} & \mathbf{M} & \ddots & \ddots & \ddots & \vdots & \vdots \\
 \vdots & \mathbf{0} & \ddots & \ddots & \ddots & \mathbf{0} & \vdots \\
 \vdots & \vdots & \ddots & \ddots & \ddots & \mathbf{M} & \mathbf{0} \\
 \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{M} & \mathbf{D} & \mathbf{M} \\
 \hline
 \mathbf{0} & \mathbf{0} & \cdots & \cdots & \cdots & \mathbf{0} & \mathbf{E}
 \end{array} \right) \mathbf{V} = \begin{pmatrix} B_{:,0} \\ B_{:,1} \\ \vdots \\ B_{:,N_y} \end{pmatrix} \quad (3.9)$$

où chaque bloc de la matrice est une matrice  $(N_x + 1)$  par  $(N_x + 1)$ . La matrice  $\mathbb{0} \in \mathcal{M}_{N_x+1}(\mathbb{R})$  est la matrice nulle. Les matrices creuses  $\mathbb{D}$ ,  $\mathbb{M}$  et  $\mathbb{E}$  ainsi que les vecteurs  $B_{:,j} \in \mathbb{R}^{N_x+1}$ , pour tout  $j \in \llbracket 0, N_y \rrbracket$ , devront être donnés explicitement. On notera  $\mathbf{B}$  le vecteur second membre de  $\mathbb{R}^N$  et  $\mathbb{A} \in \mathcal{M}_N(\mathbb{R})$  la matrice du système. ■

### 3.3 Assemblage de la matrice sans conditions aux limites

#### 3.3.1 Méthode 1 : utilisation de la bijection

Cette méthode est celle qui va, au final, se rapprocher le plus des techniques d'assemblages utilisées par les méthodes d'éléments finis et volumes finis sur des domaines  $\Omega \subset \mathbb{R}^d$  quelconques.

Nous avons vu que notre problème discret revient à résoudre le système linéaire

$$\mathbb{A}\mathbf{V} = \mathbf{B}. \quad (3.10)$$

Il nous faut donc construire/assembler la matrice  $\mathbb{A} \in \mathcal{M}_N(\mathbb{R})$  et le second membre  $\mathbf{B} \in \mathbb{R}^N$ .

Nous allons commencer par générer/assembler la matrice correspondant à la prise en compte de toutes les équations de (3.3) sans tenir compte, pour le moment, des conditions aux limites (i.e. les lignes du système correspondant à des points du bord sont nulles ... pour l'instant). Nous noterons  $-\mathbb{A}_{xy}$  cette matrice pour qu'elle corresponde à la matrice du Laplacien et non l'opposée du Laplacien.

Nous rappelons que nous avons choisi d'écrire l'équation (3.3) au point  $(x_i, y_j)$  en ligne  $k = \mathcal{F}(i, j)$  du système.

**Q. 12 (Matlab)** 1. Sans aucune vectorisation/optimisation de code, écrire la fonction `Assemblage2DBijection00` retournant la matrice creuse  $\mathbb{A}_{xy}$  du système sans tenir compte des conditions aux limites en utilisant toutes les équations de (3.3) et la fonction `bijF`. Il y aura ici deux boucles imbriquées, l'une en  $i$  et l'autre en  $j$  pour écrire chacune des équations en  $(i, j)$ .

2. Proposer au moins un programme permettant de tester/valider la matrice ainsi obtenue. ■

**Q. 13 (Matlab)** Proposer plusieurs variantes de vectorisation/optimisation de la fonction `Assemblage2DBijection00` : il faudrait supprimer au moins une des boucles en  $i$  ou en  $j$ . ■

#### 3.3.2 Méthode 2 : produits de Kronecker (facultatif)

On note  $\mathbb{I}_n$  la matrice identité de  $\mathcal{M}_n(\mathbb{R})$  et  $\mathbb{J}_n$  la matrice de  $\mathcal{M}_n(\mathbb{R})$  définie par

$$\mathbb{J}_n = \begin{pmatrix} 0 & 0 & \cdots & \cdots & 0 \\ 0 & 1 & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & 1 & 0 \\ 0 & \cdots & \cdots & 0 & 0 \end{pmatrix} \in \mathcal{M}_n(\mathbb{R}).$$

Dans cette partie nous allons générer/assembler la matrice du système (3.9) **sans tenir compte des conditions aux limites** (i.e. les lignes du système correspondant à des points du bord sont nulles ... pour l'instant).

Dans ce cas, on note  $-\mathbb{A}_{xy}$  la matrice ainsi obtenue. C'est une matrice bloc de  $\mathcal{M}_N(\mathbb{R})$  avec  $n_y$  lignes bloc composées de blocs carrés de dimension  $n_x$  qui s'écrit sous la forme :

$$\mathbb{A}_{xy} = \begin{pmatrix} 0 & 0 & \cdots & \cdots & \cdots & 0 & 0 \\ 0 & \mathbb{T}_x & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \mathbb{T}_x & \ddots & \ddots & \vdots & \vdots \\ \vdots & 0 & \ddots & \ddots & \ddots & 0 & \vdots \\ \vdots & \vdots & \ddots & \ddots & \mathbb{T}_x & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 & \mathbb{T}_x & 0 \\ 0 & 0 & \cdots & \cdots & \cdots & 0 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 & \cdots & \cdots & \cdots & 0 & 0 \\ \mathbb{S}_y & \mathbb{T}_y & \mathbb{S}_y & 0 & \cdots & 0 & 0 \\ 0 & \mathbb{S}_y & \mathbb{T}_y & \ddots & \ddots & \vdots & \vdots \\ \vdots & 0 & \ddots & \ddots & \ddots & 0 & \vdots \\ \vdots & \vdots & \ddots & \ddots & \mathbb{T}_y & \mathbb{S}_y & 0 \\ 0 & 0 & \cdots & 0 & \mathbb{S}_y & \mathbb{T}_y & \mathbb{S}_y \\ 0 & 0 & \cdots & \cdots & \cdots & 0 & 0 \end{pmatrix} \quad (3.11)$$

où

$$\mathbb{S}_y = \frac{1}{h_y^2} \mathbb{J}_{n_x}, \quad \mathbb{T}_y = -\frac{2}{h_y^2} \mathbb{J}_{n_x} \quad \text{et} \quad \mathbb{T}_x = \frac{1}{h_x^2} \begin{pmatrix} 0 & 0 & 0 & \cdots & \cdots & 0 \\ 1 & -2 & 1 & \ddots & & \vdots \\ 0 & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 \\ \vdots & & \ddots & 1 & -2 & 1 \\ 0 & \cdots & \cdots & 0 & 0 & 0 \end{pmatrix}$$

On peut noter que les matrices  $\mathbb{T}_x$ ,  $\mathbb{T}_y$  et  $\mathbb{S}_y$  sont des matrices de  $\mathcal{M}_{n_x}(\mathbb{R})$ .

En notant  $\mathbb{A}_x \in \mathcal{M}_{n_x}(\mathbb{R})$  et  $\mathbb{A}_y \in \mathcal{M}_{n_y}(\mathbb{R})$  les matrices définies par

$$\mathbb{A}_x = \frac{1}{h_x^2} \begin{pmatrix} 0 & 0 & 0 & \cdots & \cdots & 0 \\ 1 & -2 & 1 & \ddots & & \vdots \\ 0 & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 \\ \vdots & & \ddots & 1 & -2 & 1 \\ 0 & \cdots & \cdots & 0 & 0 & 0 \end{pmatrix} \quad \text{et} \quad \mathbb{A}_y = \frac{1}{h_y^2} \begin{pmatrix} 0 & 0 & 0 & \cdots & \cdots & 0 \\ 1 & -2 & 1 & \ddots & & \vdots \\ 0 & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 \\ \vdots & & \ddots & 1 & -2 & 1 \\ 0 & \cdots & \cdots & 0 & 0 & 0 \end{pmatrix}$$

on déduit de (3.11)

$$\mathbb{A}_{xy} = \mathbb{J}_{n_y} \otimes \mathbb{A}_x + \mathbb{A}_y \otimes \mathbb{J}_{n_x}. \quad (3.12)$$

**Q. 14 (Matlab)** 1. Ecrire la fonction `Lap2DAssembling` retournant la matrice bloc creuse  $\mathbb{A}$  en utilisant (3.12) et les fonctions déjà écrites.

2. Proposer au moins un programme permettant de tester/valider la matrice ainsi obtenue. ■

Sur le même principe, il est possible de généraliser en dimension quelconque la construction de la matrice du Laplacien obtenue par différences finies sans prise en compte des conditions aux limites. Par exemple en dimension 3, cette matrice s'écrit

$$\begin{aligned} \mathbb{A}_{xyz} &= \mathbb{J}_{n_z} \otimes \mathbb{A}_{xy} + (\mathbb{A}_z \otimes (\mathbb{J}_{n_y} \otimes \mathbb{J}_{n_x})) \\ &= \mathbb{J}_{n_z} \otimes \mathbb{J}_{n_y} \otimes \mathbb{A}_x + \mathbb{J}_{n_z} \otimes \mathbb{A}_y \otimes \mathbb{J}_{n_x} + \mathbb{A}_z \otimes \mathbb{J}_{n_y} \otimes \mathbb{J}_{n_x} \end{aligned}$$

## 3.4 Quelques indices

### 3.4.1 Indices des bords

La frontière  $\Gamma$  du domaine peut se décomposer en 4 sous-ensembles :

$$\Gamma = \Gamma_N \cup \Gamma_S \cup \Gamma_E \cup \Gamma_O.$$

**Q. 15 (sur feuille)** 1. Pour  $\Gamma_S$ , donner le sous-ensemble  $\mathcal{D}_S$  d'indices  $(i, j)$  dans  $\llbracket 0, N_x \rrbracket \times \llbracket 0, N_y \rrbracket$  tel que  $(x(i), y(j)) \in \Gamma_S$ .

2. En déduire le sous-ensemble  $\mathcal{I}_S$  (ordonné suivant  $\mathcal{D}_S$ ) correspondant dans la numérotation globale.

3. Même chose pour  $\Gamma_N$ ,  $\Gamma_E$  et  $\Gamma_O$ . ■

**Q. 16 (matlab)** Ecrire la fonction `BordDomaine` retournant la structure de données (`help struct`) contenant les champs suivants :

- `Sud` ayant pour valeur  $\mathcal{I}_S$ ,
- `Nord` ayant pour valeur  $\mathcal{I}_N$ ,
- `Est` ayant pour valeur  $\mathcal{I}_E$ ,
- `Ouest` ayant pour valeur  $\mathcal{I}_O$ .

Chacun des tableaux  $\mathcal{I}_x$  sera stocké sous forme ligne. ■

Une fois cette fonction écrite, il est aisé en numérotation globale de récupérer l'ensemble `idxBD` des indices des points du bords ainsi que l'ensemble `idxBDC` des indices des points strictement intérieurs (complémentaire du précédent) :

```
BD=BordDomaine (...); % ... -> a remplacer
idxBD=unique ([BD.Sud, BD.Nord, BD.Est, BD.Ouest]);
idxBDC=setdiff(1:N, idxBD); % N=nx*ny
```

### 3.4.2 Evaluation de fonctions sur une partie de la grille

Soit  $g : \Gamma_N \rightarrow \mathbb{R}$ . On suppose cette fonction vectorisée sous Matlab/Octave. Pour évaluer cette fonction aux points de la grille 2D appartenant à  $\Gamma_N$ , on utilise la structure de données retournée par la fonction `BordDomaine` et plus particulièrement le champ `Nord` de cette structure. Voici une manière de créer un vecteur `G` valant  $G(k)=g(x(i),y(j))$  sur un point du bord `Nord` avec  $k = \mathcal{F}(i, j)$  et 0 sinon.

Listing 3 – Calcul d'une fonction sur le bord Nord de la grille 2D

```
BD=BordDomaine (...); % ... -> a remplacer
[I, J]=bijRecF(1:N, nx);
X=x(I+1); Y=y(J+1); % coordonnees des points de la grille
                    % en numerotation globale
g=@(x, y) sin(x+y).*x+y.^2;
Gloc=g(X(BD.Nord), Y(BD.Nord));
G=zeros(N, 1);
G(BD.Nord)=Gloc;
```

Pour évaluer une fonction sur une autre partie de la grille, il suffit de récupérer l'ensemble des indices dans la numérotation globale des points de appartenant à cette partie et de remplacer `BD.Nord` par cet ensemble.

### 3.5 Assemblage du second membre sans conditions aux limites

Nous allons générer/assembler le vecteur second membre du système (3.9) **sans tenir compte des conditions aux limites**. Avec le choix de la numérotation globale, ce vecteur `B` de  $\mathbb{R}^N$  est défini par

$$\forall k \in \llbracket 1, N \rrbracket, \quad \mathbf{B}_k = \begin{cases} f(x(i), y(j)) & \text{si } (x(i), y(j)) \notin \Gamma \text{ et } k = \mathcal{F}(i, j) \\ 0 & \text{sinon} \end{cases}$$

En utilisant les codes donnés (et expliqués) en section 3.4, on a immédiatement un code vectorisé permettant d'initialiser ce vecteur : il est donné en Listing 4.

Listing 4 – Calcul du vecteur second membre sans conditions aux limites

```
BD=BordDomaine (...); % ... -> a remplacer
[I, J]=bijRecF(1:N, nx);
X=x(I+1); Y=y(J+1); % coordonnees des points de la grille
                    % en numerotation globale
f=@(x, y) sin(x+y).*cos(x-y);

idxBD=unique ([BD.Sud, BD.Nord, BD.Est, BD.Ouest]);
idxBDC=setdiff(1:N, idxBD); % N=nx*ny

B=zeros(N, 1);
B(idxBDC)=f(X(idxBDC), Y(idxBDC));
```

### 3.6 Prise en compte des conditions aux limites

Sans tenir compte des conditions aux limites, nous avons décrit le calcul de la matrice  $\mathbb{A}_{xy} \in \mathcal{M}_N(\mathbb{R})$  (section 3.3) et du vecteur  $\underline{\mathbf{B}} \in \mathbb{R}^N$  (section 3.5).

En reprenant les résultats et notations de **Q.11** (page 9), on en déduit que

$$\mathbb{A} = -\mathbb{A}_{xy} + \mathbb{A}_\Gamma \quad \text{et} \quad \mathbf{B} = \underline{\mathbf{B}} + \mathbf{B}_\Gamma$$

où  $\mathbb{A}_\Gamma$  et  $\mathbf{B}_\Gamma$  sont les *contributions* des conditions aux limites, c'est à dire que le système

$$\mathbb{A}_\Gamma \mathbf{U} = \mathbf{B}_\Gamma \tag{3.13}$$

contient uniquement les équations (3.4) à (3.7). Pour  $k = \mathcal{F}(i, j) \in \llbracket 1, N \rrbracket$ , si  $(x(i), y(j)) \in \Gamma$  alors la ligne  $k$  du système (3.13) correspond à

$$\mathbf{U}(k) = g(x(i), y(j)).$$

Toutes les autres lignes du système (3.13) sont nulles. On a donc pour tout  $k \in \llbracket 1, N \rrbracket$

$$\mathbb{A}_\Gamma(k, :) = \begin{cases} \mathbf{e}_k^\top & \text{si } (x(i), y(j)) \in \Gamma \text{ avec } (i, j) = \mathcal{F}^{-1}(k) \\ 0 & \text{sinon} \end{cases}$$

où  $\mathbf{e}_k^\top$  est le  $k^{\text{ième}}$  vecteur de la base canonique de  $\mathbb{R}^N$ , et

$$\mathbf{B}_\Gamma(k) = \begin{cases} g(x(i), y(j)) & \text{si } (x(i), y(j)) \in \Gamma \text{ avec } (i, j) = \mathcal{F}^{-1}(k) \\ 0 & \text{sinon} \end{cases}$$

En utilisant les codes donnés (et expliqués) en section 3.4, on calcule dans le Listing 5, le vecteur **Bcl** correspondant à  $\mathbf{B}_\Gamma$  et la matrice **Acl** correspondant à  $\mathbb{A}_\Gamma$ .

Listing 5 – Calcul des contributions du bord

```
BD=BordDomaine (...); % ... -> a remplacer
[I, J]=bijRecF(1:N, nx);
X=x(I+1); Y=y(J+1); % coordonnees des points de la grille
% en numerotation globale
g=@(x, y) sin(x+y) .* cos(x-y);

idxBD=unique([BD.Sud, BD.Nord, BD.Est, BD.Ouest]);
idxBDc=setdiff(1:N, idxBD); % N=nx*ny

Bcl=zeros(N, 1);
Bcl(idxBD)=g(X(idxBD), Y(idxBD));
Acl=sparse(idxBD, idxBD, ones(1, length(idxBD)), N, N);
```

### 3.7 Résolution du système

On doit résoudre le système  $\mathbb{A}\mathbf{U} = \mathbf{B}$  défini en **Q.11** (page 9). Or on a construit les matrices  $\mathbb{A}_{xy}$ ,  $\mathbb{A}_\Gamma$  et les vecteurs  $\underline{\mathbf{B}}$ ,  $\mathbf{B}_\Gamma$  de telle sorte que

$$\mathbb{A} = -\mathbb{A}_{xy} + \mathbb{A}_\Gamma \quad \text{et} \quad \mathbf{B} = \underline{\mathbf{B}} + \mathbf{B}_\Gamma$$

#### 3.7.1 Méthode 1

Dans cette section nous allons directement résoudre le système  $\mathbb{A}\mathbf{U} = \mathbf{B}$ .

**Q. 17** 1. *Ecrire le programme **EDP2D01** permettant de résoudre numériquement l'EDP (3.1)-(3.2) à l'aide du schéma discrétisé (3.3) à (3.7). On choisira judicieusement les données permettant sur une même figure de représenter de la solution numérique et de la solution exacte, et sur une autre figure de représenter l'erreur entre les deux solutions.*

2. *Ecrire le programme **ordreEDP2D01** permettant de retrouver numériquement l'ordre du schéma utilisé. ■*

### 3.7.2 Méthode 2

Dans cette section nous proposons une autre méthode permettant de résoudre le système  $\mathbb{A}\mathbf{U} = \mathbf{B}$  en éliminant les équations portants sur les points Dirichlet.

Pour cela, on note

$$\mathcal{I}_D = \{k \in \llbracket 1, N \rrbracket; (x(i), y(j)) \in \Gamma \text{ avec } (i, j) = \mathcal{F}^{-1}(k)\}$$

et son complémentaire

$$\mathcal{I}_D^c = \llbracket 1, N \rrbracket \setminus \mathcal{I}_D.$$

Le système linéaire à résoudre  $\mathbb{A}\mathbf{U} = \mathbf{B}$  peut se décomposer de manière équivalente en

$$\mathbb{A}\mathbf{U} = \mathbf{B} \Leftrightarrow \begin{cases} \mathbb{A}(\mathcal{I}_D, :)\mathbf{U} = \mathbf{B}(\mathcal{I}_D) \\ \mathbb{A}(\mathcal{I}_D^c, :)\mathbf{U} = \mathbf{B}(\mathcal{I}_D^c) \end{cases} \quad (3.14a)$$

$$(3.14b)$$

Nous allons réécrire les deux systèmes (3.14a) et (3.14b).

- Pour (3.14a), nous avons par construction on a  $\mathbb{A}(\mathcal{I}_D, :) = \mathbb{A}_\Gamma(\mathcal{I}_D, :)$  et donc

$$\begin{aligned} \mathbb{A}(\mathcal{I}_D, :)\mathbf{U} &= \mathbb{A}_\Gamma(\mathcal{I}_D, :)\mathbf{U} \\ &= \mathbb{A}_\Gamma(\mathcal{I}_D, \mathcal{I}_D)\mathbf{U}(\mathcal{I}_D) + \mathbb{A}_\Gamma(\mathcal{I}_D, \mathcal{I}_D^c)\mathbf{U}(\mathcal{I}_D^c) \end{aligned}$$

Comme  $\mathbb{A}_\Gamma(\mathcal{I}_D, \mathcal{I}_D^c) = \mathbb{0}$  et  $\mathbb{A}_\Gamma(\mathcal{I}_D, \mathcal{I}_D) = \mathbb{1}$  on obtient

$$\mathbb{A}(\mathcal{I}_D, :)\mathbf{U} = \mathbf{U}(\mathcal{I}_D).$$

Donc (3.14a) est équivalent à

$$\mathbf{U}(\mathcal{I}_D) = \mathbf{B}(\mathcal{I}_D). \quad (3.15)$$

- Pour (3.14b), nous avons par construction  $\mathbb{A}(\mathcal{I}_D^c, :) = -\mathbb{A}_{x,y}(\mathcal{I}_D^c, :)$  et donc

$$\begin{aligned} \mathbb{A}(\mathcal{I}_D^c, :)\mathbf{U} &= -\mathbb{A}_{x,y}(\mathcal{I}_D^c, :)\mathbf{U} \\ &= -\mathbb{A}_{x,y}(\mathcal{I}_D^c, \mathcal{I}_D)\mathbf{U}(\mathcal{I}_D) - \mathbb{A}_{x,y}(\mathcal{I}_D^c, \mathcal{I}_D^c)\mathbf{U}(\mathcal{I}_D^c) \end{aligned}$$

En utilisant (3.15), le système (3.14b) s'écrit

$$-\mathbb{A}_{x,y}(\mathcal{I}_D^c, \mathcal{I}_D)\mathbf{U}(\mathcal{I}_D) = \mathbf{B}(\mathcal{I}_D^c) + \mathbb{A}_{x,y}(\mathcal{I}_D^c, \mathcal{I}_D^c)\mathbf{B}(\mathcal{I}_D). \quad (3.16)$$

Résoudre le système linéaire  $\mathbb{A}\mathbf{U} = \mathbf{B}$  est alors équivalent à calculer  $\mathbf{U}(\mathcal{I}_D)$  par (3.15) et déterminer  $\mathbf{U}(\mathcal{I}_D^c)$  en résolvant le système linéaire (3.16).

**Q. 18** 1. *Ecrire le programme `EDP2D02` permettant de résoudre numériquement l'EDP (3.1)-(3.2) à l'aide du schéma discrétisé (3.3) à (3.7) que l'on résoudra par (3.15) et (3.16). On choisira judicieusement les données permettant sur une même figure de représenter de la solution numérique et de la solution exacte, et sur une autre figure de représenter l'erreur entre les deux solutions.*

2. *Ecrire le programme `ordreEDP2D02` permettant de retrouver numériquement l'ordre du schéma utilisé. ■*