

---

*Volumes Finis*

J.-F. Scheid

---



# Table des matières

<b>1</b>	<b>Equations elliptiques</b>	<b>5</b>
1.1	Introduction . . . . .	5
1.2	Quelques rappels sur les solutions d'équations elliptiques linéaires . . . . .	5
1.2.1	Existence, unicité et régularité des solutions . . . . .	6
1.2.2	Principes du maximum . . . . .	7
1.3	Volumes Finis pour les problèmes elliptiques 1D . . . . .	7
1.3.1	Maillage . . . . .	7
1.3.2	Formulation en Volumes Finis . . . . .	9
1.3.3	Système linéaire . . . . .	10
1.3.4	Convergence . . . . .	10
1.3.5	Equation elliptique 1D à coefficients discontinus . . . . .	10
1.4	Volumes Finis pour les problèmes elliptiques 2D . . . . .	14
1.4.1	Maillage . . . . .	14
1.4.2	Formulation en Volumes Finis . . . . .	15
1.4.3	Exemples de maillages admissibles . . . . .	17
1.4.4	Système linéaire . . . . .	18
1.4.5	Estimations d'erreurs . . . . .	20
1.4.6	Equation elliptique 2D avec coefficients discontinus . . . . .	21
<b>2</b>	<b>Equations paraboliques</b>	<b>25</b>
2.1	Introduction . . . . .	25
2.1.1	Existence et unicité des solutions . . . . .	25
2.1.2	Principes du maximum . . . . .	26
2.2	Volumes Finis pour l'équation de la chaleur en dimension 1 d'espace . . . . .	26
2.2.1	Schéma d'Euler explicite . . . . .	27
2.2.2	Schéma d'Euler implicite . . . . .	29
2.3	Equation de la chaleur en 2D d'espace et Volumes Finis . . . . .	34
2.3.1	Discretisation en espace . . . . .	34
2.3.2	Schéma explicite en temps . . . . .	35
2.3.3	Schéma implicite en temps . . . . .	36
<b>3</b>	<b>Equation de transport</b>	<b>37</b>
3.1	Introduction . . . . .	37
3.2	Maillage . . . . .	38
3.3	Formulation en Volumes Finis . . . . .	39
3.4	Système linéaire . . . . .	41
3.5	Condition de stabilité . . . . .	43
<b>4</b>	<b>Equations de Stokes</b>	<b>45</b>
4.1	Introduction . . . . .	45
4.2	Maillages . . . . .	45
4.3	Formulation en Volumes Finis . . . . .	45

4.3.1	Approximation de la divergence . . . . .	46
4.3.2	Approximation de l'équation de Stokes . . . . .	47
4.3.3	Système linéaire . . . . .	48
<b>5</b>	<b>Equations de Navier-Stokes incompressibles</b>	<b>49</b>
5.1	Introduction . . . . .	49
5.2	Semi-discrétisation en temps . . . . .	49
5.3	Formulations en Volumes Finis . . . . .	50
	<b>Appendices</b>	<b>51</b>
<b>A</b>	<b>Modélisation des équations de Navier-Stokes et équations de Stokes</b>	<b>53</b>
A.1	Introduction . . . . .	53
A.2	Conservation de la masse . . . . .	53
A.3	Loi fondamentale de la dynamique - loi de comportement . . . . .	54
A.4	Conservation du volume . . . . .	55
A.5	Adimensionalisation des équations de Navier-Stokes . . . . .	55
A.6	Réductions des équations . . . . .	56
<b>B</b>	<b>Eléments d'analyse matricielle.</b>	<b>59</b>
B.1	Matrice réductible . . . . .	59
B.2	Matrice à diagonale dominante . . . . .	60
B.3	Matrice monotone . . . . .	60
B.4	Localisation des valeurs propres . . . . .	62
<b>C</b>	<b>Quelques inégalités.</b>	<b>65</b>
C.1	Inégalité de Cauchy-Schwarz . . . . .	65
C.2	Inégalité de Young . . . . .	65
C.3	Inégalité de Gronwall . . . . .	65
C.4	Inégalité de Gronwall discrète . . . . .	65
	<b>Références</b>	<b>67</b>

# Chapitre 1

## Equations elliptiques

### 1.1 Introduction

Les méthodes de Volumes Finis sont construites à partir d'une formulation intégrale basée directement sur la forme *forte* des équations à résoudre. Les intégrales ne portent pas sur tout le domaine dans lequel sont posées les équations, mais sur des cellules disjointes appelées volumes de contrôles. En comparaison, la méthode des Éléments Finis s'appuie également sur une formulation intégrale des équations, appelée *formulation variationnelle* (ou encore formulation faible) faisant intervenir des "fonctions tests" et où les intégrales portent sur tout le domaine. Dans la méthode des Volumes Finis, les termes de divergence apparaissant dans les EDP à résoudre sont traités en utilisant le théorème de la divergence. Ainsi, les intégrales de volume d'un terme de divergence sont transformées en intégrales de surface. Ces termes de flux sont ensuite évalués aux interfaces entre les volumes de contrôle et les flux aux interfaces sont approchés par une fonction de flux numérique.

Les méthodes de Volumes Finis ont été initialement développées et mises au point pour des lois de conservation hyperboliques. Les méthodes de Volumes Finis sont conservatives car on impose que le flux entrant dans un volume de contrôle soit égal au flux sortant du volume adjacent. Ces méthodes sont par conséquent très bien adaptées à la résolution de lois de conservation. Le développement pour des équations elliptiques et paraboliques est plus récent. Un avantage de la méthode des Volumes Finis par rapport à la méthode des Différences Finies est qu'elle permet de résoudre des EDP avec des géométries complexes dans la mesure où elle utilise des maillages non-structurés. Dans la méthode des Volumes Finis, le domaine (supposé polygonal) est discrétisé par un maillage constitué de volumes de contrôles qui sont des (petits) volumes disjoints en 3D, des polygones en 2D, des segments en 1D. Les volumes de contrôles peuvent être construits autour des points d'un maillage initial par tétraédrisation/triangulation.

Dans ce chapitre, on commence par présenter la méthode des Volumes Finis pour des équations elliptiques (1D puis 2D) et on s'intéressera ensuite à l'équation de transport 2D. On développera une méthode de Volumes Finis pour les équations de Stokes et enfin on résoudra les équations de Navier-Stokes incompressibles par un schéma semi-implicite en temps et une formulation Volumes Finis *upwind* pour le terme de convection linéarisé.

### 1.2 Quelques rappels sur les solutions d'équations elliptiques linéaires

Dans cette section on rappelle quelques propriétés des solutions d'équations elliptiques linéaires. Pour un domaine (ouvert connexe)  $\Omega \subset \mathbb{R}^n$  borné et de frontière  $\partial\Omega$  *régulière*, on considère le problème aux limites suivant pour une fonction  $u$  :

$$\mathcal{L}u := -\operatorname{div}(A\nabla u) + cu = f \quad \text{dans } \Omega \quad (1.1)$$

$$u = 0 \quad \text{sur } \partial\Omega \quad (1.2)$$

où  $A$  est une matrice de taille  $n \times n$  avec  $A = A(\mathbf{x}) = (a_{ij}(\mathbf{x}))_{1 \leq i, j \leq n}$  pour  $\mathbf{x} \in \Omega$  et on suppose que les coefficients  $a_{ij} \in C^1(\bar{\Omega})$ . Les fonctions  $f$  et  $c$  sont données et on suppose que la fonction  $c \in C(\bar{\Omega})$  est

**positive ou nulle** i.e.  $c(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \Omega$ . Les opérateurs différentiels apparaissant dans (1.1) sont définis par

$$\nabla u = \text{grad } u = \left( \frac{\partial u}{\partial x_1}, \dots, \frac{\partial u}{\partial x_n} \right), \quad \text{div } \mathbf{u} = \sum_{i=1}^n \frac{\partial u_i}{\partial x_i},$$

pour  $\mathbf{x} = (x_1, \dots, x_n)$  et  $\mathbf{u} = (u_1, \dots, u_n)$ . L'opérateur  $\mathcal{L}$  apparaît sous forme divergente et on doit lire

$$\text{div}(A\nabla u) = \sum_{i,j=1}^n \frac{\partial}{\partial x_i} \left( a_{ij} \frac{\partial u}{\partial x_j} \right).$$

En particulier, avec la matrice  $A = I_d$  on obtient  $\mathcal{L} = \Delta$  l'opérateur laplacien défini par

$$\Delta u = \sum_{i=1}^n \frac{\partial^2 u}{\partial x_i^2}.$$

On dit que l'équation (1.1) est (uniformément) **elliptique** ou bien que l'opérateur  $\mathcal{L}$  est elliptique si

$$\exists \alpha > 0 \text{ tel que } \langle A(\mathbf{x})\boldsymbol{\xi}, \boldsymbol{\xi} \rangle_{\mathbb{R}^n} = \sum_{i,j=1}^n a_{ij}(\mathbf{x})\xi_i\xi_j \geq \alpha \|\boldsymbol{\xi}\|^2, \quad \forall \boldsymbol{\xi} \in \mathbb{R}^n, \forall \mathbf{x} \in \Omega, \quad (1.3)$$

où  $\langle \cdot, \cdot \rangle_{\mathbb{R}^n}$  désigne le produit scalaire dans  $\mathbb{R}^n$ . On supposera désormais que la condition d'ellipticité (1.3) est vérifiée.

### 1.2.1 Existence, unicité et régularité des solutions

Pour les quelques résultats d'existence, d'unicité et de principe du maximum rappelés ci-après, on pourra consulter par exemple [2],[5],[15],[12].

Commençons par le cas de la dimension 1 avec  $\Omega = (0, 1)$ . Dans ce cas, le problème s'écrit

$$-(au')' + cu = f \quad \text{dans } (0, 1) \quad (1.4)$$

$$u(0) = u(1) = 0. \quad (1.5)$$

- Si  $f, c \in C([0, 1])$  avec  $c \geq 0$  et si  $a \in C^1([0, 1])$  avec  $a(x) \geq \alpha > 0$  pour tout  $x \in (0, 1)$ <sup>(1)</sup>, alors le problème (1.4)-(1.5) admet une unique solution *classique*  $u \in C^2([0, 1])$ .

Dans le cas de la dimension quelconque avec  $\Omega \subset \mathbb{R}^n$  un domaine borné et *régulier*, on a le résultat suivant d'existence de solution *faible*<sup>(2)</sup>.

- Si  $f \in L^2(\Omega)$ ,  $c \in C(\bar{\Omega})$  avec  $c \geq 0$  et si les coefficients  $a_{ij} \in C(\bar{\Omega})$  vérifient (1.3), alors le problème (1.1)-(1.2) admet une unique solution *faible*  $u \in H_0^1(\Omega)$  (Lax-Milgram) et  $\|u\|_{H^1} \leq C\|f\|_{L^2}$  où  $C$  est une constante indépendante de  $u$  et  $f$ .  
Si de plus,  $a_{ij} \in C^1(\bar{\Omega})$  alors  $u \in H^2(\Omega)$ <sup>(3)</sup> et  $u$  vérifie les équations (1.1)-(1.2) au sens *presque partout*.

Donnons enfin un résultat de régularité dans les espaces de Hölder. Pour  $0 < \alpha < 1$ , on définit les espaces

$$C^{0,\alpha}(\bar{\Omega}) = \left\{ u \in C^0(\bar{\Omega}), \sup_{\mathbf{x} \neq \mathbf{y}} \frac{|u(\mathbf{x}) - u(\mathbf{y})|}{|\mathbf{x} - \mathbf{y}|^\alpha} < \infty \right\}$$

et pour  $m \in \mathbb{N}$ ,

$$C^{m,\alpha}(\bar{\Omega}) = \left\{ u \in C^m(\bar{\Omega}), D^\beta u \in C^{0,\alpha}(\bar{\Omega}) \text{ avec } |\beta| = m \right\},$$

où  $\beta$  est un multi-indice  $\beta = (\beta_1, \dots, \beta_n)$  et  $D^\beta = \frac{\partial^{|\beta|}}{\partial x_1^{\beta_1} \dots \partial x_n^{\beta_n}}$  avec  $|\beta| = \beta_1 + \beta_2 + \dots + \beta_n$ .

---

1. Cette condition ne traduit rien d'autre que la condition d'ellipticité (1.3) dans le cas 1D.  
2.  $u$  est solution faible de (1.1),(1.2) si  $u \in H_0^1(\Omega)$  et vérifie  $\int_\Omega A\nabla u \cdot \nabla v + \int_\Omega c u v = \int_\Omega f v, \quad \forall v \in H_0^1(\Omega)$ .  
3. On a supposé  $\Omega$  régulier ; en général, on n'a pas la régularité  $u \in H^2(\Omega)$  si  $\Omega$  présente des coins rentrants par exemple.

- Si  $a_{ij} \in C^{k+1,\alpha}(\overline{\Omega})$  vérifient (1.3) et  $f, c \in C^{k,\alpha}(\overline{\Omega})$  avec  $c \geq 0$ , alors il existe une unique solution  $u \in C^{k+2,\alpha}(\overline{\Omega})$ .

Ces résultats d'existence indiquent un effet "régularisant" d'un opérateur elliptique, par rapport aux données. Par exemple pour l'opérateur Laplacien ( $A = I_d$ ), si on prend un second membre  $f$  aléatoire (entre -100 et 100) sur le pavé unité  $\Omega = (0,1) \times (0,1)$  on obtient une solution "plus régulière", comme l'illustre la figure 1.1. Si à présent on prend cette solution comme second membre de l'équation de Laplace, on obtient une solution encore plus régulière (cf. Fig. 1.1).

*Remarques sur la régularité des solutions.*

L'étude de la convergence d'approximation de solution doit tenir compte de la régularité de la solution exacte. Si la solution est faible ( $L^2, H^1, \dots$ ), il faut regarder la convergence dans  $L^2, H^1, \dots$ ; si la solution est régulière ( $C^2, \dots$ ), alors on peut regarder la convergence dans  $L^2, H^1, H^2, L^\infty, C^1, C^2, \dots$ .

### 1.2.2 Principes du maximum

On donne à présent quelques résultats de principe du maximum pour des solutions classiques.

On suppose que  $a_{ij} \in C^1(\overline{\Omega})$  vérifient (1.3) et  $c \in C(\overline{\Omega})$  avec  $c \geq 0$ . Soit  $u \in C^2(\Omega) \cap C^0(\overline{\Omega})$  vérifiant  $\mathcal{L}u = f$  dans  $\Omega$ .

- i)* • On prend  $c \equiv 0$ . Si  $f \leq 0$  (resp.  $f \geq 0$ ) alors  $u$  atteint son maximum (resp. minimum) sur  $\partial\Omega$ .
- ii)* • Si  $c \equiv 0$  et  $f \equiv 0$  alors  $\inf_{\partial\Omega} u \leq u \leq \sup_{\partial\Omega} u$ . Ceci montre qu'avec  $c \equiv 0$  et  $f \equiv 0$ , si on choisit en particulier  $u|_{\partial\Omega} = 0$  alors on obtient  $u \equiv 0$ .
- iii)* • Si  $f \geq 0$  et  $u|_{\partial\Omega} \geq 0$  alors  $u \geq 0$  dans  $\Omega$ .
- iv)* • (PRINCIPE DE HOPF) Soit  $f \leq 0$  et soit  $\mathbf{x}_0 \in \partial\Omega$  tel que  $u(\mathbf{x}_0) > u(\mathbf{x}), \forall \mathbf{x} \in \Omega$ . On suppose que  $u$  est dérivable en  $\mathbf{x}_0$ . Si  $c \equiv 0$  ou bien si  $u(\mathbf{x}_0) = 0$ , alors

$$\frac{\partial u}{\partial \mathbf{n}}(\mathbf{x}_0) > 0,$$

où  $\mathbf{n}$  désigne la normale extérieure à  $\partial\Omega$ .

## 1.3 Volumes Finis pour les problèmes elliptiques 1D

On présente la méthode des Volumes Finis pour les équations elliptiques en 1D, en considérant le problème modèle suivant

$$(P) \begin{cases} -u''(x) = f(x), & x \in (0,1) \\ u(0) = u(1) = 0, \end{cases}$$

où  $f$  est une fonction donnée. A la fin de la section, on applique la méthode au cas d'une équation elliptique 1D avec des coefficients discontinus.

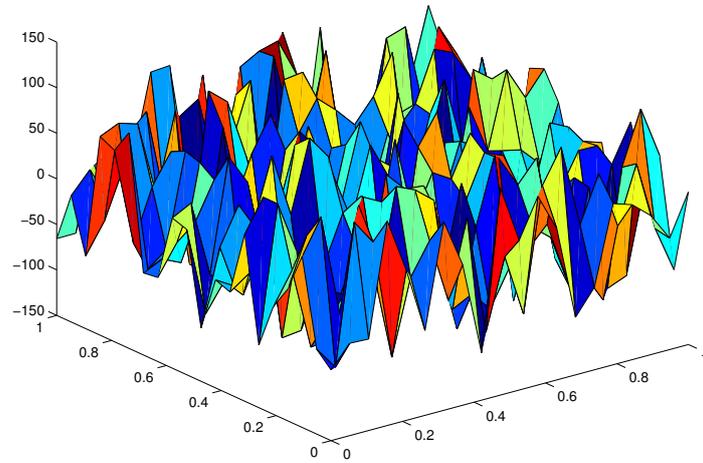
### 1.3.1 Maillage

On discrétise l'intervalle  $[0,1]$  en introduisant un maillage  $\mathcal{T}$  de l'intervalle  $[0,1]$  défini de la façon suivante :

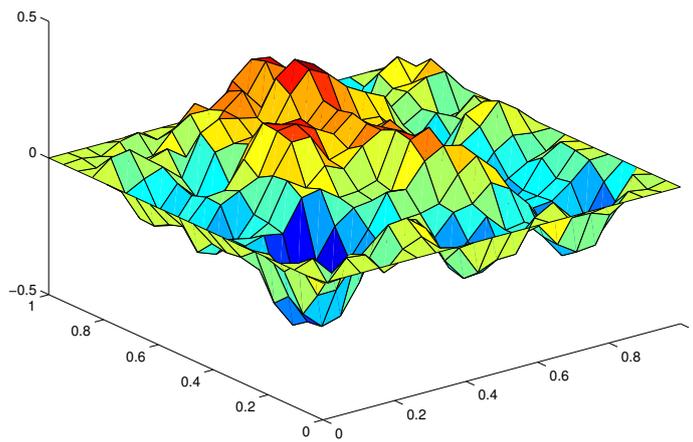
— Soient  $N$  volumes de contrôle appelés aussi *cellules*, notés  $K_i$  pour  $i = 1, \dots, N$  :

$$K_i = ]x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}[$$

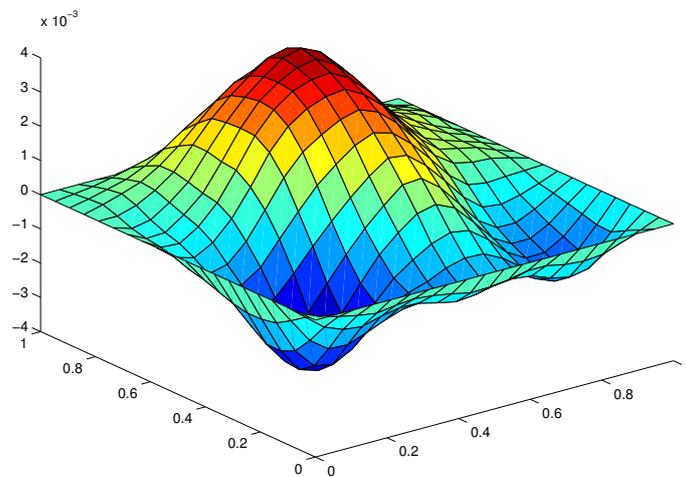
avec les points  $x_{i+\frac{1}{2}} \in [0,1]$  tels que  $0 = x_{\frac{1}{2}} < x_{\frac{3}{2}} < \dots < x_{N-\frac{1}{2}} < x_{N+\frac{1}{2}} = 1$ .



Terme source initial  $f$ .



Solution de  $-\Delta u_1 = f$ .



Solution de  $-\Delta u_2 = u_1$ .

FIGURE 1.1 – Effet “régularisant” du Laplacien.

— A chaque cellule  $K_i$ , on associe un point (centre)  $x_i \in K_i$  tel que :

$$0 = x_0 = x_{\frac{1}{2}} < x_1 < \cdots < x_{i-\frac{1}{2}} < x_i < x_{i+\frac{1}{2}} < x_{i+1} < \cdots < x_{N+\frac{1}{2}} = x_{N+1} = 1.$$

On introduit alors les pas de discrétisation

$$\begin{aligned} h_i &= |K_i| = x_{i+\frac{1}{2}} - x_{i-\frac{1}{2}} \\ h_{i+\frac{1}{2}} &= x_{i+1} - x_i \end{aligned}$$

### 1.3.2 Formulation en Volumes Finis

On considère les approximations  $u_i$  de la solution  $u$  de (P) dans chaque cellule  $K_i$ . On a donc  $N$  inconnues. Plus précisément,  $u_i$  est une approximation de la valeur moyenne de  $u$  dans  $K_i$  :

$$u_i \simeq \frac{1}{|K_i|} \int_{K_i} u(x) dx = \frac{1}{h_i} \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} u(x) dx \quad \text{pour } 1 \leq i \leq N.$$

On intègre alors l'équation différentielle de (P) sur chaque cellule  $K_i$ . On a

$$- \int_{K_i} u''(x) dx = \int_{K_i} f(x) dx,$$

ce qui donne

$$-u'(x_{i+\frac{1}{2}}) + u'(x_{i-\frac{1}{2}}) = h_i f_i \quad (1.6)$$

où  $f_i$  désigne la valeur moyenne de  $f$  dans  $K_i$ , i.e.  $f_i = \frac{1}{h_i} \int_{K_i} f(x) dx$ .

La quantité  $-u'(x_{i+\frac{1}{2}})$  (resp.  $-u'(x_{i-\frac{1}{2}})$ ) représente le *flux* rentrant (resp. *flux* sortant) associé à la cellule  $K_i$ , au point  $x = x_{i+\frac{1}{2}}$  (resp. en  $x = x_{i-\frac{1}{2}}$ ). On approche le flux  $-u'(x_{i+\frac{1}{2}})$  par différences décentrées :

$$-u'(x_{i+\frac{1}{2}}) \simeq -\frac{u(x_{i+\frac{1}{2}}) - u(x_i)}{x_{i+\frac{1}{2}} - x_i} \quad \text{ou} \quad -u'(x_{i+\frac{1}{2}}) \simeq -\frac{u(x_{i+\frac{1}{2}}) - u(x_{i+1})}{x_{i+\frac{1}{2}} - x_{i+1}}. \quad (1.7)$$

Les approximations (1.7) traduisent la **consistance des flux** numériques. Du fait des décentremets, il s'agit d'approximations d'ordre  $\mathcal{O}(h)$  avec  $h = \max(h_i)$ . Au point  $x = x_{i+\frac{1}{2}}$ , on introduit les *flux numériques*  $F_{i+\frac{1}{2}}^-$  associé à la cellule  $K_i$  et  $F_{i+\frac{1}{2}}^+$  associé à la cellule  $K_{i+1}$  :

$$F_{i+\frac{1}{2}}^- = -\frac{u_{i+\frac{1}{2}} - u_i}{x_{i+\frac{1}{2}} - x_i}, \quad F_{i+\frac{1}{2}}^+ = -\frac{u_{i+\frac{1}{2}} - u_{i+1}}{x_{i+\frac{1}{2}} - x_{i+1}}. \quad (1.8)$$

On impose alors la **conservation des flux** numériques à travers le point  $x = x_{i+\frac{1}{2}}$  :

$$F_{i+\frac{1}{2}}^- = F_{i+\frac{1}{2}}^+ \quad (1.9)$$

Cette condition correspond à la continuité du flux (exacte)  $-u'$  en  $x = x_{i+\frac{1}{2}}$ . En combinant (1.8) avec (1.9), on obtient

$$F_{i+\frac{1}{2}}^- = F_{i+\frac{1}{2}}^+ = -\frac{u_{i+1} - u_i}{x_{i+1} - x_i} = -\frac{u_{i+1} - u_i}{h_{i+\frac{1}{2}}}. \quad (1.10)$$

Le schéma numérique correspondant à l'approximation de (1.6) par (1.7) avec (1.10), s'écrit :

$$-\frac{(u_{i+1} - u_i)}{h_{i+\frac{1}{2}}} + \frac{(u_i - u_{i-1})}{h_{i-\frac{1}{2}}} = h_i f_i \quad \text{pour } 1 \leq i \leq N, \quad (1.11)$$

et on fixe

$$u_0 = u_{N+1} = 0. \quad (1.12)$$

### 1.3.3 Système linéaire

On regroupe les  $N$  inconnues dans le vecteur  $\mathbf{u} = (u_1, \dots, u_N)^\top$  et on note  $\mathbf{b} = (b_1, \dots, b_N)^\top$  avec  $b_i = h_i f_i$ . Soit  $A$  la matrice de taille  $N \times N$  définie par

$$A = \begin{pmatrix} \alpha_1 & \beta_1 & & & 0 \\ \beta_1 & \alpha_2 & \beta_2 & & \\ & \ddots & \ddots & \ddots & \\ & & \beta_{N-2} & \alpha_{N-1} & \beta_{N-1} \\ 0 & & & \beta_{N-1} & \alpha_N \end{pmatrix} \quad (1.13)$$

avec

$$\alpha_i = \frac{1}{h_{i-\frac{1}{2}}} + \frac{1}{h_{i+\frac{1}{2}}} > 0, \quad \beta_i = -\frac{1}{h_{i+\frac{1}{2}}} < 0.$$

Les relations (1.11),(1.12) s'écrivent

$$A\mathbf{u} = \mathbf{b}. \quad (1.14)$$

On remarque que  $\beta_{i-1} + \alpha_i + \beta_i = 0$  pour  $2 \leq i \leq N-1$  et que  $\alpha_1 > \beta_1$ ,  $\alpha_N > \beta_N$ . La matrice  $A$  est irréductible et à diagonale fortement dominante, elle est donc inversible (cf. Annexe B).

### 1.3.4 Convergence

On a le résultat de convergence suivant (cf. exercices).

**Théorème 1.1** *Soit  $f \in C([0,1])$  et  $u \in C^2([0,1])$  l'unique solution de (P). Soit  $\mathbf{u}$  l'unique solution du schéma "Volumes Finis" (1.14). On note l'erreur  $e_i = u(x_i) - u_i$  pour  $1 \leq i \leq N$  avec  $\mathbf{e} = (e_0, e_1, \dots, e_N, e_{N+1})$  où  $e_0 = e_{N+1} = 0$ . Alors il existe une constante  $C \geq 0$  dépendant de  $u$  mais indépendant de  $h = \max_i(h_{i+\frac{1}{2}})$ , telle que*

$$\|\mathbf{e}\|_\infty := \max_{1 \leq i \leq N} |e_i| \leq Ch, \quad (1.15)$$

$$\|\mathbf{e}\|_{1,h} := \left( \sum_{i=0}^N \frac{(e_{i+1} - e_i)^2}{h_{i+\frac{1}{2}}} \right)^{1/2} \leq Ch. \quad (1.16)$$

Dans l'estimation (1.16),  $\|\cdot\|_{1,h}$  représente une norme  $H_0^1$  discrète. L'approximation en Volumes Finis est donc (au moins) d'ordre  $\mathcal{O}(h)$ .

### 1.3.5 Equation elliptique 1D à coefficients discontinus

On décrit à présent la méthode des Volumes Finis pour un problème elliptique 1D avec des coefficients variables et discontinus. Les points de discontinuité des coefficients coïncident avec les extrémités des cellules du maillage. Pour une fonction  $f$  régulière dans  $[0, 1]$ , on considère le problème suivant

$$(P) \begin{cases} -(\beta u)' = f & \text{dans } [0, 1] \\ u(0) = u(1) = 0. \end{cases} \quad (1.17)$$

La fonction  $\beta$  est définie dans l'intervalle  $[0, 1]$  et constante par morceaux :

$$\beta(x) = \begin{cases} \beta^+ & \text{si } x > \alpha \\ \beta^- & \text{si } x < \alpha, \end{cases} \quad (1.18)$$

où  $\alpha \in ]0, 1[$  et  $\beta^+, \beta^-$  sont deux constantes strictement positives. La solution (faible)  $u \in H_0^1(0, 1)$  du problème (1.17) est continue dans  $[0, 1]$ , régulière dans chaque sous-intervalle  $[0, \alpha[$ ,  $]\alpha, 1]$  et elle vérifie

$$-\beta u'' = f \quad \text{dans } [0, \alpha[ \cup ]\alpha, 1] \quad (1.19)$$

$$u(\alpha^+) = u(\alpha^-) \quad (1.20)$$

$$[\beta u'] \equiv \beta^+ u'(\alpha^+) - \beta^- u'(\alpha^-) = 0 \quad (1.21)$$

où on a noté  $v(a^-) = \lim_{x \rightarrow a^-} v(x)$  et  $v(a^+) = \lim_{x \rightarrow a^+} v(x)$ .

On utilise le maillage de l'intervalle  $[0, 1]$  précédemment défini (cf. Section 1.3.1) avec  $h_i = |K_i| = x_{i+\frac{1}{2}} - x_{i-\frac{1}{2}}$ . On notera aussi

$$h_{i+\frac{1}{2}}^- = x_{i+\frac{1}{2}} - x_i > 0, \quad h_{i+\frac{1}{2}}^+ = x_{i+1} - x_{i+\frac{1}{2}} > 0.$$

On suppose que le maillage coïncide avec la discontinuité de la fonction  $\beta$  c'est-à-dire qu'il existe  $k$  tel que :

$$\alpha = x_{k+\frac{1}{2}}. \quad (1.22)$$

En intégrant l'équation différentielle (1.19) sur la cellule  $K_i$  (y compris la cellule  $K_k$ ), on obtient après intégration par parties et en utilisant la relation (1.21) pour la cellule  $K_k$  :

$$-\beta(x_{i+\frac{1}{2}}^-)u'(x_{i+\frac{1}{2}}^-) + \beta(x_{i-\frac{1}{2}}^+)u'(x_{i-\frac{1}{2}}^+) = h_i f_i \quad (1.23)$$

où  $f_i = \frac{1}{|K_i|} \int_{K_i} f(\mathbf{x}) d\mathbf{x}$ , pour tout  $i = 1, \dots, N$ . Soient  $F_{i+\frac{1}{2}}^-$  et  $F_{i-\frac{1}{2}}^+$  les flux numériques associés à la cellule  $K_i$  obtenus en approchant respectivement les flux  $-\beta(x_{i+\frac{1}{2}}^-)u'(x_{i+\frac{1}{2}}^-)$  et  $-\beta(x_{i-\frac{1}{2}}^+)u'(x_{i-\frac{1}{2}}^+)$  par différences décentrées (**consistance des flux**) :

$$F_{i+\frac{1}{2}}^- = -\beta(x_{i+\frac{1}{2}}^-) \frac{(u_{i+\frac{1}{2}} - u_i)}{h_{i+\frac{1}{2}}^-}, \quad F_{i+\frac{1}{2}}^+ = -\beta(x_{i+\frac{1}{2}}^+) \frac{(u_{i+1} - u_{i+\frac{1}{2}})}{h_{i+\frac{1}{2}}^+}. \quad (1.24)$$

Le schéma "Volumes Finis" s'écrit

$$F_{i+\frac{1}{2}}^- - F_{i-\frac{1}{2}}^+ = h_i f_i, \quad 1 \leq i \leq N.$$

On impose la **conservation des flux numériques** à travers les points  $x_{i+\frac{1}{2}}$  :

$$F_{i+\frac{1}{2}}^- = F_{i+\frac{1}{2}}^+$$

ce qui donne

$$-\beta(x_{i+\frac{1}{2}}^-) \frac{(u_{i+\frac{1}{2}} - u_i)}{h_{i+\frac{1}{2}}^-} = -\beta(x_{i+\frac{1}{2}}^+) \frac{(u_{i+1} - u_{i+\frac{1}{2}})}{h_{i+\frac{1}{2}}^+}. \quad (1.25)$$

qui est l'analogie discret de (1.21). On note  $F_{i+\frac{1}{2}} = F_{i+\frac{1}{2}}^-$  et le schéma "Volumes Finis" s'écrit

$$F_{i+\frac{1}{2}} - F_{i-\frac{1}{2}} = h_i f_i, \quad 1 \leq i \leq N. \quad (1.26)$$

La relation (1.25) permet d'éliminer  $u_{i+\frac{1}{2}}$  dans l'expression de  $F_{i+\frac{1}{2}}$ . On obtient

$$F_{i+\frac{1}{2}} = -\beta_{i+\frac{1}{2}}^* \left( \frac{u_{i+1} - u_i}{x_{i+1} - x_i} \right), \quad (1.27)$$

avec  $\beta_{i+\frac{1}{2}}^* = \beta^+$  pour  $i > k$  et  $\beta_{i+\frac{1}{2}}^* = \beta^-$  pour  $i < k$ . Pour  $i = k$ , le coefficient  $\beta_{k+\frac{1}{2}}^*$  est donné par :

$$\frac{1}{\beta_{k+\frac{1}{2}}^*} = \frac{1}{h_{k+\frac{1}{2}}} \left( \frac{h_{k+\frac{1}{2}}^+}{\beta^+} + \frac{h_{k+\frac{1}{2}}^-}{\beta^-} \right), \quad (1.28)$$

où  $h_{k+\frac{1}{2}} = h_{k+\frac{1}{2}}^+ + h_{k+\frac{1}{2}}^- = x_{k+1} - x_k$ . Le coefficient  $\beta_{k+\frac{1}{2}}^*$  est la moyenne harmonique de  $\beta^+$  et  $\beta^-$  pondérée par  $h_{k+\frac{1}{2}}^+$  et  $h_{k+\frac{1}{2}}^-$ .

**Etude numérique de la convergence.** On peut montrer (cf. [6, Theorem 2.3]) que le schéma (1.26)–(1.28) vérifie les estimations (1.15) et (1.16) c'est-à-dire que le schéma est au moins d'ordre  $\mathcal{O}(h)$  pour les normes  $L^\infty$  et  $H_0^1$  discrète. Dans ce paragraphe, on étudie numériquement la convergence du schéma Volumes Finis précédent et on montre qu'on obtient dans certaines situations (maillages uniformes et positions des centres) un ordre de convergence en  $h$  plus grand que 1. On compare également avec le même schéma obtenu en prenant le coefficient  $\beta_{k+\frac{1}{2}}^*$  comme la moyenne arithmétique (i.e.  $\beta_{k+\frac{1}{2}}^* = (\beta^- + \beta^+)/2$ ) au lieu de la moyenne harmonique (1.28).

Pour commencer, donnons une solution exacte de (1.19)–(1.21) dans le cas où  $f \equiv \gamma \in \mathbb{R}$ . La solution exacte  $u$  est alors donnée dans ce cas par

$$u(x) = \begin{cases} u_1(x) = \frac{\gamma}{2\beta^-} x(1 + \mu - x) & \text{pour } x \in [0, \alpha] \\ u_2(x) = \frac{\gamma}{2\beta^+} (1 - x)(x - \mu) & \text{pour } x \in [\alpha, 1] \end{cases} \quad (1.29)$$

avec  $\mu = \frac{\alpha(1 - \alpha)(\beta^- - \beta^+)}{\alpha\beta^+(1 - \alpha)\beta^-}$ . La Figure 1.2 montre la solution exacte avec  $f \equiv \gamma = 10$ ,  $\beta^- = 4$ ,  $\beta^+ = 1$  et  $\alpha = 0.4$ .

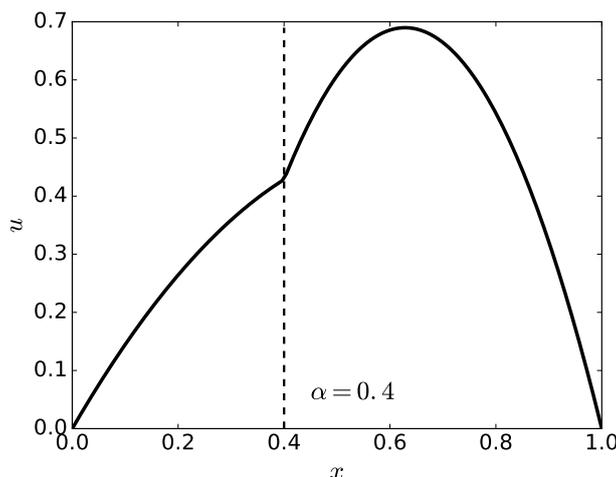


FIGURE 1.2 – Solution exacte de (1.17) avec  $\alpha = 0.4$ ,  $\beta^- = 4$ ,  $\beta^+ = 1$  et  $f \equiv \gamma = 10$ .

Avec un maillage non-uniforme (i.e. les  $h_i$  sont du même ordre mais différents) et des centres  $x_i$  qui ne sont pas aux milieux des cellules, on obtient les ordres de convergence indiqués dans le Tableau 1.1 en comparant les cas de la moyenne harmonique et arithmétique (voir aussi Figure 1.3).

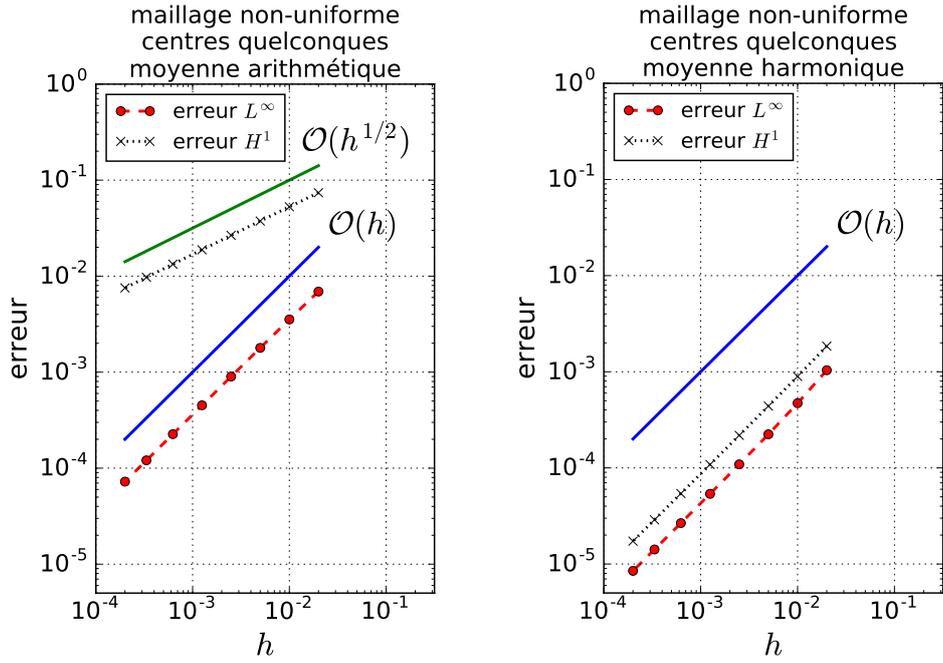


FIGURE 1.3 – Equation elliptique 1d avec coefficients discontinus : ordres de convergence pour un maillage non-uniforme avec des centres de cellules quelconques ; moyenne arithmétique (à gauche) et moyenne harmonique (à droite) des coefficients.

maillage non-uniforme	$\ \mathbf{e}\ _\infty$	$\ \mathbf{e}\ _{1,h}$
moyenne harmonique	$\mathcal{O}(h)$	$\mathcal{O}(h)$
moyenne arithmétique	$\mathcal{O}(h)$	$\mathcal{O}(h^{1/2})$

TABLE 1.1 – Ordres de convergence numériquement obtenus avec un maillage non-uniforme ; comparaison moyenne harmonique/moyenne arithmétique des coefficients

Avec la moyenne harmonique (1.28), lorsque les centres des cellules sont choisis aux milieux, on obtient numériquement un ordre de convergence en  $\mathcal{O}(h^2)$  pour la norme  $\|\cdot\|_\infty$ , que le maillage soit uniforme ou non (voir le Tableau 1.2 et la Figure 1.4). Par ailleurs, un maillage uniforme seul (i.e.  $h_i = h$  pour tout  $i$ ) ne suffit pas à obtenir un ordre en  $\mathcal{O}(h^2)$  pour la norme  $\|\cdot\|_\infty$ . En effet, un maillage uniforme avec des centres qui ne sont pas au milieu des cellules ne fournit pas un ordre en  $\mathcal{O}(h^2)$ .

moyenne harmonique centres milieux	$\ \mathbf{e}\ _\infty$	$\ \mathbf{e}\ _{1,h}$
maillage non-uniforme	$\mathcal{O}(h^2)$	$\mathcal{O}(h)$
maillage uniforme	$\mathcal{O}(h^2)$	$\mathcal{O}(h^{3/2})$

TABLE 1.2 – Ordres de convergence numériquement obtenus avec les centres aux milieux des cellules (moyenne harmonique) ; comparaison maillage uniforme/non-uniforme

L'ordre de convergence du schéma Volumes Finis avec la moyenne arithmétique des coefficients n'est pas amélioré avec un maillage uniforme ni en choisissant les centres aux milieux des cellules.

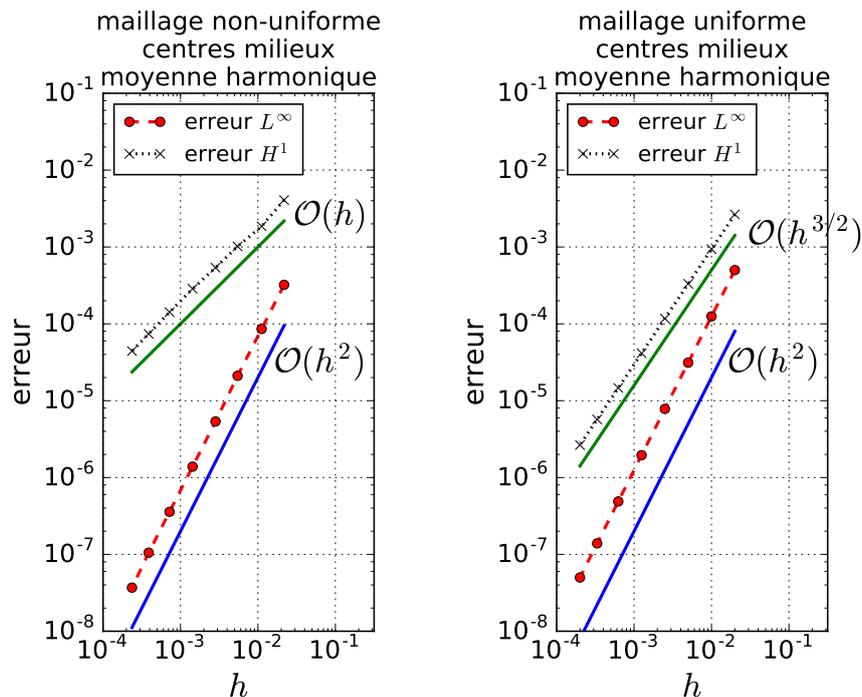


FIGURE 1.4 – Equation elliptique 1d avec coefficients discontinus : ordres de convergence avec les centres aux milieux des cellules et moyenne harmonique des coefficients ; maillage non-uniforme (à gauche) et uniforme (à droite).

## 1.4 Volumes Finis pour les problèmes elliptiques 2D

On va résoudre l'équation de Poisson par une méthode de Volumes Finis dans un domaine polygonal  $\Omega \subset \mathbb{R}^2$ . On cherche une fonction  $u = u(\mathbf{x})$  définie pour  $\mathbf{x} \in \Omega \subset \mathbb{R}^2$ , vérifiant

$$(P) \begin{cases} -\Delta u = f & \text{dans } \Omega \\ u = g & \text{sur le bord } \partial\Omega \end{cases}$$

avec des fonctions  $f, g \in L^2(\Omega)$  données. On suppose que la solution  $u$  de (P) possède la régularité  $u \in H^2(\Omega)$ . C'est vrai par exemple si  $\Omega$  est *convexe*,  $f \in L^2(\Omega)$  et  $g \in H^{3/2}(\partial\Omega)$  (voir [4]).

On va d'abord définir un maillage admissible au sens des Volumes Finis, puis on donnera la formulation en Volumes Finis du problème (P). On donnera ensuite des exemples concrets de maillages admissibles. Enfin, on terminera la section en traitant un problème elliptique à coefficients discontinus.

### 1.4.1 Maillage

On définit un maillage  $\mathcal{T}$  de  $\Omega$  par des volumes de contrôle (ou cellules)  $K$  de la façon suivante :

1. Les volumes de contrôle  $K$  sont des polygones convexes qui forment une partition de  $\Omega$  :

$$\bar{\Omega} = \cup_{K \in \mathcal{T}} \bar{K} \quad \text{et} \quad \overset{\circ}{K} \cap \overset{\circ}{L} = \emptyset, \quad \forall K, L \in \mathcal{T}, K \neq L.$$

2. Pour chaque cellule  $K$ , il existe un point  $\mathbf{x}_K \in K$  appelé *centre*, tel que les propriétés suivantes soient vérifiées :

- (a) Pour chaque cellule  $L$  adjacente à  $K$ , on a  $\mathbf{x}_K \neq \mathbf{x}_L$  et le segment de droite  $(\mathbf{x}_K, \mathbf{x}_L)$  est perpendiculaire à l'arête  $e$  commune aux deux cellules  $K$  et  $L$  (cf. Figure 1.5). On notera  $e = (K|L)$ .
- (b) Pour chaque arête  $e$  appartenant au bord  $\partial\Omega$ , la droite passant par  $\mathbf{x}_K$  et perpendiculaire à l'arête  $e$ , intersecte  $e$  (cf. Figure 1.6).

Un tel maillage sera dit *admissible* au sens des Volumes Finis.

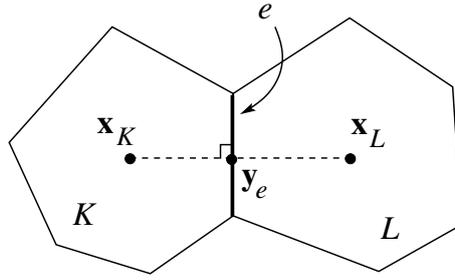


FIGURE 1.5 – Cellules admissibles d'un maillage "Volumes Finis"

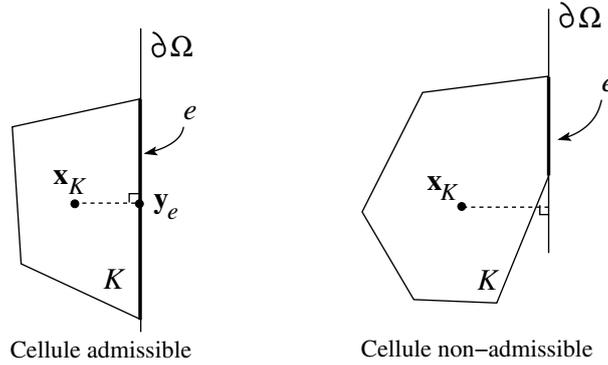


FIGURE 1.6 – Condition d'admissibilité du maillage pour les cellules du bord

Après avoir établi à la section suivante la formulation en Volumes Finis du problème ( $P$ ), on donnera des exemples de maillages *admissibles* au sens des Volumes Finis. Il s'agit de maillages triangulaires et de maillages de type Voronoï.

### 1.4.2 Formulation en Volumes Finis

On intègre l'équation de Poisson sur une cellule  $K$ .

$$\int_K -\Delta u \, d\mathbf{x} = \int_K f \, d\mathbf{x}$$

Par la formule de la divergence, on obtient

$$-\int_{\partial K} \nabla u \cdot \mathbf{n} \, d\Gamma = |K| f_K \quad (1.30)$$

où  $\mathbf{n}$  désigne la normale unitaire dirigée à l'extérieur de  $K$  et où on a noté  $f_K$  la valeur moyenne de  $f$  dans la cellule  $K$  i.e.

$$f_K = \frac{1}{|K|} \int_K f \, d\mathbf{x}. \quad (1.31)$$

On note  $\mathcal{E}_K$  l'ensemble des arêtes de la cellule  $K$  et on décompose le bord de la cellule  $K$  :

$$\partial K = \cup_{e \in \mathcal{E}_K} e.$$

La relation (1.30) s'écrit alors

$$\sum_{e \in \mathcal{E}_K} - \int_e \nabla u \cdot \mathbf{n}_{K,e} \, d\Gamma = |K| f_K \quad (1.32)$$

où on a noté  $\mathbf{n}_{K,e}$  la normale unitaire à  $e$  dirigée vers l'extérieur de  $K$ . Par ailleurs, pour toute arête intérieure  $e$  i.e. qui n'appartient pas au bord  $\partial\Omega$ , telle que  $e = (K|L)$ , la propriété suivante est satisfaite pour toute solution  $u \in H^2(\Omega)$  de ( $P$ ) :

$$\int_e \nabla u|_K \cdot \mathbf{n}_{K,e} \, d\Gamma = - \int_e \nabla u|_L \cdot \mathbf{n}_{L,e} \, d\Gamma \quad (1.33)$$

avec  $\mathbf{n}_{K,e} = -\mathbf{n}_{L,e}$ .

**Schéma VF.** On approche le flux à travers l'arête  $e$  :

$$-\int_e \nabla u \cdot \mathbf{n}_{K,e} d\Gamma \simeq F_{K,e} \quad (1.34)$$

où  $F_{K,e}$  est le *flux numérique* à travers l'arête  $e$ , associé à la cellule  $K$ . Le schéma "Volumes Finis" s'écrit

$$\sum_{e \in \mathcal{E}_K} F_{K,e} = |K| f_K, \quad \forall K \in \mathcal{T} \quad (1.35)$$

**Construction des flux numériques.** On considère les inconnues  $(u_K)_{K \in \mathcal{T}}$  associées à chaque volume de contrôle, avec les approximations  $u_K \simeq \frac{1}{|K|} \int_K u(\mathbf{x}) d\mathbf{x}$ . On désigne également par  $(u_e)_{e \in \mathcal{E}_K}$  des valeurs associées aux arêtes de la cellule  $K$ . Ces valeurs seront utilisées de façon intermédiaire et finalement éliminées.

On va distinguer les cas selon qu'une arête  $e$  appartient ou non au bord  $\partial\Omega$ .

- Soit une arête  $e$  d'une cellule  $K$  telle que  $e \not\subset \partial\Omega$  c'est-à-dire qui n'appartient pas au bord de  $\Omega$ .

(a) Pour un centre  $\mathbf{x}_K \notin e$ , le flux numérique  $F_{K,e}$  est choisi égal à :

$$F_{K,e} = -\frac{(u_e - u_K)}{d_{K,e}} |e| \quad (1.36)$$

où  $d_{K,e}$  est la distance de  $\mathbf{x}_K$  à l'arête  $e \subset \partial K$ . Le choix de  $F_{K,e}$  est donné pour  $\mathbf{x}_K \notin e$  de sorte que  $d_{K,e} \neq 0$ . L'approximation (1.34) avec (1.36) correspond à la **consistance des flux** numériques.

Pour éliminer les valeurs  $(u_e)$ , on suppose qu'il y a **conservation des flux** numériques à travers les arêtes. Précisément, pour l'arête  $e = (K|L)$  commune aux deux volumes de contrôle  $K$  et  $L$ , on impose :

$$F_{K,e} = -F_{L,e}. \quad (1.37)$$

C'est l'analogie discret de la propriété de conservation (1.33). On écrit alors :

$$F_{K,e} = -\frac{(u_e - u_K)}{d_{K,e}} |e| = \frac{(u_e - u_L)}{d_{L,e}} |e| = -F_{L,e}$$

ce qui donne

$$d_{L,e} u_K + d_{K,e} u_L = \underbrace{(d_{K,e} + d_{L,e})}_{d(\mathbf{x}_K, \mathbf{x}_L)} u_e.$$

Par conséquent,

$$u_e = \frac{d_{L,e} u_K + d_{K,e} u_L}{|\mathbf{x}_K - \mathbf{x}_L|}.$$

On obtient donc

$$F_{K,e} = -\frac{(u_L - u_K)}{|\mathbf{x}_K - \mathbf{x}_L|} |e|. \quad (1.38)$$

- (b) Pour un centre  $\mathbf{x}_K \in e$ , on choisit encore la relation (1.38) pour le flux numérique  $F_{K,e}$  (on a toujours  $\mathbf{x}_K \neq \mathbf{x}_L$ ).

- Soit une arête  $e$  d'une cellule  $K$  telle que  $e \subset \partial\Omega$  c'est-à-dire qui appartient au bord de  $\Omega$ .

(a) Si  $\mathbf{x}_K \notin e$ , alors on choisit

$$F_{K,e} = -\frac{(u_e - u_K)}{d_{K,e}} |e| \quad (1.39)$$

avec

$$u_e = \frac{1}{|e|} \int_e g d\Gamma. \quad (1.40)$$

(b) si  $\mathbf{x}_K \in e \subset \partial\Omega$ , alors on prend  $F_{K,e} = 0$  et on choisit

$$u_K = u_e = \frac{1}{|e|} \int_e g d\Gamma \quad (e \subset \partial\Omega). \quad (1.41)$$

Autrement dit, dans ce cas, la valeur  $u_K$  est connue dans la cellule  $K$ .

En résumé, le schéma "Volumes Finis" s'écrit

$$\sum_{e \in \mathcal{E}_K} F_{K,e} = |K|f_K, \quad \forall K \in \mathcal{T}$$

avec :

- si  $e \not\subset \partial\Omega$  avec  $e = (K|L)$  alors  $F_{K,e} = -\frac{(u_L - u_K)}{|\mathbf{x}_K - \mathbf{x}_L|}|e|$
- si  $e \subset \partial\Omega$ , on pose  $u_e = \frac{1}{|e|} \int_e g d\Gamma$ 
  - i) si  $\mathbf{x}_K \notin e$ , alors  $F_{K,e} = -\frac{(u_e - u_K)}{d_{K,e}}|e|$
  - ii) si  $\mathbf{x}_K \in e$ , alors  $F_{K,e} = 0$  et on prend  $u_K = u_e$

(1.42)

*Remarque 1* : Si on considère des conditions limites de Neumann homogènes  $\nabla u \cdot \mathbf{n} = 0$  sur une partie  $\Gamma \subset \partial\Omega$ , on impose alors  $F_{K,e} = 0$  pour toute arête  $e \subset \Gamma$ .

*Remarque 2* : Les intégrales dans (1.42) peuvent être calculées avec des formules de quadrature. Par exemple, on peut utiliser les formules suivantes :

$$\int_K f(\mathbf{x}) d\mathbf{x} = |K|f(\mathbf{x}_K) + \mathcal{O}(|K|^2),$$

$$\int_e g d\Gamma = |e|g(\mathbf{x}_K) + \mathcal{O}(|e|^2) \quad \text{si } \mathbf{x}_K \in e \subset \partial K \subset \partial\Omega.$$

### 1.4.3 Exemples de maillages admissibles

1. **Grille.** Pour  $\Omega = ]0, 1[ \times ]0, 1[$ , on peut utiliser un maillage cartésien (uniforme pour simplifier). Soit  $N$  le nombre de points de discrétisation dans les directions  $x$  et  $y$ . On note  $h = \frac{1}{N-1}$  le pas de discrétisation. Pour  $i, j = 1, \dots, N$ , les cellules  $K_{ij}$  sont définies par

$$K_{ij} = ]x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}[ \times ]y_{j-\frac{1}{2}}, y_{j+\frac{1}{2}}[ \quad \text{avec} \quad x_{i-\frac{1}{2}} = (i-1)h \quad \text{et} \quad y_{j-\frac{1}{2}} = (j-1)h.$$

Les centres  $\mathbf{x}_{K_{ij}}$  des cellules sont les barycentres des cellules :

$$\mathbf{x}_{K_{ij}} = (x_i, y_j)^\top \quad \text{avec} \quad x_i = x_{i-\frac{1}{2}} + \frac{h}{2} = i\frac{h}{2} \quad \text{et} \quad y_j = y_{j-\frac{1}{2}} + \frac{h}{2} = j\frac{h}{2}.$$

2. **Maillage triangulaire.** On considère une triangulation du domaine  $\Omega$ , admissible au sens des Eléments Finis : si deux triangles ont une arête commune alors ils ont deux sommets communs. On suppose que tous les angles des triangles sont plus petits que  $\pi/2$ . En conséquence, dans un triangle les médiatrices s'intersectent à l'intérieur du triangle. Les volumes de contrôle sont les triangles de la triangulation et les centres  $\mathbf{x}_K$  des cellules sont choisis comme les centres de masse des triangles i.e. l'intersection des médiatrices (cf. Figure 1.7). Avec l'hypothèse sur les angles, les centres  $\mathbf{x}_K$  se trouvent à l'intérieur des triangles.

On remarquera qu'une triangulation de Delaunay ne fournit pas nécessairement des triangles ayant tous des angles  $\leq \pi/2$ .

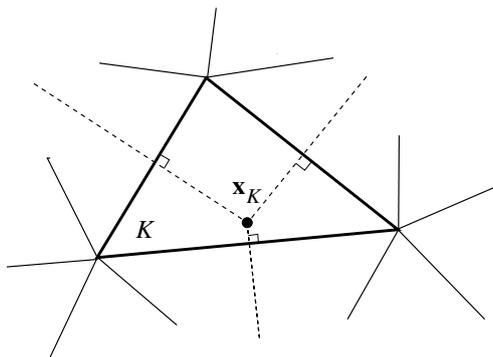


FIGURE 1.7 – Maillage triangulaire. Intersection des médiatrices à l'intérieur des triangles si les angles sont plus petits que  $\pi/2$ .

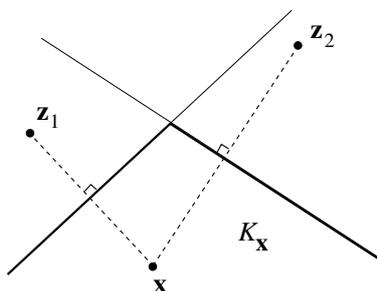


FIGURE 1.8 – Cellules de Voronoï

3. **Maillage Voronoï.** Soit  $\mathcal{P}$  un ensemble de points appartenant au domaine  $\Omega$ . On définit les cellules de Voronoï par rapport à chaque point  $\mathbf{x} \in \mathcal{P}$  par :

$$K_{\mathbf{x}} = \{\mathbf{y} \in \Omega, |\mathbf{x} - \mathbf{y}| < |\mathbf{z} - \mathbf{y}|, \forall \mathbf{z} \in \mathcal{P}, \mathbf{z} \neq \mathbf{x}\}.$$

La construction des cellules de Voronoï est basée sur la détermination des régions délimitées par les médiatrices des segments joignant les points de  $\mathcal{P}$  (cf. Figure 1.8).

Les volumes de contrôles  $K$  sont choisis comme étant les cellules de Voronoï associées à une triangulation de Delaunay du domaine  $\Omega$ . A chaque cellule de Voronoï  $K$ , on associe le centre  $\mathbf{x}_K \in K$  qui est un sommet de la triangulation (cf. Figure 1.9).

#### 1.4.4 Système linéaire

Pour fixer les idées, on choisit un maillage de **Voronoï**. En particulier, pour les cellules  $K$  ayant (au moins) une arête sur le bord de  $\Omega$ , le centre  $\mathbf{x}_K$  appartient au bord  $\partial\Omega$  (cf. Figure 1.10).

Avec le maillage de Voronoï, le schéma "Volumes Finis" s'écrit

$$\begin{aligned} & \bullet \sum_{\substack{e \in \mathcal{E}_K \\ e \not\subset \partial\Omega}} F_{K,e} = |K| f_K, \quad \forall K \in \mathcal{T} \\ & \text{avec } F_{K,e} = -\frac{(u_L - u_K)}{|\mathbf{x}_K - \mathbf{x}_L|} |e| \quad \text{si } e \not\subset \partial\Omega, e = (K|L) \\ & \bullet u_K = u_e = \frac{1}{|e|} \int_e g d\Gamma \quad \text{si } e \subset \partial K \subset \partial\Omega \end{aligned} \tag{1.43}$$

On numérote  $K_1, \dots, K_N$  les  $N$  cellules **intérieures** i.e. les cellules qui ne possèdent pas d'arêtes appartenant au bord  $\partial\Omega$ . On range les inconnues dans le vecteur  $\mathbf{u} = (u_{K_1}, \dots, u_{K_N})^\top \in \mathbb{R}^N$  et le système linéaire s'écrit alors

$$A\mathbf{u} = \mathbf{b} \tag{1.44}$$

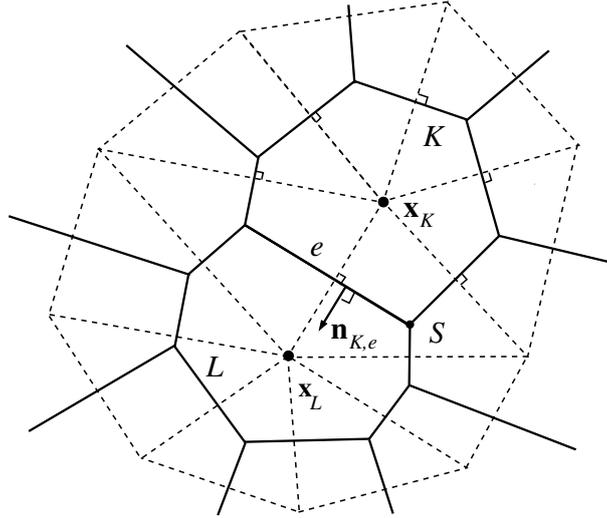


FIGURE 1.9 – Cellules de Voronoï associées à une triangulation de Delaunay

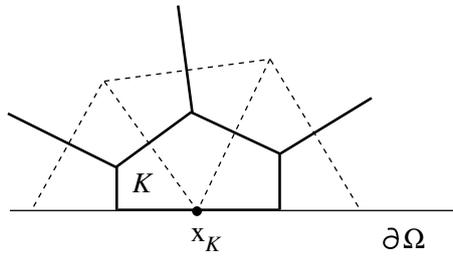


FIGURE 1.10 – Cellule d'un maillage de Voronoï sur le bord

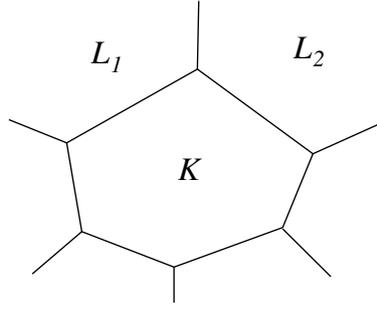
avec la matrice  $A$  de taille  $N \times N$ . On prend pour simplifier la condition limite  $g \equiv 0$  sur le bord  $\partial\Omega$ . Le second membre  $\mathbf{b} \in \mathbb{R}^N$  est alors donné par

$$\mathbf{b} = \begin{pmatrix} |K_1|f_{K_1} \\ \vdots \\ |K_N|f_{K_N} \end{pmatrix},$$

La matrice  $A$  est de la forme :

$$A = \begin{pmatrix} \dots & \sum_{\substack{e_{K,L} \subset \partial K \\ e_{K,L} \not\subset \partial\Omega}} \frac{|e_{K,L}|}{d_{K,L}} & \dots & -\frac{|e_{K,L_1}|}{d_{K,L_1}} & \dots & -\frac{|e_{K,L_2}|}{d_{K,L_2}} & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \dots & -\frac{|e_{L_1,K}|}{d_{L_1,K}} & \dots & \sum_{\substack{e_{L_1,L} \subset \partial L_1 \\ e_{L_1,L} \not\subset \partial\Omega}} \frac{|e_{L_1,L}|}{d_{L_1,L}} & \dots & -\frac{|e_{L_1,L_2}|}{d_{L_1,L_2}} & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \dots & -\frac{|e_{L_2,K}|}{d_{L_2,K}} & \dots & -\frac{|e_{L_2,L_1}|}{d_{L_2,L_1}} & \dots & \sum_{\substack{e_{L_2,L} \subset \partial L_2 \\ e_{L_2,L} \not\subset \partial\Omega}} \frac{|e_{L_2,L}|}{d_{L_2,L}} & \dots \end{pmatrix} \begin{matrix} K \\ \\ L_1 \\ L_2 \end{matrix} \quad (1.45)$$

On a noté  $e_{K,L} = (K|L) \not\subset \partial\Omega$  l'arête commune aux cellules  $K$  et  $L$  et  $d_{K,L} = |\mathbf{x}_K - \mathbf{x}_L|$ . Les cellules  $L_1, L_2, \dots$  sont les cellules **intérieures** ayant une arête en commun avec la cellule  $K$  (cf. Figure 1.11).

FIGURE 1.11 – Cellules adjacentes à la cellule  $K$ 

La matrice  $A$  est symétrique, elle est à diagonale fortement dominante et irréductible<sup>4</sup>. Elle est donc inversible (cf. Proposition B.3). En fait, on a  $A_{ij} \leq 0$  pour tous  $i \neq j$  et  $\sum_{1 \leq j \leq N} A_{ij} \geq 0$  pour  $1 \leq i \leq N$ . La matrice  $A$  est donc une  $M$ -matrice (cf. Proposition B.5).

**Assemblage de la matrice  $A$ .** Dans la pratique, pour construire la matrice  $A$ , on parcourt les arêtes du maillage de Voronoï. Pour une arête intérieure courante  $e = (K|L)$ , on ajoute les contributions des deux cellules adjacentes  $K$  et  $L$  :

$$\begin{array}{cc} & K & L \\ \left( \begin{array}{ccccc} \cdots & +\frac{|e|}{d_e} & \cdots & -\frac{|e|}{d_e} & \cdots \\ & \vdots & & \vdots & \\ \cdots & -\frac{|e|}{d_e} & \cdots & +\frac{|e|}{d_e} & \cdots \end{array} \right) & \begin{array}{l} K \\ L \end{array} \end{array}$$

avec  $d_e = |\mathbf{x}_K - \mathbf{x}_L|$ . Lorsqu'on parcourt l'arête intérieure  $e = (K|L)$ , on calcule les coefficients  $A(K, L)$ ,  $A(L, K)$  et on met à jour les coefficients  $A(K, K)$ ,  $A(L, L)$  :

$$\begin{aligned} A(K, K) &\leftarrow A(K, K) + |e|/d_e \\ A(L, L) &\leftarrow A(L, L) + |e|/d_e \\ A(K, L) &\leftarrow -|e|/d_e \\ A(L, K) &\leftarrow -|e|/d_e \end{aligned}$$

#### 1.4.5 Estimations d'erreurs

On établit des estimations d'erreurs en norme  $H^1$  et  $L^2$  dans le cas où  $u = 0$  sur le bord  $\partial\Omega$ . On introduit  $\mathcal{T}$  un maillage admissible de  $\Omega$  et on note  $\mathcal{E}$  l'ensemble des arêtes de  $\mathcal{T}$  avec la décomposition

$$\begin{aligned} \mathcal{E}_{\text{int}} &= \{e \in \mathcal{E}, e \not\subset \partial\Omega\} && \text{l'ensemble des arêtes intérieures,} \\ \mathcal{E}_{\text{ext}} &= \{e \in \mathcal{E}, e \subset \partial\Omega\} && \text{l'ensemble des arêtes du bord de } \Omega. \end{aligned}$$

On désignera par  $\mathcal{E}_K$  l'ensemble des arêtes de la cellule  $K \in \mathcal{T}$ . On note  $d_e$  la distance définie par

$$d_e = \begin{cases} |\mathbf{x}_K - \mathbf{x}_L| & \text{si } e = K|L \in \mathcal{E}_{\text{int}} \\ \text{dist}(\mathbf{x}_K, e) & \text{si } e \in \mathcal{E}_{\text{ext}} \cap \mathcal{E}_K, K \in \mathcal{T} \end{cases}$$

4. En effet, supposons que  $A$  soit réductible (cf. Annexe B). Alors il existe  $I, J$  *non-vides* formant une partition de  $\{1, \dots, N\}$  tels que  $A_{i,j} = 0$  pour tout  $(i, j) \in I \times J$ . Soit  $i_0 \in I$  et on considère  $I_0$  l'ensemble des indices des cellules *intérieures* qui sont adjacentes à la cellule  $K_{i_0}$ . On a  $A_{i_0, m} < 0$  pour tout  $m \in I_0$  et donc  $I_0 \cap J = \emptyset$ ,  $I_0 \subset I$ . Soit alors  $i_1 \in I_0$  ( $i_1 \neq i_0$ ) et on considère  $I_1$  l'ensemble des indices des cellules *intérieures* qui sont adjacentes à la cellule  $K_{i_1}$ . On a  $A_{i_1, m} < 0$  pour tout  $m \in I_1$  et donc  $I_1 \cap J = \emptyset$ ,  $I_1 \subset I$ . On continue avec un indice  $i_2 \in I_0 \cup I_1$  tel que  $i_2 \notin \{i_0, i_1\}$  et on considère  $I_2$  l'ensemble des indices des cellules *intérieures* qui sont adjacentes à la cellule  $K_{i_2}$ . On a  $A_{i_2, m} < 0$  pour tout  $m \in I_2$  et donc  $I_2 \cap J = \emptyset$ ,  $I_2 \subset I$ . On construit ainsi les ensembles  $I_k$  d'indices des cellules *intérieures* qui sont adjacentes à la cellule  $K_{i_k}$  et les indices  $i_k \in I_0 \cup I_1 \cdots \cup I_{k-1}$  avec  $i_k \notin \{i_0, i_1, \dots, i_{k-1}\}$  tels que  $(I_0 \cup I_1 \cdots \cup I_{k-1}) \cap J = \emptyset$ ,  $(I_0 \cup I_1 \cdots \cup I_{k-1}) \subset I$ . On obtient nécessairement  $(I_0 \cup I_1 \cdots \cup I_{N-1}) = \{1, \dots, N\} = I$  et donc  $J = \emptyset$  ce qui est impossible d'où la contradiction. La matrice  $A$  est donc irréductible.

et on suppose que

$$d_e \neq 0, \forall e \in \mathcal{E}. \quad (1.46)$$

Cette hypothèse implique en particulier que  $\mathbf{x}_K \notin \partial\Omega, \forall K \in \mathcal{T}$ .

Pour les estimations d'erreurs, on a besoin de définir la norme  $H_0^1$  discrète. Soit  $X(\mathcal{T})$  l'ensemble des fonctions de  $\Omega$  dans  $\mathbb{R}$  qui sont constantes sur chaque cellule de  $\mathcal{T}$ . Pour  $u \in X(\mathcal{T})$ , on définit la norme

$$\|u\|_{1,\mathcal{T}} = \left( \sum_{e \in \mathcal{E}} \frac{|e|}{d_e} (D_e u)^2 \right)^{1/2} \quad (1.47)$$

où

$$D_e u = \begin{cases} u_K - u_L & \text{si } e = K|L \in \mathcal{E}_{\text{int}}, \\ u_K & \text{si } e \in \mathcal{E}_K \cap \mathcal{E}_{\text{ext}}. \end{cases} \quad (1.48)$$

On a une inégalité de Poincaré pour la norme  $H_0^1$  discrète (voir par exemple [6]).

**Lemme 1.1 (Inégalité de Poincaré discrète)** *Pour tout  $u \in X(\mathcal{T})$ , on a*

$$\|u\|_{L^2(\Omega)} \leq \text{diam}(\Omega) \|u\|_{1,\mathcal{T}}$$

où  $\text{diam}(\Omega)$  désigne le diamètre du plus petit cercle contenant le domaine  $\Omega$ .

On note  $h$  la taille du maillage défini par  $h = \max(\text{diam}(K), K \in \mathcal{T})$  où  $\text{diam}(K)$  désigne le diamètre du plus petit cercle contenant la cellule  $K$ . Les estimations d'erreurs de la méthode des volumes finis (1.43) pour le problème (P) sont les suivantes (cf. démonstration en exercice).

**Théorème 1.2** *Soit  $\mathcal{T}$  un maillage admissible de  $\Omega$  vérifiant (1.46). On suppose que la solution  $u$  de (P) avec  $u = 0$  sur  $\partial\Omega$  vérifie  $u \in C^2(\overline{\Omega})$ . Soit  $(u_K)_{K \in \mathcal{T}}$  la solution du système (1.43) et on définit la fonction  $\delta_{\mathcal{T}} \in X(\mathcal{T})$  avec  $\delta_{\mathcal{T}}(\mathbf{x}) = u(\mathbf{x}_K) - u_K$  pour presque tout  $\mathbf{x} \in K$ . Alors, il existe une constante  $C > 0$  indépendante de  $h = \max(\text{diam}(K), K \in \mathcal{T})$  telle que*

$$\|\delta_{\mathcal{T}}\|_{L^2(\Omega)} + \|\delta_{\mathcal{T}}\|_{1,\mathcal{T}} \leq Ch.$$

### 1.4.6 Equation elliptique 2D avec coefficients discontinus

On généralise la section précédente en considérant un problème elliptique où les coefficients de l'EDP sont discontinus. On considère un domaine  $\Omega \subset \mathbb{R}^2$  polygonal convexe se décomposant en deux sous-domaines  $\Omega_1$  et  $\Omega_2$  *polygonaux* formant une partition de  $\Omega$  avec  $\overline{\Omega_2} \subset \Omega$ ,  $\Omega_1 = \Omega \setminus \overline{\Omega_2}$  (voir Figure 1.12). On cherche alors une fonction  $u = u(\mathbf{x})$  définie pour  $\mathbf{x} \in \Omega$  et vérifiant

$$\begin{cases} -\text{div}(a(\mathbf{x})\nabla u) = f & \text{dans } \Omega \\ u = 0 & \text{sur le bord } \partial\Omega \end{cases} \quad (1.49)$$

avec  $f \in L^2(\Omega)$  donnée et  $a = a(\mathbf{x})$  est une fonction constante par morceaux sur  $\Omega$  :

$$a(\mathbf{x}) = \begin{cases} a_1 & \text{si } \mathbf{x} \in \Omega_1 \\ a_2 & \text{si } \mathbf{x} \in \Omega_2, \end{cases} \quad (1.50)$$

où  $a_1$  et  $a_2$  sont des *constantes* réelles.

La solution  $u$  de (1.49) est régulière sur  $\Omega_1$  et  $\Omega_2$  (on a  $u \in H^2(\Omega_1)$ ,  $u \in H^2(\Omega_2)$  mais attention en général  $u \notin H^2(\Omega)$ ) et elle vérifie

$$-\text{div}(a_i \nabla u) = f \quad \text{dans } \Omega_i, \quad i = 1, 2 \quad (1.51)$$

$$u = 0 \quad \text{sur } \partial\Omega \quad (1.52)$$

On note  $\gamma = \partial\Omega_1 \cap \partial\Omega_2$  la frontière commune aux deux sous-domaines (cf. Figure 1.12).

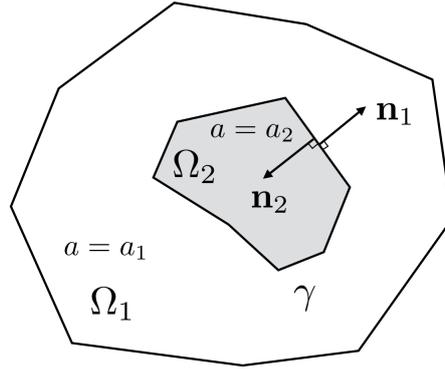


FIGURE 1.12 – Sous-domaines pour les coefficients discontinus

Pour tout  $v \in H_0^1(\Omega)$ , la solution  $u$  vérifie également la relation <sup>5</sup> :

$$\int_{\gamma} a_1 \nabla u \cdot \mathbf{n}_1 v \, d\Gamma = - \int_{\gamma} a_2 \nabla u \cdot \mathbf{n}_2 v \, d\Gamma \quad (1.54)$$

avec  $\mathbf{n}_1 = -\mathbf{n}_2$  est la normale unitaire à  $\gamma$  dirigée vers l'extérieur de  $\Omega_1$  (cf. Figure 1.12).

On considère un maillage *admissible*  $\mathcal{T} = \mathcal{T}_1 \cup \mathcal{T}_2$  de  $\Omega$  où  $\mathcal{T}_1$  et  $\mathcal{T}_2$  sont des maillages admissibles de  $\Omega_1$  et  $\Omega_2$  respectivement (cf. Section 1.4.1). En intégrant (1.51) sur une cellule  $K$  quelconque du maillage  $\mathcal{T}$ , on obtient

$$- \int_{\partial K} a(\mathbf{x}) \nabla u \cdot \mathbf{n}_K \, d\Gamma = \int_K f \, d\mathbf{x} \quad (1.55)$$

où  $\mathbf{n}_K$  est la normale unitaire extérieure à  $K$ . On déduit de (1.54) et (1.55) que, pour toutes cellules  $K$  et  $L$  adjacentes, on a la relation

$$\int_e (a(\mathbf{x}) \nabla u)|_K \cdot \mathbf{n}_{K,e} \, d\Gamma = - \int_e (a(\mathbf{x}) \nabla u)|_L \cdot \mathbf{n}_{L,e} \, d\Gamma. \quad (1.56)$$

avec  $\mathbf{n}_{K,e}$  la normale unitaire à  $e = (K|L)$ , dirigée vers l'extérieur de  $K$  et  $\mathbf{n}_{K,e} = -\mathbf{n}_{L,e}$ .

La relation (1.55) s'écrit encore

$$\sum_{e \in \mathcal{E}_K} \mathcal{F}_{K,e} = |K| f_K \quad (1.57)$$

avec le flux  $\mathcal{F}_{K,e}$  associé à la cellule  $K$  à travers l'arête  $e \subset \partial K$ , donné par

$$\mathcal{F}_{K,e} = - \int_e a(\mathbf{x}) \nabla u \cdot \mathbf{n}_{K,e} \, d\Gamma, \quad (1.58)$$

où  $\mathbf{n}_{K,e}$  est la normale unitaire à  $e$ , dirigée vers l'extérieur de  $K$ .

<sup>5</sup>. Soit  $v \in H_0^1(\Omega)$ . La solution faible  $u \in H_0^1(\Omega)$  de (1.49) vérifie

$$\int_{\Omega} a(\mathbf{x}) \nabla u \cdot \nabla v \, d\mathbf{x} = \int_{\Omega} f v \, d\mathbf{x}. \quad (1.53)$$

Par ailleurs, en multipliant (1.51) par  $v$  et en intégrant sur chaque sous-domaine  $\Omega_i$ , on obtient pour  $i = 1, 2$  :  $\int_{\Omega_i} a_i \nabla u \cdot \nabla v \, d\mathbf{x} - \int_{\gamma} a_i \nabla u \cdot \mathbf{n}_i v \, d\Gamma = \int_{\Omega_i} f v \, d\mathbf{x}$ . En sommant sur  $i$ , il vient

$$\int_{\Omega} a(\mathbf{x}) \nabla u \cdot \nabla v \, d\mathbf{x} - \int_{\gamma} a_1 \nabla u \cdot \mathbf{n}_1 v \, d\Gamma - \int_{\gamma} a_2 \nabla u \cdot \mathbf{n}_2 v \, d\Gamma = \int_{\Omega} f v \, d\mathbf{x}$$

En comparant avec (1.53), on obtient (1.54)

**Approximations des flux.** Si  $e \not\subset \partial\Omega$  est une arête commune aux deux cellules  $K$  et  $L$ , on approche le flux exacte  $\mathcal{F}_{K,e}$  par

$$\mathcal{F}_{K,e} \simeq F_{K,e} \stackrel{\text{def}}{=} -a_{K,e} \frac{(u_e - u_K)}{d_{K,e}} |e| \quad (1.59)$$

où  $d_{K,e} = \text{dist}(\mathbf{x}_K, e)$  est la distance de  $\mathbf{x}_K$  à l'arête  $e \subset \partial K$  et

$$a_{K,e} = \lim_{\substack{\mathbf{x} \rightarrow \mathbf{x}_e \\ \mathbf{x} \in K}} a(\mathbf{x}), \quad (1.60)$$

où  $\mathbf{x}_e$  désigne le point d'intersection de l'arête  $e$  avec le segment  $[\mathbf{x}_K, \mathbf{x}_L]$ . Le terme  $u_e$  représente une approximation de  $u$  sur l'arête  $e$ . On impose la conservation des flux numériques au travers de l'arête  $e = (K|L)$ . C'est l'analogie discret de la conservation (1.56) :

$$F_{K,e} = -F_{L,e}. \quad (1.61)$$

On peut alors éliminer le terme  $u_e$  dans l'expression (1.59) de  $F_{K,e}$  et on obtient ( $e \not\subset \partial\Omega$ )

$$F_{K,e} = -\frac{a_{K,e}a_{L,e}}{(d_{K,e}a_{L,e} + d_{L,e}a_{K,e})} (u_L - u_K) |e|.$$

La relation précédente s'écrit encore

$$F_{K,e} = -a_{K,L}^* \frac{(u_L - u_K)}{d_{K,L}} |e|, \quad (1.62)$$

où le coefficient  $a_{K,L}^*$  est donné par

$$\frac{1}{a_{K,L}^*} = \frac{1}{a_{K,e}} \left( \frac{d_{K,e}}{d_{K,L}} \right) + \frac{1}{a_{L,e}} \left( \frac{d_{L,e}}{d_{K,L}} \right)$$

avec  $d_{K,L} = d_{K,e} + d_{L,e} = |\mathbf{x}_K - \mathbf{x}_L|$ . Le coefficient  $a_{K,L}^*$  est obtenu par *moyenne harmonique* de  $a_{K,e}$  et  $a_{L,e}$  pondérée par les distances  $d_{K,e}$  et  $d_{L,e}$ .

**Remarque.** Si  $a$  est continue, on retrouve le coefficient  $a_{K,L}^* = a_{K,e} = a_{L,e}$ .



# Chapitre 2

## Equations paraboliques

### 2.1 Introduction

Pour un domaine (ouvert connexe)  $\Omega \subset \mathbb{R}^n$  borné, de frontière  $\partial\Omega$  régulière et pour un temps  $T > 0$  fixé, on considère le problème aux limites suivant : trouver une fonction  $u = u(\mathbf{x}, t)$  avec  $\mathbf{x} \in \Omega$  et  $t \in (0, T)$ , telle que

$$\frac{\partial u}{\partial t} - Lu = f \quad \text{dans } \Omega \times (0, T) \quad (2.1)$$

$$u = 0 \quad \text{sur } \partial\Omega \times (0, T) \quad (2.2)$$

$$u(\mathbf{x}, 0) = u_0(\mathbf{x}) \quad \text{dans } \Omega. \quad (2.3)$$

La condition (2.2) est la *condition limite* sur le bord du domaine  $\Omega$  (ici Dirichlet homogène). La condition (2.3) est la *condition initiale* à l'instant  $t = 0$ . L'opérateur  $L$  est défini par

$$Lu := \operatorname{div}(A\nabla u) - \mathbf{b} \cdot \nabla u - cu \quad (2.4)$$

$$= \sum_{i,j=1}^n \frac{\partial}{\partial x_i} \left( a_{ij} \frac{\partial u}{\partial x_j} \right) - \sum_{i=1}^n b_i \frac{\partial u}{\partial x_i} - cu \quad (2.5)$$

avec  $A = A(\mathbf{x}, t) = (a_{ij}(\mathbf{x}, t)) \in \mathbb{R}^{n,n}$ ,  $\mathbf{b} = \mathbf{b}(\mathbf{x}, t) = (b_i(\mathbf{x}, t)) \in \mathbb{R}^n$  et  $c = c(\mathbf{x}, t)$ . On suppose dans tout ce chapitre (sauf précision contraire) que les coefficients  $a_{ij} \in C^1(\bar{\Omega} \times [0, T])$ . Les fonctions  $\mathbf{b}$ ,  $c$  et  $f = f(\mathbf{x}, t)$  sont données dans  $C(\bar{\Omega} \times [0, T])$ .

On dit que l'équation (2.1) est uniformément **parabolique** si l'opérateur  $-L$  est uniformément elliptique en la variable d'espace  $\mathbf{x}$ , c'est-à-dire si la condition suivante est satisfaite :

$$\exists \alpha > 0 \text{ tel que } \langle A(\mathbf{x}, t)\xi, \xi \rangle_{\mathbb{R}^n} = \sum_{i,j=1}^n a_{ij}(\mathbf{x}, t)\xi_i\xi_j \geq \alpha\|\xi\|^2, \quad \forall \xi \in \mathbb{R}^n, \forall (\mathbf{x}, t) \in \Omega \times (0, T). \quad (2.6)$$

On supposera désormais que la condition (2.6) est vérifiée.

Voici à présent quelques résultats d'existence, d'unicité et de principe du maximum relatifs aux équations paraboliques (on renvoie par exemple aux ouvrages [2], [5], [15], [12]).

#### 2.1.1 Existence et unicité des solutions

Commençons par le cas de l'équation de la chaleur avec  $L = \Delta$ .

- Si  $f \in L^2(\Omega \times (0, T))$  et  $u_0 \in H_0^1(\Omega)$  alors le problème (2.1)-(2.3) avec  $L = \Delta$  admet une unique solution

$$u \in L^2(0, T; H^2(\Omega) \cap H_0^1(\Omega)) \cap C([0, T]; H_0^1(\Omega)) \quad (2.7)$$
$$\frac{\partial u}{\partial t} \in L^2(0, T; L^2(\Omega))$$

vérifiant les équations (2.1), (2.2) et (2.3) au sens *presque partout*.

- Si  $f \in C^\infty(\overline{\Omega} \times [0, T])$  et  $u_0 \in L^2(\Omega)$  alors il existe une unique solution

$$u \in C^\infty(\overline{\Omega} \times [\varepsilon, T]) \quad \forall \varepsilon > 0. \quad (2.8)$$

Si de plus,  $u_0 \in C^\infty(\overline{\Omega})$  et  $(f, u_0)$  vérifie certaines relations de compatibilité<sup>(1)</sup> sur  $\partial\Omega$  alors  $u \in C^\infty(\overline{\Omega} \times [0, T])$ .

D'après (2.8), on voit que l'équation de la chaleur a un effet fortement régularisant sur la donnée initiale  $u_0$ . En effet, la solution  $u$  est  $C^\infty$  en  $x$  pour chaque  $t > 0$ , même si la donnée initiale  $u_0$  est discontinue.

Ces résultats se généralisent au cas d'un opérateur  $L$  à coefficients *réguliers* :

- Si  $a_{ij}, b_i, c \in C^\infty(\overline{\Omega} \times [0, T])$  et si les coefficients  $a_{ij}$  vérifient la condition d'ellipticité (2.6) alors tous les résultats précédents restent vrais. On remarquera qu'il n'y a pas d'hypothèse sur le signe de la fonction  $c$ , contrairement au cas elliptique (cf. Chap. 2).

### 2.1.2 Principes du maximum

Soit  $u \in C^2(\Omega \times (0, T]) \cap C^0(\overline{\Omega} \times [0, T])$  vérifiant  $u_t - Lu = f$ .

- On prend  $c \equiv 0$ . Si  $f \leq 0$  (resp.  $f \geq 0$ ) alors  $u$  atteint son maximum (resp. minimum) sur  $\Sigma = \partial\Omega \times [0, T] \cup \Omega \times \{t = 0\}$ . En particulier, si  $u$  est solution de (2.1)-(2.3) avec  $c \equiv 0$ ,  $f \geq 0$  et  $u_0 \geq 0$ , alors  $u \geq 0$  dans  $\overline{\Omega} \times [0, T]$ .
- (PRINCIPE DE HOPF) Soit  $f \leq 0$  et soient  $\mathbf{x}_0 \in \partial\Omega$ ,  $t_0 \in (0, T)$  tels que  $u(\mathbf{x}, t_0) < u(\mathbf{x}_0, t_0)$  pour tout  $\mathbf{x} \in \Omega$ . On suppose que  $u$  est dérivable en  $\mathbf{x}_0$ . Si  $c \equiv 0$  ou bien si  $u(\mathbf{x}_0, t_0) = 0$  alors

$$\frac{\partial u}{\partial \mathbf{n}}(\mathbf{x}_0, t_0) > 0,$$

où  $\mathbf{n}$  désigne la normale extérieure à  $\partial\Omega$ .

Pour l'équation de la chaleur avec  $f \equiv 0$  et  $c \equiv 0$ , si  $u_0 \geq 0$  avec  $u_0 \not\equiv 0$  alors la solution  $u$  de (2.1)-(2.3) vérifie  $u(\mathbf{x}, t) > 0$ ,  $\forall \mathbf{x} \in \Omega$ ,  $\forall t > 0$ . Autrement dit, l'effet d'une petite perturbation initiale est ressenti immédiatement partout : la chaleur se propage avec une vitesse infinie.

## 2.2 Volumes Finis pour l'équation de la chaleur en dimension 1 d'espace

On considère l'équation en 1D d'espace avec  $\Omega = (0, 1)$  : trouver  $u = u(x, t)$  pour  $(x, t) \in (0, 1) \times (0, T)$  telle que

$$(P) \begin{cases} \frac{\partial u}{\partial t} - \gamma \frac{\partial^2 u}{\partial x^2} = f & \text{dans } Q_T := (0, 1) \times (0, T) \\ u(0, t) = u(1, t) = 0, & t \in (0, T) \\ u(x, 0) = u_0(x), & x \in (0, 1) \end{cases}$$

avec  $\gamma > 0$  une constante donnée. On suppose que le problème (P) admet une unique solution  $u \in C^2(Q_T) \cap C^0(\overline{Q}_T)$ . On notera qu'on a l'estimation suivante (avec  $u_0 \in H_0^1(0, 1)$  et  $f \in L^2(Q_T)$ ) :

$$\|u(t)\|_{H^1(0,1)} \leq C \left( \|u_0\|_{H^1(0,1)} + \|f\|_{L^2(Q_T)} \right) \quad \text{pour tout } t \in (0, T). \quad (2.9)$$

---

1. Les relations de compatibilité sont des conditions nécessaires pour que  $u \in C^\infty(\overline{\Omega} \times [0, T])$ . Elles s'obtiennent en écrivant que toutes les dérivées de  $u$  par rapport à  $t$ , sont nulles sur  $\partial\Omega \times [0, T]$  i.e.  $u = \partial_t u = \dots = \partial_t^j u = \dots = 0$  sur  $\partial\Omega \times [0, T]$ . On dérive de façon itérée l'équation (2.1) par rapport à  $t$  et on utilise les relations précédentes. Par exemple, l'équation  $u_t = \Delta u + f$  dans  $\overline{\Omega} \times [0, T]$ , fournit sur  $\partial\Omega \times \{t = 0\}$  la relation  $-\Delta u_0 = f(\mathbf{x}, 0)$  pour  $\mathbf{x} \in \partial\Omega$ . Puis l'équation  $u_{tt} = \Delta u_t + f_t = \Delta^2 u + \Delta f + f_t$  dans  $\overline{\Omega} \times [0, T]$ , fournit sur  $\partial\Omega \times \{t = 0\}$  la relation  $-\Delta^2 u_0 = \Delta f(\mathbf{x}, 0) + f_t(\mathbf{x}, 0)$  pour  $\mathbf{x} \in \partial\Omega$ , et ainsi de suite ... On peut prendre par exemple, les relations suivantes

$$\begin{aligned} u_0 &= \Delta u_0 = \dots = \Delta^j u_0 = \dots = 0 \text{ sur } \partial\Omega, \\ f &= \Delta f = \dots = \Delta^j f = \dots = 0 \text{ sur } \partial\Omega \times (0, T). \end{aligned}$$

On remarquera enfin, que l'hypothèse  $u_0 \in C^\infty(\overline{\Omega})$  avec  $u_0 = 0$  sur  $\partial\Omega$  ne suffit pas pour avoir  $u \in C^\infty(\overline{\Omega} \times [0, T])$ ...

On discrétise l'intervalle  $[0, 1]$  par un maillage  $\mathcal{T}$  défini par les cellules de contrôle  $K_i = ]x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}[$  de centres  $x_i$ ,  $i = 1, \dots, N$  (cf. Section 1.2.1 du Chapitre 1). On note  $h_i = |K_i|$  et  $u_i(t) = \frac{1}{h_i} \int_{K_i} u(x, t) dx$  pour  $t \in [0, T]$ . En intégrant l'équation différentielle de  $(P)$  sur la cellule  $K_i$ , on obtient

$$\frac{d}{dt} \int_{K_i} u(x, t) dx - \gamma \left( \frac{\partial u}{\partial x}(x_{i+\frac{1}{2}}, t) - \frac{\partial u}{\partial x}(x_{i-\frac{1}{2}}, t) \right) = h_i f_i(t),$$

où  $f_i(t) = \frac{1}{h_i} \int_{K_i} f(x, t) dx$ . On obtient ainsi, pour tout  $i = 1, \dots, N$  et  $t \in [0, T]$ ,

$$\frac{d}{dt} u_i(t) + \frac{\mathcal{F}_{i+\frac{1}{2}}(t) - \mathcal{F}_{i-\frac{1}{2}}(t)}{h_i} = f_i(t), \quad (2.10)$$

où on a noté  $\mathcal{F}_{i+\frac{1}{2}}(t) = -\gamma \frac{\partial u}{\partial x}(x_{i+\frac{1}{2}}, t)$  le flux exact de  $u$  en  $x_{i+\frac{1}{2}}$  à l'instant  $t$ .

### 2.2.1 Schéma d'Euler explicite

On discrétise en espace et en temps l'équation (2.10). Soit  $\Delta t > 0$  le pas de discrétisation en temps et on introduit les instants  $t^n = n\Delta t$  pour  $n = 0, \dots, M$  avec  $T = M\Delta t$ . On considère l'approximation  $u_i^n \simeq u_i(t^n) = \frac{1}{h_i} \int_{K_i} u(x, t) dx$  et le flux numérique  $F_{i+\frac{1}{2}}^n \simeq -\gamma \frac{\partial u}{\partial x}(x_{i+\frac{1}{2}}, t^n)$  donné par (cf. Chapitre 1, Section 1.3.2)

$$F_{i+\frac{1}{2}}^n = -\gamma \frac{u_{i+1}^n - u_i^n}{h_{i+\frac{1}{2}}}. \quad (2.11)$$

Le schéma d'Euler explicite s'écrit en approchant la dérivée en temps

$$\frac{\partial u_i}{\partial t}(t^n) \simeq \frac{u_i^{n+1} - u_i^n}{\Delta t}$$

et en écrivant l'équation (2.10) au temps  $t^n$ . On obtient, pour  $i = 1, \dots, N$ ,  $n = 0, \dots, M$ ,

$$\frac{u_i^{n+1} - u_i^n}{\Delta t} + \frac{F_{i+\frac{1}{2}}^n - F_{i-\frac{1}{2}}^n}{h_i} = f_i^n \quad (2.12)$$

où  $f_i^n = f_i(t^n) = \frac{1}{h_i} \int_{K_i} f(x, t^n) dx$ .

#### Forme matricielle.

En combinant (2.11),(2.12), on a, pour  $i = 1, \dots, N$  et  $n = 0, \dots, M$ ,

$$\frac{u_i^{n+1} - u_i^n}{\Delta t} - \gamma \frac{u_{i+1}^n - u_i^n}{h_i h_{i+\frac{1}{2}}} + \gamma \frac{u_i^n - u_{i-1}^n}{h_{i-\frac{1}{2}} h_i} = f_i^n,$$

qu'on peut écrire sous la forme

$$u_i^{n+1} = \gamma \frac{\Delta t}{h_i h_{i-\frac{1}{2}}} u_{i-1}^n + \left( 1 - \gamma \frac{\Delta t}{h_i} \left( \frac{1}{h_{i-\frac{1}{2}}} + \frac{1}{h_{i+\frac{1}{2}}} \right) \right) u_i^n + \gamma \frac{\Delta t}{h_i h_{i+\frac{1}{2}}} u_{i+1}^n + \Delta t f_i^n.$$

Le schéma "Volumes Finis" Euler explicite s'écrit donc

$$u_i^{n+1} = -\gamma \Delta t \frac{\beta_{i-1}}{h_i} u_{i-1}^n + \left( 1 - \gamma \Delta t \frac{\alpha_i}{h_i} \right) u_i^n - \gamma \Delta t \frac{\beta_i}{h_i} u_{i+1}^n + \Delta t f_i^n \quad (2.13)$$

avec

$$\alpha_i = \frac{1}{h_{i-\frac{1}{2}}} + \frac{1}{h_{i+\frac{1}{2}}} > 0, \quad \beta_i = -\frac{1}{h_{i+\frac{1}{2}}} < 0, \quad (2.14)$$

pour  $i = 1, \dots, N$  et  $n = 0, \dots, M$ . On pose également  $\beta_0 = \beta_N = 0$ . Les conditions limites et initiales sont données par

$$\begin{aligned} u_0^n &= u_{N+1}^n = 0, \quad n = 0, \dots, M \\ u_i^0 &= \frac{1}{h_i} \int_{K_i} u_0(x) dx, \quad i = 1, \dots, N \end{aligned} \quad (2.15)$$

On pose  $\mathbf{u}^n = (u_1^n, \dots, u_N^n)^\top \in \mathbb{R}^N$  et  $\mathbf{f}^n = (f_1^n, \dots, f_N^n)^\top \in \mathbb{R}^N$ . Le schéma d'Euler explicite s'écrit vectoriellement

$$\mathbf{u}^{n+1} = (I_N - \gamma \Delta t H^{-1} A) \mathbf{u}^n + \Delta t \mathbf{f}^n, \quad (2.16)$$

où  $A$  est la matrice du schéma Volumes Finis pour l'équation elliptique 1d du Chapitre 1 (cf. (1.13)) :

$$A = \begin{pmatrix} \alpha_1 & \beta_1 & & & 0 \\ \beta_1 & \alpha_2 & \beta_2 & & \\ & \ddots & \ddots & \ddots & \\ & & \beta_{N-2} & \alpha_{N-1} & \beta_{N-1} \\ 0 & & & \beta_{N-1} & \alpha_N \end{pmatrix} \quad (2.17)$$

avec  $\alpha_i, \beta_i$  définis par (2.14). La matrice  $H$  est la matrice diagonale telle que  $H_{ii} = h_i$ ,  $1 \leq i \leq N$  i.e.

$$H = \begin{pmatrix} h_1 & & & 0 \\ & h_2 & & \\ & & \ddots & \\ 0 & & & h_N \end{pmatrix}, \quad H^{-1} = \begin{pmatrix} 1/h_1 & & & 0 \\ & 1/h_2 & & \\ & & \ddots & \\ 0 & & & 1/h_N \end{pmatrix}. \quad (2.18)$$

Le système (2.16) est *explicite* dans la mesure où il n'y a pas de système linéaire à résoudre : le vecteur  $\mathbf{u}^{n+1}$  est calculé directement à partir de  $\mathbf{u}^n$  sans résolution de système linéaire.

### Stabilité.

On note  $\|\mathbf{v}\|_\infty = \max_{1 \leq i \leq N} |v_i|$  pour  $\mathbf{v} = (v_1, \dots, v_N)^\top$ .

**Proposition 2.1** Soit  $\lambda = \gamma \Delta t \max_{1 \leq i \leq N} \frac{1}{h_i} \left( \frac{1}{h_{i-\frac{1}{2}}} + \frac{1}{h_{i+\frac{1}{2}}} \right)$ .

1. Principe du maximum discret.

Soit  $f \geq 0$  et  $u_0 \geq 0$ . Si  $\lambda \leq 1$  alors  $\mathbf{u}^n \geq \mathbf{0}$  pour tout  $1 \leq n \leq N$ .

2. Stabilité  $L^\infty$ .

Si  $\lambda \leq 1$  alors

$$\|\mathbf{u}^n\|_\infty \leq \|\mathbf{u}^0\|_\infty + T \|f\|_{L^\infty(\Omega \times (0, T))}, \quad (2.19)$$

pour tout  $n = 1, \dots, M$ .

*Démonstration.* Sous l'hypothèse  $\lambda \leq 1$ , on a  $1 - \gamma \Delta t \frac{\alpha_i}{h_i} = 1 - \gamma \frac{\Delta t}{h_i} \left( \frac{1}{h_{i-\frac{1}{2}}} + \frac{1}{h_{i+\frac{1}{2}}} \right) \geq 0$  pour tout  $i = 1, \dots, N$ .

1. Par récurrence sur  $n$ , on montre alors facilement à partir de (2.13) que  $\mathbf{u}^n \geq \mathbf{0} \Rightarrow \mathbf{u}^{n+1} \geq \mathbf{0}$ .

2. D'après (2.13), pour tout  $i = 1, \dots, N$ , on a ( $\beta_i < 0$ )

$$|u_i^{n+1}| \leq \left(1 - \gamma \Delta t \frac{\alpha_i}{h_i}\right) |u_i^n| - \gamma \Delta t \frac{\beta_{i-1}}{h_i} |u_{i-1}^n| - \gamma \Delta t \frac{\beta_i}{h_i} |u_{i+1}^n| + \Delta t |f_i^n|$$

soit

$$|u_i^{n+1}| \leq \left(1 - \gamma \Delta t \frac{\alpha_i}{h_i}\right) \|\mathbf{u}^n\|_\infty - \gamma \Delta t \frac{\beta_{i-1}}{h_i} \|\mathbf{u}^n\|_\infty - \gamma \Delta t \frac{\beta_i}{h_i} \|\mathbf{u}^n\|_\infty + \Delta t \|f^n\|_\infty,$$

pour tout  $i = 1, \dots, N$ . On a  $\alpha_i + \beta_{i-1} + \beta_i = 0$  pour  $1 < i < N$ ,  $\beta_0 = \beta_N = 0$  et  $\alpha_1 + \beta_1 > 0$ ,  $\alpha_N + \beta_N > 0$ . On obtient donc

$$\|\mathbf{u}^{n+1}\|_\infty \leq \|\mathbf{u}^n\|_\infty + \Delta t \|f^n\|_\infty,$$

ce qui implique (2.19).  $\square$

### Convergence

On note  $h = \max_{1 \leq i \leq N} h_i$  et on introduit la norme  $L^2$  discrète définie par

$$\|\mathbf{v}\|_{0,h} = \left( \sum_{i=1}^N h_i v_i^2 \right)^{1/2}, \quad (2.20)$$

pour  $\mathbf{v} = (v_1, \dots, v_N)^\top$ . La norme  $L^2$  discrète  $\|\cdot\|_{0,h}$  possède la propriété<sup>1</sup> :  $\|\mathbf{v}\|_{0,h} \xrightarrow{h \rightarrow 0} \|v\|_{L^2(0,1)}$  pour toute fonction  $v \in L^2(0,1)$  continue sur  $(0,1)$  avec  $\mathbf{v} = (v(x_1), \dots, v(x_N))^\top$ .

**Proposition 2.2** Soit  $u_0 \in C^2([0,1])$  et  $f \in C^0([0,1])$ . On suppose que (P) admet une solution  $u \in C^2([0,1] \times [0,T])$ . On note  $(u_i^n)_{\substack{1 \leq i \leq N \\ 0 \leq n \leq M}}$  la solution du schéma Volumes Finis explicite (2.16) et on pose  $e_i^n = u(x_i, t^n) - u_i^n$  pour tout  $1 \leq i \leq N$ ,  $0 \leq n \leq M$ . On suppose que  $\lambda = \gamma \Delta t \max_{1 \leq i \leq N} \frac{1}{h_i} \left( \frac{1}{h_{i-\frac{1}{2}}} + \frac{1}{h_{i+\frac{1}{2}}} \right) \leq 1$ . Alors il existe une constante  $C > 0$  indépendante de  $h$  et  $\Delta t$  telle que

$$\|\mathbf{e}^n\|_{0,h} \leq C(h + \Delta t)$$

avec  $\mathbf{e}^n = (e_1^n, \dots, e_N^n)^\top$ , pour tout  $n = 0, \dots, M$ .

*Démonstration.* A faire.

#### 2.2.2 Schéma d'Euler implicite

On écrit l'équation (2.10) à l'instant  $t^{n+1}$  et on approche la dérivée en temps

$$\frac{\partial u_i}{\partial t}(t^{n+1}) \simeq \frac{u_i^{n+1} - u_i^n}{\Delta t}.$$

Le schéma Volumes Finis s'écrit alors, pour  $1 \leq i \leq N$ ,  $0 \leq n \leq M$  :

$$\frac{u_i^{n+1} - u_i^n}{\Delta t} + \frac{F_{i+\frac{1}{2}}^{n+1} - F_{i-\frac{1}{2}}^{n+1}}{h_i} = f_i^{n+1} \quad (2.21)$$

1. En effet, on a

$$\begin{aligned} \|v\|_{L^2(0,1)}^2 &= \sum_{i=1}^N \int_{K_i} v^2(x) dx = \sum_{i=1}^N |K_i| v^2(x_i) + \mathcal{O}(h) \quad (\text{formule des rectangles pour } v \text{ de classe } C^1 \text{ sur } [0,1]) \\ &= \sum_{i=1}^N h_i v_i^2 + \mathcal{O}(h) \end{aligned}$$

où  $f_i^{n+1} = f_i(t^{n+1}) = \frac{1}{h_i} \int_{K_i} f(x, t^{n+1}) dx$  et avec (cf. (2.11))

$$F_{i+\frac{1}{2}}^{n+1} = -\gamma \frac{u_{i+1}^{n+1} - u_i^{n+1}}{h_{i+\frac{1}{2}}}.$$

### Forme matricielle

Pour  $1 \leq i \leq N$ ,  $0 \leq n \leq M$ , on obtient

$$-\gamma \frac{\Delta t}{h_i h_{i-\frac{1}{2}}} u_{i-1}^{n+1} + \left( 1 + \gamma \frac{\Delta t}{h_i} \left( \frac{1}{h_{i-\frac{1}{2}}} + \frac{1}{h_{i+\frac{1}{2}}} \right) \right) u_i^{n+1} - \gamma \frac{\Delta t}{h_i h_{i+\frac{1}{2}}} u_{i+1}^{n+1} = u_i^n + \Delta t f_i^{n+1}.$$

Les conditions limites et initiales sont données par

$$\begin{aligned} u_0^n &= u_{N+1}^n = 0, \quad n = 0, \dots, M \\ u_i^0 &= \frac{1}{h_i} \int_{K_i} u_0(x) dx, \quad i = 1, \dots, N \end{aligned} \quad (2.22)$$

Le schéma "Volumes Finis" Euler implicite s'écrit donc

$$\gamma \Delta t \frac{\beta_{i-1}}{h_i} u_{i-1}^{n+1} + \left( 1 + \gamma \Delta t \frac{\alpha_i}{h_i} \right) u_i^{n+1} + \gamma \Delta t \frac{\beta_i}{h_i} u_{i+1}^{n+1} = u_i^n + \Delta t f_i^n \quad (2.23)$$

avec (cf. (2.14))

$$\alpha_i = \frac{1}{h_{i-\frac{1}{2}}} + \frac{1}{h_{i+\frac{1}{2}}} > 0, \quad \beta_i = -\frac{1}{h_{i+\frac{1}{2}}} < 0, \quad (2.24)$$

pour  $1 \leq i \leq N$  et  $0 \leq n \leq M$ .

On pose  $\mathbf{u}^n = (u_1^n, \dots, u_N^n)^\top \in \mathbb{R}^N$  et  $\mathbf{f}^n = (f_1^n, \dots, f_N^n)^\top \in \mathbb{R}^N$ . Le schéma d'Euler implicite s'écrit vectoriellement, pour  $0 \leq n \leq M$  :

$$(I_N + \gamma \Delta t H^{-1} A) \mathbf{u}^{n+1} = \mathbf{u}^n + \Delta t \mathbf{f}^{n+1}, \quad (2.25)$$

où  $A$  est la matrice définie par (2.17) et la matrice diagonale  $H$  est donnée par (2.18). On remarquera que les matrices  $A$  et  $H^{-1}$  sont symétriques mais le produit  $H^{-1}A$  ne l'est pas. On préfère alors écrire le système (2.25) sous la forme

$$(H + \gamma \Delta t A) \mathbf{u}^{n+1} = H (\mathbf{u}^n + \Delta t \mathbf{f}^{n+1}), \quad (2.26)$$

La matrice  $H + \gamma \Delta t A$  (de même que la matrice  $I_N + \gamma \Delta t H^{-1}A$ ) est à diagonale strictement dominante ( $\gamma > 0$ ) donc elle est inversible (cf. Annexe B). Par ailleurs, le schéma est implicite : une fois qu'on connaît  $\mathbf{u}^n$ , on détermine  $\mathbf{u}^{n+1}$  en résolvant le système (2.26).

### Stabilité

#### Proposition 2.3

1. Principe du maximum discret. La matrice  $H + \gamma \Delta t A$  est une  $M$ -matrice. Si  $f \geq 0$  et  $u_0 \geq 0$  alors  $\mathbf{u}^n \geq \mathbf{0}$  pour tout  $1 \leq n \leq M$ .

2. Stabilité  $L^\infty$ .

Pour tout  $1 \leq n \leq M$ ,

$$\|\mathbf{u}^n\|_\infty \leq \|\mathbf{u}^0\|_\infty + T \|f\|_{L^\infty((0,1) \times (0,T))} \quad (2.27)$$

*Remarque.* Contrairement au schéma explicite, il n'y a pas de condition de stabilité sur  $\Delta t$  et  $h$ . En revanche, il faut résoudre un système linéaire.

*Démonstration.*

1. Soit  $\mathcal{M} = H + \gamma\Delta t A$ . On note  $(m_{ij})$  les coefficients de  $\mathcal{M}$ . Pour tout  $1 \leq i, j \leq N$ , on a clairement  $m_{ij} \leq 0$  pour  $i \neq j$  et  $\sum_{j=1}^N m_{ij} > 0$ . Par conséquent (cf. Annexe B),  $\mathcal{M}$  est une M-matrice et en particulier elle est monotone. Par récurrence sur  $n$ , on en déduit que  $\mathbf{u}^n \geq \mathbf{0}$  pour tout  $n$ , dès que  $f \geq 0$  et  $u_0 \geq 0$ .
2. A partir de (2.23), on obtient, pour tout  $1 \leq i \leq N$ ,

$$\begin{aligned} \left(1 + \gamma\Delta t \frac{\alpha_i}{h_i}\right) |u_i^{n+1}| &\leq -\gamma\Delta t \frac{\beta_{i-1}}{h_i} |u_{i-1}^{n+1}| - \gamma\Delta t \frac{\beta_i}{h_i} |u_{i+1}^{n+1}| + u_i^n + \Delta t |f_i^{n+1}| \\ &\leq -\gamma\Delta t \frac{\beta_{i-1}}{h_i} \|\mathbf{u}^{n+1}\|_\infty - \gamma\Delta t \frac{\beta_i}{h_i} \|\mathbf{u}^{n+1}\|_\infty + \|\mathbf{u}^n\|_\infty + \Delta t \|\mathbf{f}^{n+1}\|_\infty \end{aligned}$$

Soit  $i_0$  tel que  $|u_{i_0}^{n+1}| = \max_{1 \leq i \leq N} |u_i^{n+1}| = \|\mathbf{u}^{n+1}\|_\infty$ . On a alors

$$\left(1 + \gamma\Delta t \frac{\alpha_{i_0}}{h_{i_0}}\right) \underbrace{|u_{i_0}^{n+1}|}_{=\|\mathbf{u}^{n+1}\|_\infty} \leq -\gamma\Delta t \frac{\beta_{i_0-1}}{h_{i_0}} \|\mathbf{u}^{n+1}\|_\infty - \gamma\Delta t \frac{\beta_{i_0}}{h_{i_0}} \|\mathbf{u}^{n+1}\|_\infty + \|\mathbf{u}^n\|_\infty + \Delta t \|\mathbf{f}^{n+1}\|_\infty.$$

On a toujours  $\alpha_{i_0} + \beta_{i_0-1} + \beta_{i_0} > 0$ , on obtient donc

$$\|\mathbf{u}^{n+1}\|_\infty \leq \|\mathbf{u}^n\|_\infty + \Delta t \|\mathbf{f}^{n+1}\|_\infty$$

pour tout  $0 \leq n \leq M$ . On en déduit que  $\|\mathbf{u}^n\|_\infty \leq \|\mathbf{u}^0\|_\infty + T\|f\|_{L^\infty((0,1) \times (0,T))}$

□

## Convergence

On note  $h = \max_{1 \leq i \leq N} h_i$ .

**Proposition 2.4** *Soit  $u_0 \in C^2([0, 1])$  et  $f \in C^0([0, 1] \times [0, T])$ . On suppose que (P) admet une solution  $u \in C^2([0, 1] \times [0, T])$ . On note  $(u_i^n)_{\substack{1 \leq i \leq N \\ 0 \leq n \leq M}}$  la solution du schéma Volumes Finis implicite (2.25) et on pose  $e_i^n = u(x_i, t^n) - u_i^n$  pour tout  $1 \leq i \leq N$ ,  $0 \leq n \leq M$ . Alors il existe une constante  $C > 0$  indépendante de  $h$  et  $\Delta t$  telle que*

$$\|\mathbf{e}^n\|_{0,h} \leq C(h + \Delta t)$$

avec  $\mathbf{e}^n = (e_1^n, \dots, e_N^n)^\top$  pour tout  $n = 0, \dots, M$ .

*Remarque.* Il n'y a pas de condition de stabilité portant sur les pas  $\Delta t$  et  $h$ .

*Démonstration.* D'après (2.10) à l'instant  $t = t^{n+1}$ , on a :

$$\frac{d}{dt} u_i(t^{n+1}) + \frac{\mathcal{F}_{i+\frac{1}{2}}(t^{n+1}) - \mathcal{F}_{i-\frac{1}{2}}(t^{n+1})}{h_i} = f_i(t^{n+1}), \quad (2.28)$$

avec  $u_i(t^{n+1}) = \frac{1}{h_i} \int_{K_i} u(x, t^{n+1}) dx$  et  $\mathcal{F}_{i+\frac{1}{2}}(t^{n+1}) = -\gamma \frac{\partial u}{\partial x}(x_{i+\frac{1}{2}}, t^{n+1})$  le flux exact de  $u$  en  $x_{i+\frac{1}{2}}$  à l'instant  $t^{n+1}$ . On a  $\frac{d}{dt} u_i(t^{n+1}) = \frac{1}{h_i} \int_{K_i} \frac{\partial u}{\partial t}(x, t^{n+1}) dx$  et, pour  $x \in K_i$  :

$$\begin{aligned} \frac{\partial u}{\partial t}(x, t^{n+1}) &= \frac{\partial u}{\partial t}(x_i, t^{n+1}) + (x - x_i) \frac{\partial^2 u}{\partial t}(\theta, t^{n+1}), \quad \theta \text{ compris entre } x \text{ et } x_i \\ &= \frac{\partial u}{\partial t}(x_i, t^{n+1}) + \mathcal{O}(h), \end{aligned}$$

ce qui donne

$$\frac{d}{dt} u_i(t^{n+1}) = \frac{\partial u}{\partial t}(x_i, t^{n+1}) + \mathcal{O}(h).$$

On obtient donc

$$\frac{d}{dt}u_i(t^{n+1}) = \frac{u(x_i, t^{n+1}) - u(x_i, t^n)}{\Delta t} + \mathcal{O}(h + \Delta t). \quad (2.29)$$

Par ailleurs, on montre que (voir aussi Chapitre 1, exercice 1, série I)

$$\begin{aligned} \mathcal{F}_{i+\frac{1}{2}}(t^{n+1}) &= -\gamma \frac{\partial u}{\partial x}(x_{i+\frac{1}{2}}, t^{n+1}) \\ &= -\gamma \frac{u(x_{i+1}, t^{n+1}) - u(x_i, t^{n+1})}{h_{i+\frac{1}{2}}} + R_{i+\frac{1}{2}} \end{aligned}$$

avec

$$|R_{i+\frac{1}{2}}| \leq h_{i+\frac{1}{2}} \left\| \frac{\partial^2 u}{\partial x^2} \right\|_{L^\infty((0,1) \times (0,T))}. \quad (2.30)$$

L'équation (2.28) s'écrit donc

$$\begin{aligned} u(x_i, t^{n+1}) - u(x_i, t^n) + \gamma \frac{\Delta t}{h_i} \left( -\frac{(u(x_{i+1}, t^{n+1}) - u(x_i, t^{n+1}))}{h_{i+\frac{1}{2}}} + \frac{(u(x_i, t^{n+1}) - u(x_{i-1}, t^{n+1}))}{h_{i-\frac{1}{2}}} \right) \\ = \Delta t f_i^{n+1} + \frac{\Delta t}{h_i} (R_{i+\frac{1}{2}} - R_{i-\frac{1}{2}}) + \mathcal{O}(h + \Delta t)\Delta t. \end{aligned} \quad (2.31)$$

Par ailleurs, le schéma "Volumes Finis" implicite (2.13) s'écrit sous la forme

$$u_i^{n+1} - u_i^n + \gamma \frac{\Delta t}{h_i} \left( -\frac{(u_{i+1}^{n+1} - u_i^{n+1})}{h_{i+\frac{1}{2}}} + \frac{(u_i^{n+1} - u_{i-1}^{n+1})}{h_{i-\frac{1}{2}}} \right) = \Delta t f_i^{n+1}. \quad (2.32)$$

En considérant l'erreur  $e_i^n = u(x_i, t^n) - u_i^n$ , on obtient à partir de la différence entre (2.31) et (2.32) :

$$e_i^{n+1} - e_i^n + \gamma \frac{\Delta t}{h_i} \left( -\frac{(e_{i+1}^{n+1} - e_i^{n+1})}{h_{i+\frac{1}{2}}} + \frac{(e_i^{n+1} - e_{i-1}^{n+1})}{h_{i-\frac{1}{2}}} \right) = \frac{\Delta t}{h_i} (R_{i+\frac{1}{2}} - R_{i-\frac{1}{2}}) + \mathcal{O}(h + \Delta t)\Delta t, \quad (2.33)$$

pour  $1 \leq i \leq N$  et  $0 \leq n \leq M$ . On multiplie l'équation (2.33) par  $h_i e_i^{n+1}$  et on somme sur  $i = 1, \dots, N$  :

$$\begin{aligned} \sum_{i=1}^N h_i (e_i^{n+1} - e_i^n) e_i^{n+1} + \gamma \Delta t \left( -\sum_{i=1}^N \frac{(e_{i+1}^{n+1} - e_i^{n+1})}{h_{i+\frac{1}{2}}} e_i^{n+1} + \sum_{i=1}^N \frac{(e_i^{n+1} - e_{i-1}^{n+1})}{h_{i-\frac{1}{2}}} e_i^{n+1} \right) \\ = \Delta t \sum_{i=1}^N (R_{i+\frac{1}{2}} - R_{i-\frac{1}{2}}) e_i^{n+1} + \mathcal{O}(h + \Delta t)\Delta t \sum_{i=1}^N h_i e_i^{n+1}. \end{aligned} \quad (2.34)$$

A partir de l'identité  $2ab = a^2 + b^2 - (a - b)^2$ , on obtient

$$2(e_i^{n+1} - e_i^n) e_i^{n+1} = |e_i^{n+1} - e_i^n|^2 + |e_i^{n+1}|^2 - |e_i^n|^2.$$

En utilisant le fait que  $e_0^{n+1} = e_{N+1}^{n+1} = 0$ , l'équation (2.34) devient

$$\begin{aligned} \|\mathbf{e}^{n+1}\|_{0,h}^2 + \|\mathbf{e}^{n+1} - \mathbf{e}^n\|_{0,h}^2 - \|\mathbf{e}^n\|_{0,h}^2 + 2\gamma \Delta t \left( -\sum_{i=0}^N \frac{(e_{i+1}^{n+1} - e_i^{n+1})}{h_{i+\frac{1}{2}}} e_i^{n+1} + \sum_{i=0}^N \frac{(e_{i+1}^{n+1} - e_i^{n+1})}{h_{i+\frac{1}{2}}} e_{i+1}^{n+1} \right) \\ = 2\Delta t \left( \sum_{i=0}^N R_{i+\frac{1}{2}} e_i^{n+1} - \sum_{i=1}^N R_{i-\frac{1}{2}} e_i^{n+1} \right) + \mathcal{O}(h + \Delta t)\Delta t \sum_{i=1}^N h_i e_i^{n+1}. \end{aligned}$$

Ainsi on a

$$\begin{aligned} \|\mathbf{e}^{n+1}\|_{0,h}^2 + \|\mathbf{e}^{n+1} - \mathbf{e}^n\|_{0,h}^2 - \|\mathbf{e}^n\|_{0,h}^2 + 2\gamma \Delta t (D_{n+1})^2 \\ = 2\Delta t \left( \sum_{i=0}^N R_{i+\frac{1}{2}} (e_i^{n+1} - e_{i+1}^{n+1}) \right) + \mathcal{O}(h + \Delta t)\Delta t \sum_{i=1}^N h_i e_i^{n+1}. \end{aligned} \quad (2.35)$$

avec

$$(D_{n+1})^2 = \sum_{i=0}^N \frac{(e_{i+1}^{n+1} - e_i^{n+1})^2}{h_{i+\frac{1}{2}}}. \quad (2.36)$$

Par Cauchy-Schwarz,

$$\bullet \sum_{i=1}^N h_i e_i^{n+1} \leq \left( \sum_{i=1}^N h_i \right)^{1/2} \left( \sum_{i=1}^N h_i (e_i^{n+1})^2 \right)^{1/2}. \text{ Or, on a } \sum_{i=1}^N h_i = 1, \text{ donc}$$

$$\sum_{i=1}^N h_i e_i^{n+1} \leq \|\mathbf{e}^{n+1}\|_{0,h}.$$

$$\begin{aligned} \bullet \sum_{i=0}^N R_{i+\frac{1}{2}}(e_i^{n+1} - e_{i+1}^{n+1}) &\leq \left( \sum_{i=0}^N h_{i+\frac{1}{2}} R_{i+\frac{1}{2}}^2 \right)^{1/2} \left( \sum_{i=0}^N \frac{(e_{i+1}^{n+1} - e_i^{n+1})^2}{h_{i+\frac{1}{2}}} \right)^{1/2} \\ &\leq C \left( \sum_{i=0}^N h_{i+\frac{1}{2}}^3 \right)^{1/2} D_{n+1} \text{ d'après (2.30) et (2.36)} \\ &\leq Ch \underbrace{\left( \sum_{i=0}^N h_{i+\frac{1}{2}} \right)^{1/2}}_{=1} D_{n+1} \\ &\leq Ch D_{n+1} \end{aligned}$$

On déduit alors de (2.35),

$$\|\mathbf{e}^{n+1}\|_{0,h}^2 - \|\mathbf{e}^n\|_{0,h}^2 + 2\gamma\Delta t(D_{n+1})^2 \leq C\Delta t \left( hD_{n+1} + (h + \Delta t)\|\mathbf{e}^{n+1}\|_{0,h} \right)$$

où  $C > 0$  est une constante indépendante de  $h$  et  $\Delta t$ . En utilisant l'inégalité de Young (cf. Annexe C), on obtient pour tout  $\varepsilon_1, \varepsilon_2 > 0$ ,

$$\|\mathbf{e}^{n+1}\|_{0,h}^2 - \|\mathbf{e}^n\|_{0,h}^2 + 2\gamma\Delta t(D_{n+1})^2 \leq \varepsilon_1\Delta t(D_{n+1})^2 + \frac{C}{\varepsilon_1}\Delta t h^2 + \varepsilon_2\|\mathbf{e}^{n+1}\|_{0,h}^2 + \frac{C}{\varepsilon_2}\Delta t^2(h + \Delta t)^2$$

où  $C > 0$  est une constante indépendante de  $h$  et  $\Delta t$ . On choisit  $\varepsilon_1 = 2\gamma > 0$  et on obtient, pour tout  $\varepsilon_2 > 0$ ,

$$(1 - \varepsilon_2)\|\mathbf{e}^{n+1}\|_{0,h}^2 \leq \|\mathbf{e}^n\|_{0,h}^2 + C\Delta t h^2 + \frac{C}{\varepsilon_2}\Delta t^2(h + \Delta t)^2. \quad (2.37)$$

On choisit alors  $1 - \varepsilon_2 = \frac{1}{1+\Delta t} > 0$  i.e.  $\varepsilon_2 = \frac{\Delta t}{1+\Delta t} > 0$ , de sorte que pour tout  $0 \leq n \leq M$ ,

$$\|\mathbf{e}^{n+1}\|_{0,h}^2 \leq (1 + \Delta t)\|\mathbf{e}^n\|_{0,h}^2 + C\Delta t(h + \Delta t)^2, \quad (2.38)$$

$C > 0$  est une constante indépendante de  $h$  et  $\Delta t$ . D'après le lemme de Gronwall (cf. Annexe C, (C.2)), on obtient, pour tout  $1 \leq n \leq M$ ,

$$\|\mathbf{e}^n\|_{0,h}^2 \leq \left( \|\mathbf{e}^0\|_{0,h}^2 + Cn\Delta t(h + \Delta t)^2 \right) \exp(n\Delta t),$$

Puisque  $n\Delta t \leq T$ , on en déduit

$$\|\mathbf{e}^n\|_{0,h} \leq C' \left( \|\mathbf{e}^0\|_{0,h} + h + \Delta t \right), \quad (2.39)$$

avec  $C' > 0$  indépendante de  $h$  et  $\Delta t$ .

De plus, on a choisi la donnée initiale  $u_i^0 = \frac{1}{h_i} \int_{K_i} u_0(x) dx = u_0(x_i) + \mathcal{O}(h_i)$  et donc  $e_i^0 = \mathcal{O}(h_i)$  pour tout  $1 \leq i \leq N$ . Par conséquent, on a  $\|\mathbf{e}^0\|_{0,h} = \mathcal{O}(h)$  et on obtient

$$\|\mathbf{e}^n\|_{0,h} \leq C(h + \Delta t),$$

avec  $C > 0$  indépendante de  $h$  et  $\Delta t$ . □

## 2.3 Equation de la chaleur en 2D d'espace et Volumes Finis

On considère un domaine polygonal  $\Omega \subset \mathbb{R}^2$ . Pour  $T > 0$ , on cherche une fonction  $u = u(\mathbf{x}, t)$  pour  $\mathbf{x} \in \Omega$  et  $t \in [0, T]$  vérifiant l'équation de la chaleur suivante :

$$\frac{\partial u}{\partial t} - \nu \Delta u = f \quad \text{dans } \Omega \times ]0, T[ \quad (2.40)$$

$$u = g \quad \text{sur } \partial\Omega \times ]0, T[ \quad (2.41)$$

$$u(\cdot, 0) = u_0 \quad \text{dans } \Omega. \quad (2.42)$$

Les fonctions  $f$ ,  $g$ ,  $u_0$  et le paramètre  $\nu > 0$  sont donnés. On associe au domaine  $\Omega$  un maillage "Volumes Finis" admissibles (cf. Chapitre 1, section 1.4). Pour fixer les idées, on considère un maillage dont les volumes de contrôle sont les cellules de Voronoï associées à une triangulation de Delaunay de  $\Omega$ . A chaque cellule  $K$ , on associe les centres  $x_K$  qui sont les sommets des triangles (cf. section 1.4.3).

### 2.3.1 Discrétisation en espace

On suppose la solution  $u$  de (2.40)–(2.42) régulière (cf. section 2.1.1). On intègre l'équation (2.40) sur une cellule  $K$  et en utilisant la formule de la divergence, on obtient :

$$\frac{d}{dt} \int_K u(\mathbf{x}, t) d\mathbf{x} - \nu \int_{\partial K} \nabla u \cdot \mathbf{n}_K d\Gamma = |K| f_K$$

avec  $f_K(t) = \frac{1}{|K|} \int_K f(\mathbf{x}, t) d\mathbf{x}$  et  $\mathbf{n}_K$  désigne la normale unitaire dirigée vers l'extérieur de  $K$ . On pose  $u_K(t) = \frac{1}{|K|} \int_K u(\mathbf{x}, t) d\mathbf{x}$  et on a, pour  $t \in [0, T]$ ,

$$|K| \frac{du_K}{dt}(t) - \nu \sum_{e \in \mathcal{E}_K} \int_e \nabla u \cdot \mathbf{n}_{K,e} d\Gamma = |K| f_K(t), \quad (2.43)$$

où  $\mathcal{E}_K$  désigne l'ensemble des arêtes de la cellule  $K$  et  $\mathbf{n}_{K,e}$  est la normale unitaire à  $e$  dirigée vers l'extérieur de  $K$ .

On approche le flux  $-\int_e \nabla u \cdot \mathbf{n}_{K,e} d\Gamma \simeq F_{K,e}$  avec

$$F_{K,e} = \frac{(u_K - u_L)}{|\mathbf{x}_K - \mathbf{x}_L|} |e| \quad \text{si } e \not\subset \partial\Omega. \quad (2.44)$$

Si  $K$  possède une arête sur le bord  $\partial\Omega$ , on impose

$$u_K = u_e = \frac{1}{|K|} \int_e g(\mathbf{x}, t) d\Gamma \quad \text{si } e \subset \partial\Omega. \quad (2.45)$$

Le schéma Volumes Finis de discrétisation en espace s'écrit ainsi, pour  $t \in [0, T]$ ,

$$|K| \frac{du_K}{dt}(t) + \sum_{\substack{e \in \mathcal{E}_K \\ e \not\subset \partial\Omega}} \nu F_{K,e}(t) = |K| f_K(t). \quad (2.46)$$

### 2.3.2 Schéma explicite en temps

Soit  $n \in \mathbb{N}^*$  et  $\Delta t = \frac{T}{n} > 0$  le pas de discrétisation en temps avec  $t^k = k\Delta t$  pour  $k = 0, 1, \dots, n$ . On note l'approximation  $u_K^k \simeq u_K(t^k) = \frac{1}{|K|} \int_K u(\mathbf{x}, t^k) d\mathbf{x}$ .

Le schéma explicite associé à (2.46) s'écrit, pour tout  $K \in \mathcal{T}$ ,  $k = 0, \dots, n$ ,

$$|K| \frac{(u_K^{k+1} - u_K^k)}{\Delta t} + \sum_{\substack{e \in \mathcal{E}_K \\ e \not\subset \partial\Omega}} \nu F_{K,e}^k = |K| f_K^k, \quad (2.47)$$

avec

$$F_{K,e}^k = \frac{(u_K^k - u_L^k)}{|\mathbf{x}_K - \mathbf{x}_L|} |e| \quad (e \not\subset \partial\Omega). \quad (2.48)$$

On écrit les relations (2.47) pour toute cellule  $K$  *intérieure* i.e. une cellule ne possédant pas d'arête sur le bord  $\partial\Omega$ . Pour une cellule  $K$  ayant une arête sur  $\partial\Omega$ , on impose  $u_K = u_e$  donnée par (2.45). On obtient ainsi le schéma :

— Pour toute cellule  $K$  *intérieure*,

$$u_K^{k+1} = u_K^k - \frac{\nu\Delta t}{|K|} \sum_{\substack{e \in \mathcal{E}_K \\ e=(K|L)}} \frac{|e|}{|\mathbf{x}_K - \mathbf{x}_L|} (u_K^k - u_L^k) + \Delta t f_K^k \quad (2.49)$$

— Pour toute cellule  $K$  possédant une arête  $e \subset \partial\Omega$ , on impose  $u_K^{k+1} = u_e$  où  $u_e$  est donnée par (2.45).

### Système linéaire

On prend (pour simplifier)  $g = 0$ . Soit  $N$  le nombre de cellules de contrôles *intérieures*. On regroupe les inconnues dans le vecteur  $\mathbf{u}^k = (u_{K_1}^k, \dots, u_{K_N}^k)^\top$  et on note  $\mathbf{f}^k = (f_{K_1}^k, \dots, f_{K_N}^k)^\top$ .

Le système linéaire correspondant à (2.49) s'écrit

$$\mathbf{u}^{k+1} = (I_N - \nu\Delta t H^{-1}A)\mathbf{u}^k + \Delta t \mathbf{f}^k, \quad (2.50)$$

où  $I_N$  est la matrice identité d'ordre  $N$  et  $A$  est la matrice de taille  $N \times N$  définie par (1.45) et correspondant au Laplacien. La matrice *diagonale*  $H$  de taille  $N \times N$  est définie par

$$H = \begin{pmatrix} |K_1| & & & 0 \\ & |K_2| & & \\ & & \ddots & \\ 0 & & & |K_N| \end{pmatrix}. \quad (2.51)$$

Le système (2.50) est *explicite* : il n'y a pas de système linéaire à résoudre pour déterminer  $\mathbf{u}^{k+1}$  à partir de  $\mathbf{u}^k$ . Le schéma (2.50) est  $L^\infty$ -stable si  $\|\mathbf{u}^k\|_\infty \leq C$  pour tout  $0 \leq k \leq n$  avec  $C > 0$  une constante indépendante de  $k$  et  $N$ . On montre (cf. exercice) que le schéma (2.50) est  $L^\infty$ -stable sous la condition

$$\lambda = \nu\Delta t \max_{1 \leq i \leq N} \left( \frac{1}{|K_i|} \sum_{\substack{e \in \mathcal{E}_{K_i} \\ e=(K_i|L)}} \frac{|e|}{|\mathbf{x}_{K_i} - \mathbf{x}_L|} \right) \leq 1. \quad (2.52)$$

### 2.3.3 Schéma implicite en temps

Le schéma implicite associé à (2.46) s'écrit, pour tout  $K \in \mathcal{T}$ ,  $k = 0, \dots, n$ ,

$$|K| \frac{(u_K^{k+1} - u_K^k)}{\Delta t} + \sum_{\substack{e \in \mathcal{E}_K \\ e \not\subset \partial\Omega}} \nu F_{K,e}^{k+1} = |K| f_K^{k+1}, \quad (2.53)$$

avec

$$F_{K,e}^{k+1} = \frac{(u_K^{k+1} - u_L^{k+1})}{|\mathbf{x}_K - \mathbf{x}_L|} |e| \quad (e \not\subset \partial\Omega). \quad (2.54)$$

On écrit les relations (2.53) pour toute cellule  $K$  *intérieure* i.e. une cellule ne possédant pas d'arête sur le bord  $\partial\Omega$ . Pour une cellule  $K$  ayant une arête sur  $\partial\Omega$ , on impose  $u_K = u_e$  donnée par (2.45). On obtient ainsi le schéma :

— Pour toute cellule  $K$  *intérieure*,

$$u_K^{k+1} + \frac{\nu \Delta t}{|K|} \sum_{\substack{e \in \mathcal{E}_K \\ e=(K|L)}} \frac{|e|}{|\mathbf{x}_K - \mathbf{x}_L|} (u_K^{k+1} - u_L^{k+1}) = u_K^k + \Delta t f_K^{k+1} \quad (2.55)$$

— Pour toute cellule  $K$  possédant une arête sur le bord  $\partial\Omega$ , on impose  $u_K^{k+1} = u_e$  où  $u_e$  est donnée par (2.45).

### Système linéaire

On prend toujours  $g = 0$  (pour simplifier). Le système linéaire correspondant à (2.55) s'écrit

$$(I_N + \nu \Delta t H^{-1} A) \mathbf{u}^{k+1} = \mathbf{u}^k + \Delta t \mathbf{f}^{k+1} \quad (2.56)$$

avec  $A$  et  $H$  définies par (1.45) et (2.51) respectivement. On écrit le système (2.56) plutôt sous la forme

$$(H + \nu \Delta t A) \mathbf{u}^{k+1} = H (\mathbf{u}^k + \Delta t \mathbf{f}^{k+1}), \quad (2.57)$$

la matrice  $\mathcal{M} = H + \nu \Delta t A$  étant *symétrique* (ceci peut être avantageux pour la résolution numérique du système linéaire (2.57)). On montre également que  $\mathcal{M}$  est une  $M$ -matrice.

## Chapitre 3

# Equation de transport

### 3.1 Introduction

On considère un problème de transport posé dans un domaine  $\Omega \subset \mathbb{R}^2$ . Etant donné un champ de vitesse  $\mathbf{V} : \Omega \rightarrow \mathbb{R}^2$ , on cherche une fonction  $u : \Omega \times \mathbb{R}^+ \rightarrow \mathbb{R}$  vérifiant

$$\frac{\partial u}{\partial t} + \operatorname{div}(\mathbf{V}u) = 0 \quad \text{dans } \Omega \times \mathbb{R}^+ \quad (3.1)$$

$$u(\mathbf{V} \cdot \mathbf{n}) = 0 \quad \text{sur } \partial\Omega \times \mathbb{R}^+ \quad (3.2)$$

$$u(\cdot, 0) = u_0 \quad \text{dans } \Omega \quad (3.3)$$

où  $u_0$  est une fonction donnée et  $\mathbf{n}$  désigne la normale unitaire dirigée vers l'*extérieur* de  $\Omega$ . On suppose que le champ de vitesse  $\mathbf{V}$  est tel que

$$\operatorname{div} \mathbf{V} = 0, \quad (3.4)$$

de sorte que l'équation de transport (3.1) s'écrit encore

$$\frac{\partial u}{\partial t} + \mathbf{V} \cdot \nabla u = 0.$$

La condition limite (3.2) est satisfaite si  $u = 0$  sur  $\partial\Omega$  ou bien si  $\mathbf{V} \cdot \mathbf{n} = 0$  sur  $\partial\Omega$ , auquel cas il n'y a pas besoin de condition limite pour  $u$  sur  $\partial\Omega$ . En fait, pour que l'équation de transport (3.1) avec la donnée initiale (3.3), il suffit d'imposer une condition limite pour  $u$  uniquement sur la partie du bord  $\Gamma = \{\mathbf{x} \in \partial\Omega \text{ tel que } \mathbf{V} \cdot \mathbf{n}(\mathbf{x}) < 0\}$ .

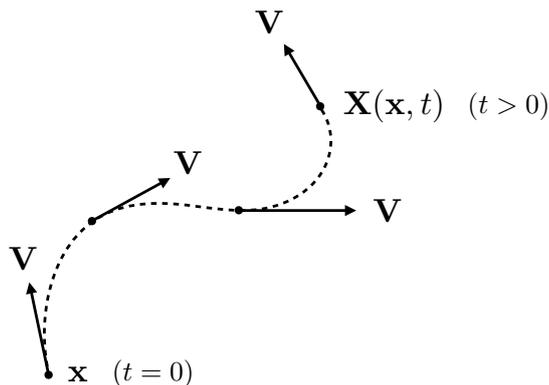
L'équation de transport (3.1) modélise par exemple l'évolution en temps et en espace de la concentration  $u$  (ou densité) d'un composé dans un fluide en mouvement par un champ de vitesse  $\mathbf{V}$  donné. La donnée  $u_0$  représente la concentration initiale dans le fluide à l'instant initial  $t = 0$ . Le système (3.1)–(3.3) représente une *loi de conservation* au sens où la masse totale  $\int_{\Omega} u(\mathbf{x}, t) d\mathbf{x}$  est conservée au cours du temps. En effet, en intégrant l'équation de transport (3.1) sur  $\Omega$  et en utilisant la formule de la divergence avec la condition limite (3.2), on obtient pour tout  $t > 0$ ,

$$\frac{d}{dt} \int_{\Omega} u(\mathbf{x}, t) d\mathbf{x} = 0.$$

Par ailleurs, la terminologie "transport" vient du fait que l'équation (3.1) "transporte" la fonction  $u$  sous l'action du champ de vitesse  $\mathbf{V}$ . En effet, considérons le flot  $\mathbf{X} : \Omega \times \mathbb{R}^+ \rightarrow \mathbb{R}^2$  associé à  $\mathbf{V}$ , solution de

$$\begin{aligned} \frac{d\mathbf{X}}{dt}(\mathbf{x}, t) &= \mathbf{V}(\mathbf{X}(\mathbf{x}, t)), \quad t > 0 \\ \mathbf{X}(\mathbf{x}, 0) &= \mathbf{x} \in \Omega \end{aligned} \quad (3.5)$$

Sous l'action de  $\mathbf{V}$ ,  $\mathbf{X}(\mathbf{x}, t)$  représente la position à l'instant  $t$  d'une particule qui était en  $\mathbf{x} \in \Omega$  à l'instant  $t = 0$  (cf. Figure 3.1).

FIGURE 3.1 – Flot associé au champ de vitesse  $\mathbf{V}$ 

La fonction  $u$  est alors constante le long du flot. En effet, on a

$$\begin{aligned} \frac{d}{dt}[u(\mathbf{X}(\mathbf{x}, t), t)] &= \frac{\partial u}{\partial t}(\mathbf{X}(\mathbf{x}, t), t) + \frac{d\mathbf{X}}{dt} \cdot \nabla u(\mathbf{X}(\mathbf{x}, t), t) \\ &= \left( \frac{\partial u}{\partial t} + \mathbf{V} \cdot \nabla u \right) (\mathbf{X}(\mathbf{x}, t), t) = 0 \end{aligned}$$

Par conséquent, on obtient bien

$$u(\mathbf{X}(\mathbf{x}, t), t) = u(\mathbf{x}, 0) = u_0(\mathbf{x}), \quad \text{pour tout } t \geq 0.$$

La valeur initiale  $u_0(\mathbf{x})$  est ainsi transportée le long de la caractéristique  $\mathbf{X}(\mathbf{x}, \cdot)$

## 3.2 Maillage

Le maillage  $\mathcal{T}$  "Volumes Finis" est constitué de volumes de contrôle  $K$  qui sont des polygones formant une partition de  $\Omega$  :

$$\bar{\Omega} = \cup_{K \in \mathcal{T}} \bar{K} \quad \text{et} \quad \overset{\circ}{K} \cap \overset{\circ}{L} = \emptyset, \quad \forall K, L \in \mathcal{T}, K \neq L.$$

### 1. Schéma "cell center"

Le maillage  $\mathcal{T}$  est une triangulation du domaine  $\Omega$ . Les volumes de contrôle  $K$  sont donc les triangles de cette triangulation. Les centres  $\mathbf{x}_K$  sont les barycentres des triangles.

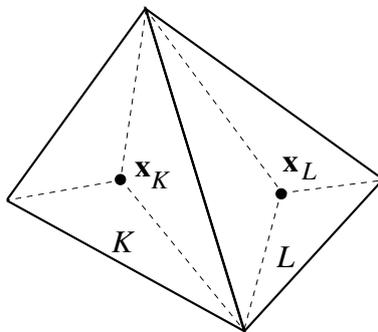


FIGURE 3.2 – Maillage "cell center"

La triangulation  $\mathcal{T}$  de  $\Omega$  doit être admissible au sens des Eléments Finis. En revanche, contrairement au cas des équations elliptiques, on ne demande pas que les angles des triangles soient plus petits que  $\pi/2$ , puisque les centres de cellules (les triangles) sont ici les barycentres.

## 2. Schéma "cell vertex"

Dans de nombreuses applications, les quantités sont définies *aux sommets des triangles* d'une triangulation du domaine  $\Omega$ , plutôt qu'en leur barycentre. Dans le schéma "cell vertex" on utilise toujours une triangulation de  $\Omega$ . Les centres des cellules sont choisis comme les sommets des triangles de cette triangulation. Les cellules  $K$  du maillage "Volumes Finis" de centre  $\mathbf{x}_K$ , sont obtenues en joignant - pour chaque triangle ayant en commun le sommet  $\mathbf{x}_K$  - le point milieu de chaque arête issue de  $\mathbf{x}_K$  au barycentre.

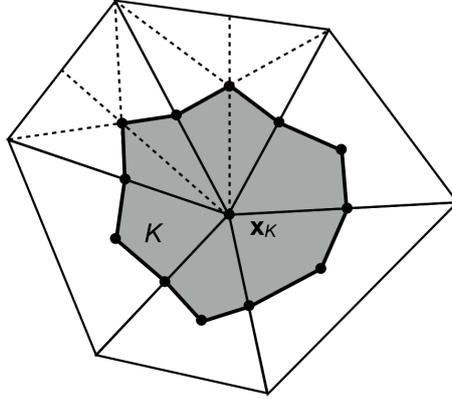


FIGURE 3.3 – Maillage "cell vertex"

### Variantes :

- Les cellules  $K$  sont obtenues en joignant les barycentres des triangles ayant le centre  $\mathbf{x}_K$  comme sommet.
- Au lieu de considérer les barycentres, on considère les centres des cercles circonscrits que l'on joint entre eux. On obtient ainsi des cellules de Voronoï.

## 3.3 Formulation en Volumes Finis

On notera dorénavant

$$\mathbf{F}(u) = \mathbf{V}u. \quad (3.6)$$

Soit  $K$  une cellule d'un maillage  $\mathcal{T}$  "Volumes Finis". On intègre l'équation différentielle (3.1) sur  $K$ .

$$\int_K \left( \frac{\partial u}{\partial t} + \operatorname{div} \mathbf{F}(u) \right) d\mathbf{x} = 0.$$

Par la formule de la divergence, on obtient

$$\frac{d}{dt} \int_K u(\mathbf{x}, t) d\mathbf{x} + \int_{\partial K} \mathbf{F}(u(\mathbf{x}, t)) \cdot \mathbf{n}_K d\Gamma = 0 \quad (3.7)$$

où  $\mathbf{n}_K$  est la normale unitaire dirigée vers l'extérieur de  $K$ . On remarque que la condition limite (3.2) implique  $\mathbf{F}(u) \cdot \mathbf{n}_K = 0$  sur une arête du bord  $e \subset \partial K \cap \partial\Omega$ . Par conséquent, les contributions des termes de bord dans (3.7) proviennent uniquement des arêtes de  $K$  adjacentes à deux cellules. On obtient ainsi

$$\frac{d}{dt} \int_K u(\mathbf{x}, t) d\mathbf{x} + \sum_{\substack{e_{K,L} \subset \partial K \\ e_{K,L} = (K|L)}} \int_{e_{K,L}} \mathbf{F}(u(\mathbf{x}, t)) \cdot \mathbf{n}_{K,L} d\Gamma = 0 \quad (3.8)$$

où  $\mathbf{n}_{K,L}$  désigne la normale unitaire à  $e_{K,L} \not\subset \partial\Omega$  dirigée de  $K$  vers  $L$ . La figure ci-dessous montre un exemple de Volumes de contrôle par un maillage "cell center".

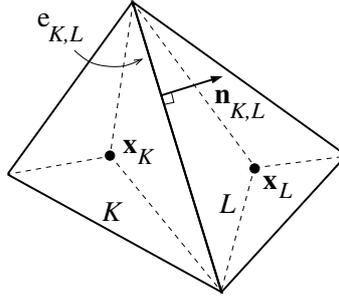


FIGURE 3.4 – Volumes de contrôle pour l'équation de transport.

On note

$$u_K(t) = \frac{1}{|K|} \int_K u(\mathbf{x}, t) d\mathbf{x}$$

et on a

$$\frac{d}{dt} \int_K u(\mathbf{x}, t) d\mathbf{x} = |K| \frac{du_K}{dt}(t). \quad (3.9)$$

Le schéma (3.8) s'écrit

$$\frac{du_K}{dt}(t) + \frac{1}{|K|} \sum_{\substack{e_{K,L} \subset \partial K \\ e_{K,L} = (K|L)}} \int_{e_{K,L}} \mathbf{F}(u(\mathbf{x}, t)) \cdot \mathbf{n}_{K,L} d\Gamma = 0 \quad (3.10)$$

Le flux à travers une arête  $e_{K,L}$  est approchée par :

$$\int_{e_{K,L}} \mathbf{F}(u) \cdot \mathbf{n}_{K,L} d\Gamma \simeq |e_{K,L}| \Phi(u_K, u_L, n_{K,L}) \quad (3.11)$$

où  $\Phi(u_K, u_L, n_{K,L})$  est le *flux numérique* à travers l'arête  $e_{K,L}$ , associé à la cellule  $K$  et  $\Phi(u_K, u_L, n_{K,L})$  est une approximation de  $\mathbf{F}(u) \cdot \mathbf{n}_{K,L}$  sur l'arête  $e_{K,L}$ .

On introduit les instants  $t^n = n\Delta t$  avec  $n$  entier positif et  $\Delta t > 0$  le pas de discrétisation en temps. On considère alors les approximations

$$u_K^n \simeq u_K(t^n) = \frac{1}{|K|} \int_K u(\mathbf{x}, t^n) d\mathbf{x}.$$

En approchant la dérivée en temps dans (3.9) par un schéma d'Euler *explícite*, on obtient le schéma "Volumes Finis" suivant :

- $(u_K^{n+1} - u_K^n) + \frac{\Delta t}{|K|} \sum_{\substack{e_{K,L} \subset \partial K \\ e_{K,L} = (K|L)}} |e_{K,L}| \Phi(u_K^n, u_L^n, \mathbf{n}_{K,L}) = 0, \quad \forall K \in \mathcal{T}, \forall n \in \mathbb{N}$
- $u_K^0 = \frac{1}{|K|} \int_K u_0(\mathbf{x}) d\mathbf{x}, \quad \forall K \in \mathcal{T}$

(3.12)

Le flux numérique  $\Phi$  est choisi par décentrement (*upwind*) :

$$\Phi(u_K, u_L, \mathbf{n}_{K,L}) = u_K (\bar{\mathbf{V}} \cdot \mathbf{n}_{K,L})^+ + u_L (\bar{\mathbf{V}} \cdot \mathbf{n}_{K,L})^- \quad (3.13)$$

avec

$$(\bar{\mathbf{V}} \cdot \mathbf{n}_e)^+ = \max(\bar{\mathbf{V}} \cdot \mathbf{n}_e, 0), \quad (\bar{\mathbf{V}} \cdot \mathbf{n}_e)^- = \min(\bar{\mathbf{V}} \cdot \mathbf{n}_e, 0)$$

et

$$\bar{\mathbf{V}} = \frac{1}{|e|} \int_e \mathbf{V} d\Gamma. \quad (3.14)$$

La relation (3.13) qui définit le flux numérique est équivalente à

$$\Phi(u_K, u_L, \mathbf{n}_{K,L}) = \begin{cases} u_K(\bar{\mathbf{V}} \cdot \mathbf{n}_{K,L}) & \text{si } \bar{\mathbf{V}} \cdot \mathbf{n}_{K,L} > 0 \\ u_L(\bar{\mathbf{V}} \cdot \mathbf{n}_{K,L}) & \text{sinon.} \end{cases} \quad (3.15)$$

Pour calculer le flux numérique, on utilise l'information venant de la cellule  $K$  ou  $L$ , en fonction de la direction du champ de vitesse  $\mathbf{V}$ . C'est le décentrement (cf. Figure3.5).

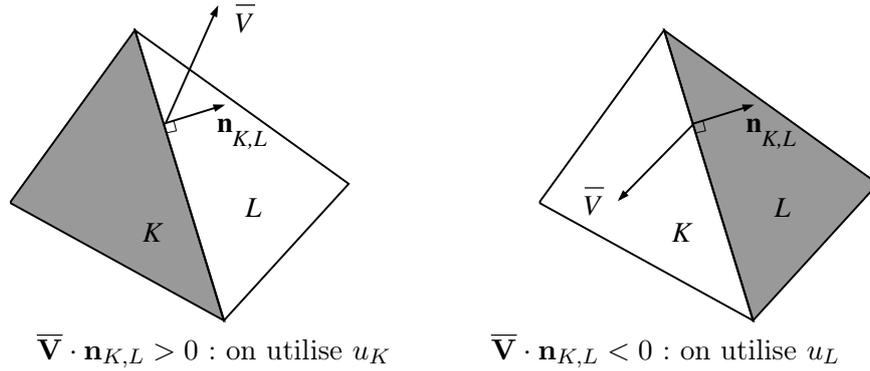


FIGURE 3.5 – Décentrement du flux numérique.

Sur chaque arête  $e$ , on approche la vitesse  $\bar{\mathbf{V}}$  par :

$$\bar{\mathbf{V}} \simeq \frac{\mathbf{V}(\mathbf{x}_1) + \mathbf{V}(\mathbf{x}_2)}{2} \quad \text{où } \mathbf{x}_1 \text{ et } \mathbf{x}_2 \text{ sont les deux extrémités de l'arête } e \quad (3.16)$$

Le flux numérique vérifie les propriétés suivantes :

1. *Conservation* des flux à travers les arêtes.

$$\Phi(u_K, u_L, \mathbf{n}_{K,L}) = -\Phi(u_L, u_K, -\mathbf{n}_{K,L}). \quad (3.17)$$

2. *Consistance* des flux.

$$\Phi(s, s, \mathbf{n}_e) = \frac{1}{|e|} \int_e F(s) \cdot \mathbf{n}_e d\Gamma \quad \text{pour tout } s \in \mathbb{R}. \quad (3.18)$$

### 3.4 Système linéaire

On choisit un maillage "cell center" (cf. Section 3.2). Les cellules de contrôle sont les triangles  $K_i = T$  d'une triangulation du domaine  $\Omega$  admissible au sens des Eléments Finis. Les centres  $\mathbf{x}_{K_i}$  des cellules sont les barycentres des triangles. On note pour simplifier,  $u_i^n = u_{K_i}^n$  l'inconnue associée à la cellule (triangle)  $K_i$ . De même, on note désormais  $e_{ij}$  l'arête commune à deux cellules (triangles)  $K_i$  et  $K_j$  et  $\mathbf{n}_{ij}$  désigne la normale unitaire à l'arête  $e_{ij}$  dirigée vers l'extérieur de  $K_i$ .

Soit  $N$  le nombre de cellules de contrôles (= nombre de triangles de la triangulation). Le schéma (3.12) avec (3.13) s'écrit alors, pour  $i = 1, \dots, N$ ,

$$u_i^{n+1} = \left(1 - \frac{\Delta t}{|K_i|} \sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i | K_j)}} |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^+ \right) u_i^n - \frac{\Delta t}{|K_i|} \sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i | K_j)}} |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^- u_j^n \quad (3.19)$$

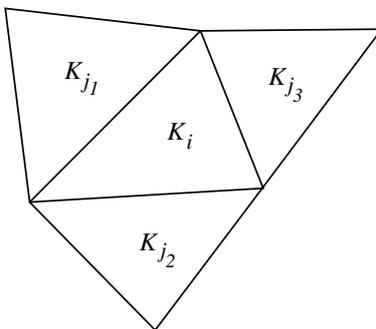
On pose  $\mathbf{u}^n = (u_1^n, \dots, u_N^n)^\top$  et le système linéaire correspondant au schéma explicite, s'écrit :

$$\mathbf{u}^{n+1} = (I_N - \Delta t A) \mathbf{u}^n \quad (3.20)$$

où  $A$  est la matrice de taille  $N \times N$ , de la forme

$$A = \left( \begin{array}{cccc} & j_1 & j_2 & i & j_3 \\ \cdots & \frac{\alpha_{ij_1}^-}{|K_i|} & \cdots & \frac{\alpha_{ij_2}^-}{|K_i|} & \cdots & \frac{1}{|K_i|} \sum_{e_{ij} \subset \partial K_i} \alpha_{ij}^+ & \cdots & \frac{\alpha_{ij_3}^-}{|K_i|} & \cdots \end{array} \right) \leftarrow \text{ligne } i \quad (3.21)$$

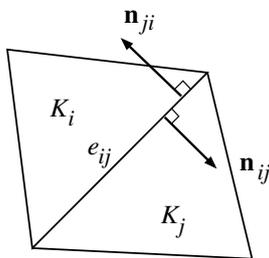
avec  $\alpha_{ij}^\pm = |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^\pm$  et  $e_{ij} = (K_i | K_j)$ . Les trois cellules  $K_{j_1}$ ,  $K_{j_2}$ ,  $K_{j_3}$  sont les trois triangles adjacents au triangle  $K_i$  (cf. Figure ci-dessous).



**Assemblage** Pour construire la matrice  $A$ , on boucle sur les arêtes *intérieures* des triangles. Pour une arête courante  $e \not\subset \partial\Omega$ , on ajoute les contributions des deux triangles  $K_i$  et  $K_j$  ayant  $e$  comme arête commune. On considère ainsi la matrice élémentaire  $A_{\text{elem}}$  suivante :

$$A_{\text{elem}} = \begin{pmatrix} i & j \\ +\frac{|e_{ij}|}{|K_i|} (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^+ & +\frac{|e_{ij}|}{|K_i|} (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^- \\ +\frac{|e_{ij}|}{|K_j|} (\bar{\mathbf{V}} \cdot \mathbf{n}_{ji})^- & +\frac{|e_{ij}|}{|K_j|} (\bar{\mathbf{V}} \cdot \mathbf{n}_{ji})^+ \end{pmatrix} \begin{matrix} i \\ j \end{matrix}$$

*Remarque* : Puisque  $\mathbf{n}_{ji} = -\mathbf{n}_{ij}$ , on a les relations  $(\bar{\mathbf{V}} \cdot \mathbf{n}_{ji})^- = -(\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^+$  et  $(\bar{\mathbf{V}} \cdot \mathbf{n}_{ji})^+ = -(\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^-$ .



L'assemblage de la matrice  $A$  se fait alors de la façon suivante :

$$\begin{aligned} A(i, i) &\leftarrow A(i, i) + \frac{|e_{ij}|}{|K_i|} (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^+, & A(i, j) &\leftarrow \frac{|e_{ij}|}{|K_i|} (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^-, \\ A(j, i) &\leftarrow -\frac{|e_{ij}|}{|K_j|} (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^+, & A(j, j) &\leftarrow A(j, j) - \frac{|e_{ij}|}{|K_j|} (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^- \end{aligned}$$

### 3.5 Condition de stabilité

Comme dans la section précédente, on considère un maillage "cell center" et le système linéaire (3.20) correspondant s'écrit

$$\mathbf{u}^{n+1} = M\mathbf{u}^n$$

avec la matrice  $M = I_N - \Delta t A$  où  $A$  est donnée par (3.21).

On suppose dans cette section que  $\mathbf{V} \cdot \mathbf{n} = 0$  sur le bord  $\partial\Omega$ . La condition limite (3.2) du problème de transport est en particulier vérifiée, de sorte qu'il n'y a pas de conditions limites sur  $u$ . On suppose de plus que le champ de vitesse  $\mathbf{V}$  vérifie la condition d'incompressibilité  $\text{div } \mathbf{V} = 0$  dans  $\Omega$ .

Le schéma est *stable* (au sens  $L^\infty$ ) si  $\|M\|_\infty \leq 1$  avec  $\|M\|_\infty = \max_{1 \leq i \leq N} \sum_{j=1}^N |M_{ij}|$ .

On a

$$\begin{aligned} \|M\|_\infty &= \max_{1 \leq i \leq N} \sum_{j=1}^N |M_{ij}| \\ &= \max_{1 \leq i \leq N} \left( \left| 1 - \frac{\Delta t}{|K_i|} \sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i|K_j)}} |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^+ \right| + \frac{\Delta t}{|K_i|} \sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i|K_j)}} |e_{ij}| |(\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^-| \right) \end{aligned}$$

La condition de stabilité  $\|M\|_\infty \leq 1$  s'écrit, pour  $1 \leq i \leq N$ ,

$$\left| 1 - \frac{\Delta t}{|K_i|} \sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i|K_j)}} |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^+ \right| + \frac{\Delta t}{|K_i|} \sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i|K_j)}} |e_{ij}| |(\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^-| \leq 1 \quad (3.22)$$

Or, on a  $|(\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^-| = -(\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^- \geq 0$ . Par conséquent, (3.22) est équivalent à

$$\left| 1 - \frac{\Delta t}{|K_i|} \sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i|K_j)}} |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^+ \right| \leq 1 + \frac{\Delta t}{|K_i|} \sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i|K_j)}} |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^- \quad \text{pour } 1 \leq i \leq N \quad (3.23)$$

Par ailleurs, comme on a supposé une condition d'incompressibilité  $\text{div } \mathbf{V} = 0$  pour le champ de vitesse  $\mathbf{V}$ , on a

$$0 = \int_{K_i} \text{div } \mathbf{V} \, dx = \int_{\partial K_i} \mathbf{V} \cdot \mathbf{n} \, d\Gamma = \sum_{e_{ij} \subset \partial K_i} \int_{e_{ij}} \mathbf{V} \cdot \mathbf{n}_{ij} \, d\Gamma.$$

De plus, on a (cf. (3.14))  $\bar{\mathbf{V}} = \frac{1}{|e_{ij}|} \int_{e_{ij}} \mathbf{V} \, d\Gamma$  et donc  $\int_{e_{ij}} \mathbf{V} \cdot \mathbf{n}_{ij} \, d\Gamma = |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})$ . On obtient ainsi, pour  $1 \leq i \leq N$ ,

$$\boxed{\sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i|K_j)}} |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij}) = 0} \quad (3.24)$$

Par ailleurs, on a  $v = v^+ + v^-$  et  $|v| = v^+ - v^-$ , par conséquent la relation (3.24) s'écrit encore, pour  $1 \leq i \leq N$ ,

$$\sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i|K_j)}} |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^- = - \sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i|K_j)}} |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^+.$$

La condition (3.23) s'écrit donc, pour  $1 \leq i \leq N$ ,

$$\left| 1 - \frac{\Delta t}{|K_i|} \sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i|K_j)}} |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^+ \right| \leq 1 - \frac{\Delta t}{|K_i|} \sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i|K_j)}} |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^+,$$

ce qui est équivalent à la condition suivante :

$$\frac{\Delta t}{|K_i|} \sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i | K_j)}} |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^+ \leq 1, \quad \forall 1 \leq i \leq N.$$

Ainsi, le schéma "Volumes Finis" (3.12), (3.13) qui est explicite en temps, est stable si la condition de stabilité CFL suivante, est vérifiée :

$$\boxed{\Delta t \max_{1 \leq i \leq N} \left( \frac{1}{|K_i|} \sum_{\substack{e_{ij} \subset \partial K_i \\ e_{ij} = (K_i | K_j)}} |e_{ij}| (\bar{\mathbf{V}} \cdot \mathbf{n}_{ij})^+ \right) \leq 1.} \quad (3.25)$$

# Chapitre 4

## Equations de Stokes

### 4.1 Introduction

On va résoudre les équations de Stokes dans un domaine  $\Omega \subset \mathbb{R}^2$  par une méthode de Volumes Finis. On utilise des maillages décalés pour la vitesse et la pression ("staggered finite volume scheme"). On suppose que le domaine  $\Omega \subset \mathbb{R}^2$  est polygonal. On cherche la vitesse  $\mathbf{u} : \Omega \rightarrow \mathbb{R}^2$  et la pression  $p : \Omega \rightarrow \mathbb{R}$  telles que

$$\eta \mathbf{u} - \nu \Delta \mathbf{u} + \nabla p = \mathbf{f} \text{ dans } \Omega \quad (4.1)$$

$$\operatorname{div} \mathbf{u} = 0 \text{ dans } \Omega \quad (4.2)$$

$$\mathbf{u} = \mathbf{g} \text{ sur } \partial\Omega \quad (4.3)$$

avec les paramètres  $\eta > 0$ ,  $\nu > 0$  et les fonctions  $\mathbf{f}$  et  $\mathbf{g}$  données. Compte tenu de la condition d'incompressibilité (4.2), la fonction  $\mathbf{g}$  doit vérifier la condition de compatibilité suivante (par la formule de la divergence) :

$$\int_{\partial\Omega} \mathbf{g} \cdot \mathbf{n} \, ds = 0 \quad (4.4)$$

où  $\mathbf{n}$  est la normale unitaire extérieure au bord  $\partial\Omega$ .

### 4.2 Maillages

On choisit deux maillages différents du domaine  $\Omega$  pour la vitesse et pour la pression : un maillage de type Voronoï pour la vitesse et un maillage triangulaire pour la pression (cf. Section 1.4.3). En particulier les angles des triangles sont  $\leq \pi/2$ . On notera  $\mathcal{T}_h$  une triangulation du domaine polygonal  $\Omega$  vérifiant cette contrainte sur les angles des triangles. On associe à cette triangulation, les cellules de Voronoï  $K$  qui sont des polygones convexes formant une partition de  $\Omega$  (i.e.  $\bar{\Omega} = \cup_K \bar{K}$  et les cellules  $K$  sont deux à deux disjointes). A chaque volume  $K$ , on associe un centre  $\mathbf{x}_K$  (cf. Figure 1.9 et Figure 4.1) et on note  $d_{K,e}$  la distance de  $\mathbf{x}_K$  à l'arête  $e \subset \partial K$ . L'ensemble des arêtes d'un volume  $K$  est noté  $\mathcal{E}_K$ . On désigne par  $\mathbf{x}_T$  le centre d'un triangle  $T$  de la triangulation  $\mathcal{T}_h$ . Le centre  $\mathbf{x}_T$  est aussi un sommet de cellules de Voronoï (cf. Figure 4.1).

### 4.3 Formulation en Volumes Finis

Dans le système de Stokes (4.1)–(4.3), la pression est connue à une constante additive près. Pour résoudre le problème de Stokes, on introduit une pénalisation de l'équation de la divergence en remplaçant l'équation d'incompressibilité (4.2),  $\operatorname{div} \mathbf{u} = 0$ , par l'équation

$$\operatorname{div} \mathbf{u} = -\lambda h p \text{ dans } \Omega, \quad (4.5)$$

où  $\lambda > 0$  est un paramètre (à choisir) et  $h = \max_K(\operatorname{diam}(K))$ .

L'équation de la divergence pénalisée (4.5) avec la condition (4.4) sur la donnée au bord  $\mathbf{g}$  implique que la pression est à moyenne nulle sur  $\Omega$  i.e.  $\int_{\Omega} p \, d\mathbf{x} = 0$  (utiliser le théorème de la divergence). C'est équivalent à fixer la constante pour la pression. Le problème de Stokes pénalisé avec (4.5) est ainsi bien posé.

Dans le schéma Volumes Finis, on cherche la vitesse constante par cellule de Voronoï  $K$  et la pression constante par triangle  $T$ . Les approximations sont

$$\mathbf{u}_K \simeq \frac{1}{|K|} \int_K \mathbf{u} \, d\mathbf{x}, \quad p_T \simeq \frac{1}{|T|} \int_T p \, d\mathbf{x}.$$

### 4.3.1 Approximation de la divergence

On intègre la relation  $\operatorname{div} \mathbf{u} = -\lambda h p$  sur un triangle  $T \in \mathcal{T}_h$ . On obtient

$$\int_T \operatorname{div} \mathbf{u} \, d\mathbf{x} = -\lambda h \int_T p \, d\mathbf{x} \simeq -\lambda h |T| p_T \quad (4.6)$$

Par la formule de la divergence, on a :

$$\int_T \operatorname{div} \mathbf{u} \, d\mathbf{x} = \int_{\partial T} \mathbf{u} \cdot \mathbf{n}_T \, ds = \sum_{K \in \mathcal{M}_T} \int_{\partial T \cap K} \mathbf{u} \cdot \mathbf{n}_T \, ds$$

où  $\mathcal{M}_T$  désigne l'ensemble des volumes de contrôle  $K$  (cellules de Voronoï) ayant  $\mathbf{x}_T$  comme sommet. On approche

$$\int_{\partial T \cap K} \mathbf{u} \cdot \mathbf{n}_T \, ds \simeq \mathbf{u}_K \cdot \int_{\partial T \cap K} \mathbf{n}_T \, ds$$

En utilisant la relation  $\sum_e |e| \mathbf{n}_e = 0$  pour tout polygone, on obtient

$$\int_{\partial T \cap K} \mathbf{n}_T \, ds = -|e'| \mathbf{n}_{e'}$$

où  $e'$  est le segment  $[\mathbf{x}_1 \mathbf{x}_2]$  et  $\mathbf{x}_1$  et  $\mathbf{x}_2$  sont les milieux des arêtes du triangle  $T$  intersectant la cellule  $K$  (cf. Figure 4.1);  $\mathbf{n}_{e'}$  est la normale unitaire au segment  $e'$  dirigée vers l'extérieur de la cellule  $K$ .

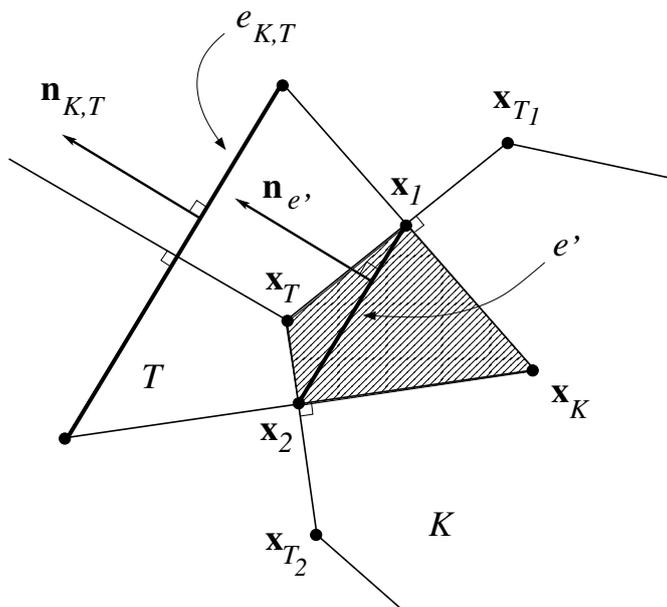


FIGURE 4.1 – Triangulation et cellule de Voronoï

On note  $e_{K,T}$  l'arête du triangle  $T$  opposé au sommet  $\mathbf{x}_K$  du triangle  $T$ , centre de la cellule  $K$ ;  $\mathbf{n}_{K,T}$  désigne la normale unitaire à l'arête  $e_{K,T}$ . On a  $|e_{K,T}| = 2|e'|$  et  $\mathbf{n}_{K,T} = \mathbf{n}_{e'}$ . La formulation en Volumes Finis de l'équation de la divergence pénalisée s'écrit

$$- \sum_{K \in \mathcal{M}_T} \mathbf{B}_{K,T} \cdot \mathbf{u}_K = -\lambda h|T|p_T \quad \text{pour tout } \mathbf{x}_T \notin \partial\Omega \quad (4.7)$$

avec

$$\boxed{\mathbf{B}_{K,T} = \frac{|e_{K,T}|}{2} \mathbf{n}_{K,T}} \quad (4.8)$$

### 4.3.2 Approximation de l'équation de Stokes

En intégrant l'équation de Stokes (4.1) sur un volume de contrôle  $K$ , on obtient :

$$\eta \int_K \mathbf{u} \, d\mathbf{x} - \nu \int_K \Delta \mathbf{u} \, d\mathbf{x} + \int_K \nabla p \, d\mathbf{x} = \int_K \mathbf{f} \, d\mathbf{x}. \quad (4.9)$$

En utilisant la formule de la divergence, on a

$$\int_K \Delta \mathbf{u} \, d\mathbf{x} = \int_{\partial K} \nabla \mathbf{u} \mathbf{n}_K \, ds = \sum_{e \in \mathcal{E}_K} \int_e \nabla \mathbf{u} \mathbf{n}_{K,e} \, ds,$$

et

$$\int_K \nabla p \, d\mathbf{x} = \int_{\partial K} p \mathbf{n}_K \, ds = \sum_{T \in \mathcal{T}_h} \int_{\partial K \cap T} p \mathbf{n}_K \, ds,$$

où  $\mathbf{n}_{K,e}$  (resp.  $\mathbf{n}_K$ ) désigne la normale unitaire à  $e$  (resp. à  $\partial K$ ) dirigée vers l'extérieur de la cellule  $K$ . L'équation (4.9) s'écrit donc

$$\eta \int_K \mathbf{u} \, d\mathbf{x} - \nu \sum_{e \in \mathcal{E}_K} \int_e \nabla \mathbf{u} \mathbf{n}_{K,e} \, ds + \sum_{T \in \mathcal{T}_h} \int_{\partial K \cap T} p \mathbf{n}_K \, ds = \int_K \mathbf{f} \, d\mathbf{x}.$$

• On approche

$$\boxed{- \int_e \nabla \mathbf{u} \mathbf{n}_{K,e} \, ds \simeq \mathbf{F}_{K,e}} \quad (4.10)$$

Pour l'arête  $e = (K|L)$  commune aux deux volumes de contrôle  $K$  et  $L$  avec  $e \not\subset \partial\Omega$ , le flux numérique  $\mathbf{F}_{K,e}$  est défini par (cf. Section 1.4, "Volumes Finis" pour le Laplacien) :

$$\mathbf{F}_{K,e} = - \frac{(\mathbf{u}_L - \mathbf{u}_K)}{|\mathbf{x}_K - \mathbf{x}_L|} |e| \quad (e \not\subset \partial\Omega). \quad (4.11)$$

Dans le cas où un volume  $K$  possède une arête  $e$  sur le bord  $\partial\Omega$ , on impose

$$\mathbf{u}_K = \mathbf{u}_e = \frac{1}{|e|} \int_e \mathbf{g}(\mathbf{x}) \, ds \quad \text{si } e \subset \partial\Omega. \quad (4.12)$$

• On approche

$$\int_{\partial K \cap T} p \mathbf{n}_K \, ds \simeq p_T \int_{\partial K \cap T} \mathbf{n}_K \, ds.$$

En utilisant la relation  $\sum_e |e| \mathbf{n}_e = 0$  pour tout polygone, on obtient

$$\int_{\partial K \cap T} \mathbf{n}_K \, ds = +|e'| \mathbf{n}_{e'}$$

où  $e'$  est le segment  $[\mathbf{x}_1\mathbf{x}_2]$  et  $\mathbf{x}_1, \mathbf{x}_2$  sont les milieux des arêtes du triangle  $T$  intersectant la cellule  $K$ ;  $\mathbf{n}_{e'}$  est la normale unitaire au segment  $e'$  dirigée vers l'extérieur de la cellule  $K$  (cf. Figure 4.1). Le terme de pression est ainsi approché par :

$$\boxed{\int_K \nabla p \, d\mathbf{x} = \sum_{T \in \mathcal{T}_h} \int_{\partial K \cap T} p \mathbf{n}_K \, ds \simeq \sum_{\mathbf{x}_T \in \mathcal{V}_K} \mathbf{B}_{K,T} p_T} \quad (4.13)$$

où  $\mathcal{V}_K$  est l'ensemble des sommets de  $K$  n'appartenant pas au bord  $\partial\Omega$  et le vecteur  $\mathbf{B}_{K,T}$  est donné par (4.8).

En combinant (4.9) avec (4.10) et (4.13), on obtient la formulation Volumes Finis de l'équation de Stokes :

$$\eta |K| \mathbf{u}_K + \nu \sum_{\substack{e \in \mathcal{E}_K \\ e \notin \partial\Omega}} \mathbf{F}_{K,e} + \sum_{\mathbf{x}_T \in \mathcal{V}_K} \mathbf{B}_{K,T} p_T = |K| \mathbf{f}_K,$$

où  $\mathbf{f}_K = \frac{1}{|K|} \int_K \mathbf{f} \, d\mathbf{x}$ .

★ En résumé, le schéma Volumes Finis pour les équations de Stokes s'écrit

$$\boxed{\begin{aligned} &\bullet \quad \eta |K| \mathbf{u}_K + \nu \sum_{\substack{e \in \mathcal{E}_K \\ e \notin \partial\Omega}} \mathbf{F}_{K,e} + \sum_{\mathbf{x}_T \in \mathcal{V}_K} \mathbf{B}_{K,T} p_T = |K| \mathbf{f}_K \quad \text{pour tout } K \text{ intérieure} \\ &\bullet \quad \sum_{K \in \mathcal{M}_T} \mathbf{B}_{K,T} \cdot \mathbf{u}_K = \lambda h |T| p_T \quad \text{pour tout } \mathbf{x}_T \notin \partial\Omega \\ &\bullet \quad \mathbf{u}_K = \frac{1}{|e|} \int_e \mathbf{g}(\mathbf{x}) \, ds \quad \text{si } e \subset \partial K \cap \partial\Omega \quad (\text{cellule } K \text{ ayant l'arête } e \text{ sur le bord}). \end{aligned}}$$

### 4.3.3 Système linéaire

On prend  $\eta = 0$  et  $\mathbf{g} = \mathbf{0}$  pour simplifier. Soit  $N$  le nombre de cellules de Voronoï *intérieures* (c'est-à-dire les cellules ne possédant pas d'arête sur le bord  $\partial\Omega$ ) et  $M$  le nombre de triangles de  $\mathcal{T}_h$ . On note  $\mathcal{U} = (u_{K_1}^1, u_{K_2}^1, \dots, u_{K_N}^1, u_{K_1}^2, u_{K_2}^2, \dots, u_{K_N}^2)^\top$  la vitesse avec  $\mathbf{u}_K = (u_K^1, u_K^2)$  et  $\mathcal{P} = (p_{T_1}, p_{T_2}, \dots, p_{T_M})^\top$  la pression. Les inconnues  $\mathcal{U}$  et  $\mathcal{P}$  vérifient le système linéaire

$$\begin{pmatrix} A & B \\ B^\top & -\lambda D \end{pmatrix} \begin{pmatrix} \mathcal{U} \\ \mathcal{P} \end{pmatrix} = \begin{pmatrix} \mathcal{F} \\ 0 \end{pmatrix} \quad (4.14)$$

La matrice  $A$  est de taille  $2N \times 2N$  avec  $A = \begin{pmatrix} \mathcal{A} & 0 \\ 0 & \mathcal{A} \end{pmatrix}$  et  $\mathcal{A}$  est la matrice du Laplacien de taille  $N \times N$  donnée par (1.45). De plus,  $D = \begin{pmatrix} |T_1| & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & |T_M| \end{pmatrix}$  est une matrice diagonale de taille  $M \times M$ .

On note  $\mathbf{B}_{K,T} = \begin{pmatrix} B_{K,T}^1 \\ B_{K,T}^2 \end{pmatrix} = \frac{|e_{K,T}|}{2} \begin{pmatrix} n_{K,T}^1 \\ n_{K,T}^2 \end{pmatrix}$ . La matrice correspondante  $B$  est de taille  $2N \times M$  avec  $B = \begin{pmatrix} B_1 \\ B_2 \end{pmatrix}$  et  $B_l = (B_{K_i, T_j}^l)_{\substack{1 \leq i \leq N \\ 1 \leq j \leq M}}$  ( $l = 1, 2$ ) sont des matrices de taille  $N \times M$ .

## Chapitre 5

# Equations de Navier-Stokes incompressibles

### 5.1 Introduction

On s'intéresse à présent aux équations de Navier-Stokes incompressibles. Il s'agit d'un problème d'évolution en temps obtenu à partir des équations de Stokes avec un terme convectif non-linéaire. Le terme convectif nonlinéaire sera linéarisé et traité de façon semi-implicite par un schéma *upwind*, de façon analogue à ce qui a été fait pour l'équation de transport (cf. Chapitre 3).

On cherche la vitesse  $\mathbf{u} = \mathbf{u}(\mathbf{x}, t)$  définie de  $\Omega \times (0, T)$  dans  $\mathbb{R}^2$  et la pression  $p = p(\mathbf{x}, t)$  définie de  $\Omega \times (0, T)$  dans  $\mathbb{R}$  telles que

$$(P_{\text{NS}}) \begin{cases} \mathbf{u}_t + (\mathbf{u} \cdot \nabla) \mathbf{u} - \nu \Delta \mathbf{u} + \nabla p = \mathbf{f} & \text{dans } \Omega \times (0, T) \\ \operatorname{div} \mathbf{u} = 0 & \text{dans } \Omega \times (0, T) \\ \mathbf{u} = \mathbf{g} & \text{sur } \partial\Omega \times (0, T) \\ \mathbf{u}(0) = \mathbf{u}_0 & \text{dans } \Omega \end{cases}$$

### 5.2 Semi-discrétisation en temps

On note  $t^n = n\Delta t$  et on considère les approximations en temps de la vitesse  $\mathbf{u}^n \simeq \mathbf{u}(\cdot, t^n)$  et de la pression  $p^n \simeq p(\cdot, t^n)$ . Le schéma semi-discrétisé en temps s'écrit

$$\begin{aligned} \frac{\mathbf{u}^n - \mathbf{u}^{n-1}}{\Delta t} + (\mathbf{u}^{n-1} \cdot \nabla) \mathbf{u}^n - \nu \Delta \mathbf{u}^n + \nabla p^n &= \mathbf{f}^n & \text{dans } \Omega \\ \operatorname{div} \mathbf{u}^n &= 0 & \text{dans } \Omega \\ \mathbf{u}^n &= \mathbf{g} & \text{sur } \partial\Omega \end{aligned}$$

Compte tenu du fait que le champ de vitesse  $\mathbf{u}^{n-1}$  est à divergence nulle, on écrit le terme de convection (linéarisé) de la façon suivante :

$$(\mathbf{u}^{n-1} \cdot \nabla) \mathbf{u}^n = \operatorname{div}(\mathbf{u}^{n-1} \otimes \mathbf{u}^n)$$

où  $(\mathbf{u} \otimes \mathbf{v})_{ij} = u_i v_j$  et  $(\operatorname{div}(\mathbf{u} \otimes \mathbf{v}))_j = \sum_i \frac{\partial(u_i v_j)}{\partial x_i}$  avec  $\mathbf{u} = (u_1, u_2)$  et  $\mathbf{v} = (v_1, v_2)$ .

Le problème semi-discrétisé s'écrit alors

$$\begin{aligned} \frac{1}{\Delta t} \mathbf{u}^n + \operatorname{div}(\mathbf{u}^{n-1} \otimes \mathbf{u}^n) - \nu \Delta \mathbf{u}^n + \nabla p^n &= \mathbf{f}^n + \frac{1}{\Delta t} \mathbf{u}^{n-1} & \text{dans } \Omega \\ \operatorname{div} \mathbf{u}^n &= 0 & \text{dans } \Omega \\ \mathbf{u}^n &= \mathbf{g} & \text{sur } \partial\Omega \end{aligned}$$

Pour déterminer  $\mathbf{u}^n$  et  $p^n$  à partir de  $\mathbf{u}^{n-1}$ , il s'agit donc de résoudre un problème de Stokes avec un terme de convection (linéaire).

### 5.3 Formulations en Volumes Finis

Le terme convectif nonlinéaire des équations de Navier-Stokes sera linéarisé et traité de façon semi-implicite par un schéma *upwind*. On décrit ici uniquement le traitement du terme de convection. Pour plus de détails, on renvoie au Chapitre 3 "Volumes Finis pour l'équation de transport". On intègre le terme de convection sur un volume de contrôle  $K$ . Par la formule de la divergence, on obtient

$$\int_K \operatorname{div}(\mathbf{u}^{n-1} \otimes \mathbf{u}^n) d\mathbf{x} = \int_{\partial K} (\mathbf{u}^{n-1} \cdot \mathbf{n}) \mathbf{u}^n ds = \sum_{\substack{e_j \subset \partial K \\ e_j = (K|K_j)}} \int_{e_j} (\mathbf{u}^{n-1} \cdot \mathbf{n}_{e_j}) \mathbf{u}^n ds$$

où  $e_j$  désigne l'arête commune aux volumes  $K$  et  $K_j$  et  $\mathbf{n}_{e_j}$  est la normale unitaire à  $e_j$  dirigée vers l'extérieur de  $K$  (cf. Figure 5.1).

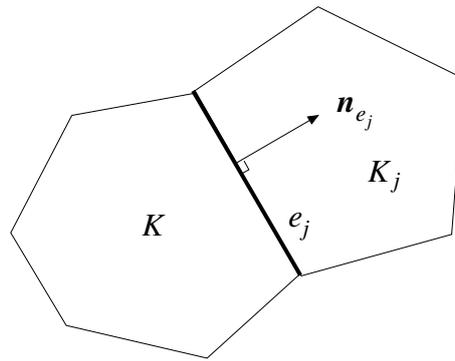


FIGURE 5.1 – Volumes de contrôle pour le traitement de la convection

On approche alors l'intégrale en introduisant un flux numérique  $\Phi$  :

$$\int_K \operatorname{div}(\mathbf{u}^{n-1} \otimes \mathbf{u}^n) d\mathbf{x} \simeq \sum_{\substack{e_j \subset \partial K \\ e_j = (K|K_j)}} |e_j| \Phi(\mathbf{u}_K^n, \mathbf{u}_{K_j}^n, \mathbf{n}_{e_j}).$$

Le flux numérique est choisi par décentrement *upwind* :

$$\Phi(\mathbf{u}_K^n, \mathbf{u}_{K_j}^n, \mathbf{n}_{e_j}) = \mathbf{u}_K^n (\mathbf{u}_K^{n-1} \cdot \mathbf{n}_{e_j})^+ + \mathbf{u}_{K_j}^n (\mathbf{u}_K^{n-1} \cdot \mathbf{n}_{e_j})^-$$

où  $(\cdot)^+$  et  $(\cdot)^-$  désignent respectivement les parties positives et négatives.

# Appendices



## Annexe A

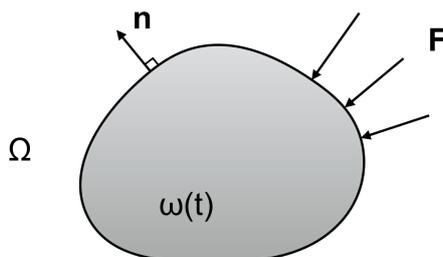
# Modélisation des équations de Navier-Stokes et équations de Stokes

### A.1 Introduction

Un fluide visqueux incompressible est contenu dans un domaine  $\Omega \subset \mathbb{R}^d$  ( $d \leq 3$ ). On introduit la densité du fluide  $\rho(\mathbf{x}, t) \in \mathbb{R}$ , la vitesse  $\mathbf{u}(\mathbf{x}, t) \in \mathbb{R}^d$  et la pression  $p(\mathbf{x}, t) \in \mathbb{R}$  en un point  $\mathbf{x} \in \Omega$  du fluide à l'instant  $t$ . Le fluide est caractérisé par un tenseur des contraintes  $\sigma(\mathbf{u}, p)$ . La résultante des forces exercées par le fluide sur un sous-domaine  $\omega(t) \subset \Omega$  est alors donnée par

$$\mathbf{F} = \int_{\partial\omega(t)} \sigma(\mathbf{u}, p) \mathbf{n} d\Gamma \quad (\text{A.1})$$

où  $\mathbf{n}$  désigne la normale unitaire extérieure à  $\omega(t)$ .



La force  $\mathbf{F}$  est une force surfacique s'appliquant sur le bord  $\partial\omega$ . Le fluide est également soumis à une force volumique extérieure  $\mathbf{f} = \mathbf{f}(\mathbf{x}, t) \in \mathbb{R}^d$  s'appliquant en chaque point  $\mathbf{x}$  du domaine  $\Omega$ , à l'instant  $t$ .

### A.2 Conservation de la masse

La masse du fluide contenu dans le domaine  $\omega(t)$  est conservée à chaque instant  $t$ . On a donc

$$\frac{d}{dt} \int_{\omega(t)} \rho d\mathbf{x} = 0. \quad (\text{A.2})$$

En utilisant la formule de transport de Reynolds<sup>1</sup>, on obtient

$$0 = \int_{\omega(t)} \frac{\partial \rho}{\partial t} d\mathbf{x} + \int_{\partial\omega(t)} \rho \mathbf{u} \cdot \mathbf{n} d\Gamma$$

et par le théorème de la divergence, il vient

$$\int_{\omega(t)} \frac{\partial \rho}{\partial t} + \operatorname{div}(\rho \mathbf{u}) d\mathbf{x} = 0 \quad (\text{A.3})$$

La relation (A.3) étant vraie pour tout sous-domaine  $\omega \subset \Omega$ , on en déduit *l'équation de continuité*

$$\frac{\partial \rho}{\partial t} + \operatorname{div}(\rho \mathbf{u}) = 0 \quad \text{dans } \Omega \times \mathbb{R}^+. \quad (\text{A.4})$$

### A.3 Loi fondamentale de la dynamique - loi de comportement

La variation de la quantité de mouvement dans le domaine  $\omega(t)$  est égale à chaque instant à la somme des forces extérieures s'appliquant sur  $\omega(t)$ . Ainsi, on a

$$\frac{d}{dt} \int_{\omega(t)} \rho \mathbf{u} d\mathbf{x} = \mathbf{F} + \int_{\omega(t)} \rho \mathbf{f} d\mathbf{x}. \quad (\text{A.5})$$

Par le théorème de la divergence, on obtient

$$\mathbf{F} = \int_{\partial\omega(t)} \sigma(\mathbf{u}, p) \mathbf{n} d\Gamma = \int_{\omega(t)} \operatorname{div} \sigma(\mathbf{u}, p) d\mathbf{x}. \quad (\text{A.6})$$

D'autre part, la formule de transport de Reynolds donne

$$\frac{d}{dt} \int_{\omega(t)} \rho \mathbf{u} d\mathbf{x} = \int_{\omega(t)} \frac{\partial}{\partial t} (\rho \mathbf{u}) d\mathbf{x} + \int_{\partial\omega(t)} (\rho \mathbf{u}) \mathbf{u} \cdot \mathbf{n} d\Gamma$$

et par application du théorème de la divergence, on obtient

$$\begin{aligned} \frac{d}{dt} \int_{\omega(t)} \rho \mathbf{u} d\mathbf{x} &= \int_{\omega(t)} \left( \rho \frac{\partial \mathbf{u}}{\partial t} + \mathbf{u} \frac{\partial \rho}{\partial t} \right) d\mathbf{x} + \int_{\omega(t)} (\operatorname{div}(\rho \mathbf{u}) \mathbf{u} + \rho (\mathbf{u} \cdot \nabla) \mathbf{u}) d\mathbf{x} \\ &= \int_{\omega(t)} \left( \rho \left( \frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} \right) + \underbrace{\mathbf{u} \left( \frac{\partial \rho}{\partial t} + \operatorname{div}(\rho \mathbf{u}) \right)}_{=0 \text{ d'après (A.4)}} \right) d\mathbf{x} \end{aligned}$$

et par conséquent

$$\frac{d}{dt} \int_{\omega(t)} \rho \mathbf{u} d\mathbf{x} = \int_{\omega(t)} \rho \left( \frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} \right) d\mathbf{x} \quad (\text{A.7})$$

En combinant (A.5),(A.6) et (A.7), on obtient

$$\int_{\omega(t)} \rho \left( \frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} \right) d\mathbf{x} = \int_{\omega(t)} (\operatorname{div} \sigma(\mathbf{u}, p) + \rho \mathbf{f}) d\mathbf{x}.$$

1. *Formule de transport de Reynolds.* Soit  $\Omega_0 \subset \mathbb{R}^d$  un domaine borné régulier. On considère  $\mathbf{X} : \Omega_0 \times \mathbb{R}^+ \rightarrow \mathbb{R}^d$  régulière telle que  $\mathbf{X}(\cdot, t)$  est un difféomorphisme sur  $\Omega_0$  avec  $\det \nabla \mathbf{X}(\cdot, t) > 0$  pour tout  $t \geq 0$ . On définit alors  $\Omega(t) = \mathbf{X}(\Omega_0, t)$  et pour  $\mathbf{x} \in \Omega(t)$ , on note  $\mathbf{V}(\mathbf{x}, t) = \frac{\partial \mathbf{X}}{\partial t}(\mathbf{X}^{-1}(\mathbf{x}, t), t)$ ,  $t \geq 0$ . Alors pour toute fonction  $F : \mathbb{R}^d \times \mathbb{R}^+ \rightarrow \mathbb{R}$  régulière, on a, pour tout  $t \geq 0$

$$\frac{d}{dt} \int_{\Omega(t)} F(\mathbf{x}, t) d\mathbf{x} = \int_{\Omega(t)} \frac{\partial}{\partial t} F(\mathbf{x}, t) d\mathbf{x} + \int_{\partial\Omega(t)} F(\mathbf{x}, t) \mathbf{V} \cdot \mathbf{n} d\Gamma$$

où  $\mathbf{n}$  désigne la normale unitaire dirigée vers l'extérieur de  $\Omega(t)$ .

La relation précédente étant vraie pour tout sous-domaine  $\omega \subset \Omega$ , on en déduit l'équation

$$\frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} - \operatorname{div}(\sigma(\mathbf{u}, p)) = \rho \mathbf{f} \quad \text{dans } \Omega \times \mathbb{R}^+ \quad (\text{A.8})$$

**Loi de comportement.** Pour un fluide dit *newtonien*, le tenseur des contraintes  $\sigma(\mathbf{u}, p)$  est une fonction linéaire du tenseur symétrique des déformations

$$D(\mathbf{u}) = \frac{1}{2} (\nabla \mathbf{u} + \nabla \mathbf{u}^\top).$$

Plus précisément, on a

$$\sigma(\mathbf{u}, p) = 2\nu D(\mathbf{u}) - pI_d, \quad (\text{A.9})$$

où  $\nu > 0$  est la viscosité du fluide. Dans ces conditions, on obtient  $\operatorname{div}(\sigma(\mathbf{u}, p)) = \nu \Delta \mathbf{u} - \nabla p$  et l'équation (A.8) devient

$$\frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} - \nu \Delta \mathbf{u} + \nabla p = \rho \mathbf{f} \quad \text{dans } \Omega \times \mathbb{R}^+ \quad (\text{A.10})$$

## A.4 Conservation du volume

Le fluide est supposé incompressible c'est-à-dire qu'il y a conservation du volume du domaine  $\omega(t)$  :

$$\frac{d}{dt} \int_{\omega(t)} d\mathbf{x} = 0.$$

En appliquant la formule de transport de Reynolds combinée à la formule de la divergence, on obtient

$$0 = \int_{\partial\omega(t)} \mathbf{u} \cdot \mathbf{n} d\Gamma = \int_{\omega(t)} \operatorname{div} \mathbf{u} d\mathbf{x}.$$

La relation précédente étant vraie pour tout sous-domaine  $\omega \subset \Omega$ , on en déduit la *condition d'incompressibilité*

$$\operatorname{div} \mathbf{u} = 0 \quad \text{dans } \Omega \times \mathbb{R}^+. \quad (\text{A.11})$$

En résumé, les équations de Navier-Stokes incompressibles s'écrivent

$$\frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} - \nu \Delta \mathbf{u} + \nabla p = \rho \mathbf{f} \quad \text{dans } \Omega \times \mathbb{R}^+ \quad (\text{A.12})$$

$$\frac{\partial \rho}{\partial t} + \operatorname{div}(\rho \mathbf{u}) = 0 \quad \text{dans } \Omega \times \mathbb{R}^+ \quad (\text{A.13})$$

$$\operatorname{div} \mathbf{u} = 0 \quad \text{dans } \Omega \times \mathbb{R}^+. \quad (\text{A.14})$$

## A.5 Adimensionalisation des équations de Navier-Stokes

Afin de mettre en évidence le rapport des forces de viscosité  $\nu \Delta \mathbf{u}$  sur les forces d'inertie  $(\mathbf{u} \cdot \nabla) \mathbf{u}$ , il est nécessaire d'écrire les équations de Navier-Stokes sous forme adimensionnée. La forme adimensionnée permet d'obtenir les équations de Stokes et les équations d'Euler comme cas *limites* (formelles) des équations de Navier-Stokes. On introduit une vitesse caractéristique  $U \in \mathbb{R}$  de l'écoulement étudié (par exemple,  $U$  peut être liée à une condition limite non-homogène, ou à une vitesse maximale dans le domaine ...) ainsi qu'une longueur caractéristique  $L$  (par exemple le diamètre de  $\Omega$ ). On considère alors le temps caractéristique  $T = L/U$  et on pose

$$\tilde{x} = \frac{x}{L}, \quad \tilde{t} = \frac{t}{T}, \quad \tilde{\mathbf{u}}(\tilde{\mathbf{x}}, \tilde{t}) = \frac{\mathbf{u}(\mathbf{x}, t)}{U}, \quad \tilde{p}(\tilde{\mathbf{x}}, \tilde{t}) = \frac{p(\mathbf{x}, t)}{\rho U^2}. \quad (\text{A.15})$$

Les nouvelles vitesse et pression  $\tilde{\mathbf{u}}$  et  $\tilde{p}$  vérifient alors

$$\rho \left( \frac{U}{L} \tilde{\mathbf{u}}_t + \frac{U^2}{L} (\tilde{\mathbf{u}} \cdot \nabla) \tilde{\mathbf{u}} \right) - \nu \frac{U}{L^2} \Delta \tilde{\mathbf{u}} + \frac{\rho U^2}{L} \nabla \tilde{p} = \mathbf{f} \quad \text{dans } \tilde{\Omega} \times \mathbb{R}^+.$$

Les (nouveaux) opérateurs différentiels  $\nabla$  et  $\Delta$  ci-dessus sont relatifs à la (nouvelle) variable  $\tilde{\mathbf{x}}$ . On obtient ainsi

$$\tilde{\mathbf{u}}_t + (\tilde{\mathbf{u}} \cdot \nabla) \tilde{\mathbf{u}} - \frac{1}{Re} \Delta \tilde{\mathbf{u}} + \nabla \tilde{p} = \tilde{\mathbf{f}} \quad \text{dans } \tilde{\Omega} \times \mathbb{R}^+, \quad (\text{A.16})$$

$$\operatorname{div} \tilde{\mathbf{u}} = 0 \quad \text{dans } \tilde{\Omega} \times \mathbb{R}^+, \quad (\text{A.17})$$

avec  $\tilde{\mathbf{f}} = \frac{L}{\rho U^2} \mathbf{f}$  et  $Re$  est le nombre de Reynolds défini par

$$Re = \frac{LU}{\nu} \rho. \quad (\text{A.18})$$

Le nombre  $\tilde{\nu} = \nu/\rho$  représente la viscosité cinématique. Par exemple, on a

$$\tilde{\nu} = 0.15 \cdot 10^{-4} \text{ m/s pour l'air}$$

$$\tilde{\nu} = 10^{-6} \text{ m/s pour l'eau.}$$

Le tableau suivant indique quelques valeurs du nombre de Reynolds.

	$U$	$L$	$Re = LU/\tilde{\nu}$
bactérie (dans l'eau)	100 $\mu\text{m/s}$	0.1 $\mu\text{m}$	$10^{-5}$
protozoaire	$10^{-1}$ $\text{cm/s}$	$10^{-2}$ $\text{cm}$	$10^{-1}$
guêpe	2 $\text{cm/s}$	2 $\text{cm}$	26
papillon	1 $\text{m/s}$	5 $\text{cm}$	3333
pigeon	5 $\text{m/s}$	30 $\text{cm}$	$10^5$
poisson (hareng)	1.67 $\text{m/s}$	30 $\text{cm}$	$5 \cdot 10^5$
poisson (saumon)	12.5 $\text{m/s}$	1 $\text{m}$	$1.25 \cdot 10^7$
automobile	100 $\text{km/h}$	3 $\text{m}$	$5 \cdot 10^6$
avion (airbus A330)	860 $\text{km/h}$	60 $\text{m}$	$\simeq 10^9$

Le nombre de Reynolds caractérise le type d'écoulement étudié. Plus le nombre de Reynolds est petit, plus les forces de viscosité sont importantes et les effets inertiels négligeables. À l'inverse, plus le nombre de Reynolds est grand, plus les forces d'inertie sont importantes.

## A.6 Réductions des équations

À partir des équations de Navier-Stokes, on obtient les équations de Stokes et les équations d'Euler selon que le nombre de Reynolds  $Re$  est petit ou grand.

★ Pour  $Re \ll 1$ , les effets dus à la viscosité sont dominants. Si on pose  $p' = LU\rho\tilde{p} = \nu Re\tilde{p}$  et  $\mathbf{f}' = \nu Re\tilde{\mathbf{f}}$ , l'équation (A.16) devient

$$\mathbf{u}_t + (\mathbf{u} \cdot \nabla) \mathbf{u} - \frac{1}{Re} \Delta \mathbf{u} + \frac{1}{\nu Re} \nabla p' = \frac{1}{\nu Re} \mathbf{f}'.$$

En faisant tendre le nombre  $Re$  vers 0, on obtient alors les équations de Stokes stationnaires c'est-à-dire que la vitesse  $\mathbf{u}$  et la pression  $p$  ne dépendent plus du temps  $t$ . En oubliant les *primes*, on obtient :

$$-\nu \Delta \mathbf{u} + \nabla p = \mathbf{f} \quad \text{dans } \Omega, \quad (\text{A.19})$$

$$\operatorname{div} \mathbf{u} = 0 \quad \text{dans } \Omega. \quad (\text{A.20})$$

★ Pour  $Re \gg 1$ , le terme de convection nonlinéaire  $(\mathbf{u} \cdot \nabla)\mathbf{u}$  est dominant ; dans ce cas, en faisant tendre  $Re$  vers  $+\infty$  dans l'équation (A.16) (ou bien en prenant directement  $\nu = 0$  dans (A.12)), on obtient les équations d'Euler (sans les *tildes*) :

$$\mathbf{u}_t + (\mathbf{u} \cdot \nabla)\mathbf{u} + \nabla p = \mathbf{f} \quad \text{dans } \Omega \times \mathbb{R}^+, \quad (\text{A.21})$$

$$\operatorname{div} \mathbf{u} = 0 \quad \text{dans } \Omega \times \mathbb{R}^+. \quad (\text{A.22})$$



## Annexe B

# Eléments d'analyse matricielle.

Pour les rappels d'analyse matricielle, on pourra consulter [9],[13],[3],[8].

### B.1 Matrice réductible

**Définition B.1** Une matrice  $A \in \mathbb{R}^{N \times N}$  est **réductible** s'il existe une partition de  $\{1, \dots, N\}$  en deux sous-ensembles non-vides  $I = \{i_1, \dots, i_p\}$ ,  $J = \{j_1, \dots, j_q\}$  tels que  $A_{i,j} = 0$  pour tout  $(i, j) \in I \times J$ . Une matrice non réductible est dite **irréductible**.

**Exemples.** La matrice  $A_1 = \begin{pmatrix} 12 & -1 & 3 \\ 0 & 2 & 0 \\ 3 & 1 & 9 \end{pmatrix}$  est réductible avec la partition  $I = \{2\}$ ,  $J = \{1, 3\}$ . La

matrice  $A_2 = \begin{pmatrix} 3 & 1 & 3 & -2 \\ 0 & -1 & 7 & 0 \\ 0 & 3 & 2 & 0 \\ 2 & -1 & 2 & 5 \end{pmatrix}$  est réductible avec  $I = \{2, 3\}$ ,  $J = \{1, 4\}$ .

Une caractérisation des matrices réductibles est donnée par le résultat suivant.

**Proposition B.1** Une matrice  $A \in \mathbb{R}^{N \times N}$  est réductible si et seulement s'il existe une matrice de permutation  $P$  telle que  $P^\top A P = \begin{pmatrix} B & C \\ 0 & D \end{pmatrix}$  où  $B$  et  $D$  sont des matrices carrées.

Dans l'exemple précédent, on obtient  $P^\top A_1 P = \begin{pmatrix} 12 & 3 & -1 \\ 3 & 9 & 1 \\ 0 & 0 & 2 \end{pmatrix}$  avec  $P = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$ . De même,

$P^\top A_2 P = \begin{pmatrix} 3 & -2 & 3 & 1 \\ 2 & 5 & 2 & -1 \\ 0 & 0 & 2 & 3 \\ 0 & 0 & 7 & -1 \end{pmatrix}$  avec  $P = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}$ .

**Proposition B.2** Soit  $A \in \mathbb{R}^{N \times N}$  la matrice tridiagonale donnée par

$$A = \begin{pmatrix} a_1 & b_1 & & & 0 \\ c_1 & a_2 & b_2 & & \\ & \ddots & \ddots & \ddots & \\ & & c_{N-2} & a_{N-1} & b_{N-1} \\ 0 & & & c_{N-1} & a_N \end{pmatrix}$$

On suppose que  $a_i, b_i, c_i \neq 0$  pour tout  $1 \leq i \leq N$ . Alors  $A$  est irréductible.

*Démonstration.* On suppose que  $A$  est réductible. Il existe donc deux sous-ensembles non vides  $I$  et  $J$  qui forment une partition de  $\{1, \dots, N\}$  tels que  $A_{i,j} = 0, \forall (i,j) \in I \times J$ . Soit  $(i_1, j_1) \in I \times J$ . On a  $A_{i_1, j_1} = 0$ .

— Si  $i_1 \neq 1$  et  $i_1 \neq N$ , alors  $\{i_1 - 1, i_1 + 1\} \not\subset J$  car  $A_{i_1, i_1-1} = c_{i_1-1} \neq 0$  et  $A_{i_1, i_1+1} = b_{i_1} \neq 0$ .

On note  $i_2 \in \{i_1 - 1, i_1 + 1\}$ .

— Si  $i_1 = 1$  alors  $i_2 = 2 \notin J$  car  $A_{1,2} = b_1 \neq 0$ .

— Si  $i_1 = N$  alors  $i_2 = N - 1 \notin J$  car  $A_{N-1, N} = c_{N-1} \neq 0$ .

On recommence avec  $i_2$  et ensuite de suite ... On obtient nécessairement que  $J = \emptyset$ , d'où la contradiction. La matrice  $A$  est donc irréductible.  $\square$

## B.2 Matrice à diagonale dominante

**Définition B.2** Soit  $A = (a_{ij})_{1 \leq i, j \leq N}$  une matrice carrée de taille  $N \times N$  à coefficients réels ou complexes.

1. La matrice  $A$  est à **diagonale dominante** si

$$|a_{ii}| \geq \sum_{\substack{1 \leq j \leq N \\ j \neq i}} |a_{ij}|, \text{ pour tout } 1 \leq i \leq N.$$

2. La matrice  $A$  est à **diagonale strictement dominante** si

$$|a_{ii}| > \sum_{\substack{1 \leq j \leq N \\ j \neq i}} |a_{ij}|, \text{ pour tout } 1 \leq i \leq N.$$

3. La matrice  $A$  est à **diagonale fortement dominante** si

(a)  $|a_{ii}| \geq \sum_{\substack{1 \leq j \leq N \\ j \neq i}} |a_{ij}|, \text{ pour tout } 1 \leq i \leq N$

(b) il existe  $i_0$  tel que  $|a_{i_0 i_0}| > \sum_{\substack{1 \leq j \leq N \\ j \neq i_0}} |a_{i_0 j}|.$

**Proposition B.3** Soit  $A$  une matrice carrée.

1. Si  $A$  est à diagonale strictement dominante alors  $A$  est inversible.

2. Si  $A$  est une matrice à diagonale fortement dominante et irréductible alors  $A$  est inversible.

## B.3 Matrice monotone

**Définition B.3** Une matrice  $A \in \mathbb{R}^{N \times N}$  est dite **monotone** si  $A$  est inversible et si  $A^{-1} \geq 0$  i.e. tous les coefficients de  $A^{-1}$  sont positifs ou nuls.

La terminologie est en fait justifiée par le résultat suivant.

**Proposition B.4** Une matrice  $A$  est monotone si et seulement si  $(A\mathbf{x} \geq 0 \Rightarrow \mathbf{x} \geq 0)$ .

*Démonstration.*

- Supposons tout d'abord que  $A$  soit une matrice monotone et considérons  $\mathbf{x} \in \mathbb{R}^N$  tel que  $A\mathbf{x} \geq 0$ . Montrons que  $\mathbf{x} \geq 0$ . La matrice  $A$  étant inversible, on écrit  $\mathbf{x} = A^{-1}A\mathbf{x}$ . Comme  $A^{-1} \geq 0$  et que  $A\mathbf{x} \geq 0$ , on en déduit que  $\mathbf{x} \geq 0$ .

- Supposons maintenant que la propriété  $(A\mathbf{x} \geq 0 \Rightarrow \mathbf{x} \geq 0)$  soit vraie. Montrons d'abord que  $A$  est inversible. Soit  $\mathbf{x} \in \mathbb{R}^N$  tel que  $A\mathbf{x} = 0$ . On déduit de l'hypothèse que  $\mathbf{x} \geq 0$ . De plus, on a  $A(-\mathbf{x}) = 0$  d'où l'on déduit également que  $-\mathbf{x} \geq 0$ . On a ainsi établi que  $\mathbf{x} = 0$ . Il reste à montrer que  $A^{-1} \geq 0$ . Soit alors  $\mathbf{y} \in \mathbb{R}^N$  tel que  $\mathbf{y} \geq 0$  et on pose  $\mathbf{x} = A^{-1}\mathbf{y}$  et par conséquent  $\mathbf{y} = A\mathbf{x} \geq 0$ . L'hypothèse implique que  $\mathbf{x} \geq 0$  et donc que  $A^{-1}\mathbf{y} \geq 0$ . Par des choix convenables de  $\mathbf{y}$  ( $\mathbf{y} = (1, 0, \dots, 0)^\top, \mathbf{y} = (0, 1, 0, \dots, 0)^\top, \dots, \mathbf{y} = (0, \dots, 0, 1)^\top$ ), on en déduit que  $A^{-1} \geq 0$ .  $\square$

**Définition B.4** Une matrice  $A \in \mathbb{R}^{N \times N}$  est une **M-matrice** si elle est monotone et si  $a_{ij} \leq 0$  pour  $i \neq j$ .

Voici quelques critères pratiques de détermination de M-matrices.

**Proposition B.5** Soit une matrice  $A \in \mathbb{R}^{N \times N}$  telle que

i)  $a_{ij} \leq 0$  pour tous  $i \neq j$ .

ii)  $\sum_{1 \leq j \leq N} a_{ij} > 0$  pour  $i = 1, \dots, N$ .

Alors  $A$  est une M-matrice.

*Démonstration.*

On a clairement  $a_{ii} > 0$ , pour tout  $i$ . On introduit alors la matrice diagonale  $D = \text{diag}(a_{ii})$  qui est inversible avec  $D^{-1} = \text{diag}(1/a_{ii})$  et par conséquent  $D^{-1} \geq 0$ . On décompose alors  $A$  sous la forme  $A = D(I - M)$  où  $I \in \mathbb{R}^{N \times N}$  est la matrice identité et  $M = (m_{ij})$  avec  $m_{ij} = -\frac{a_{ij}}{a_{ii}}$  pour  $i \neq j$  et  $m_{ii} = 0$ . Par l'hypothèse i), on a  $M \geq 0$ . On va montrer que  $I - M$  est inversible puis monotone. Pour cela, on va utiliser un résultat bien connu sur le rayon spectral de matrice (cf. [3]).

**Lemme B.1** Soit  $M \in \mathbb{R}^{N \times N}$  une matrice dont on note  $\lambda_i(M)$ ,  $i = 1, \dots, N$  les valeurs propres et  $\rho(M) = \max_i |\lambda_i(M)|$  le rayon spectral<sup>(1)</sup>.

i) Pour toute norme matricielle  $\|\cdot\|$  subordonnée<sup>(2)</sup>, on a  $\rho(M) \leq \|M\|$ .

ii) Si  $\rho(M) < 1$  alors  $I - M$  est inversible et  $(I - M)^{-1} = \sum_{k=0}^{\infty} M^k$ .

En utilisant le point i) de ce lemme, on a  $\rho(M) \leq \|M\|_{\infty} = \sup_{1 \leq i \leq N} \sum_{j=1}^N |m_{ij}|$ . Or, d'après les hypothèses i) et ii) de la Proposition, il vient

$$\sum_{j=1}^N |m_{ij}| = \sum_{j \neq i} \left| -\frac{a_{ij}}{a_{ii}} \right| = -\frac{1}{a_{ii}} \sum_{j \neq i} a_{ij} < 1,$$

ce qui implique d'après le point ii) du lemme, que la matrice  $I - M$  est inversible et  $(I - M)^{-1} = \sum_{k=0}^{\infty} M^k$ .

Comme  $M \geq 0$ , on a aussi  $M^k \geq 0$  et donc  $(I - M)^{-1} \geq 0$ . Ainsi,  $A = D(I - M)$  est inversible et  $A^{-1} = (I - M)^{-1} D^{-1} \geq 0$ . Les coefficients  $a_{ij}$  étant tous négatifs ou nuls par hypothèse, la matrice  $A$  est bien une M-matrice.  $\square$

On peut "relacher" la contrainte de positivité stricte de la somme des coefficients et obtenir le résultat suivant.

1. Si  $M$  est symétrique, on a  $\|M\|_2 := \sup_{\mathbf{x} \neq 0} \frac{\|M\mathbf{x}\|_2}{\|\mathbf{x}\|_2} = \rho(M)$ .

2.  $\|M\| = \sup_{\mathbf{x} \neq 0} \frac{\|M\mathbf{x}\|}{\|\mathbf{x}\|}$

**Proposition B.6** Soit une matrice  $A \in \mathbb{R}^{N \times N}$  telle que

- i)  $a_{ij} \leq 0$  pour tous  $i \neq j$ .
- ii)  $\sum_{1 \leq j \leq N} a_{ij} \geq 0$  pour  $i = 1, \dots, N$ .
- iii)  $A$  est inversible.

Alors  $A$  est une M-matrice.

*Démonstration.*

- Montrons d'abord que  $A + \varepsilon I$  est monotone, quelque soit  $\varepsilon > 0$ . On pose  $A + \varepsilon I = (a_{ij}^\varepsilon)_{i,j}$  avec  $a_{ij}^\varepsilon = a_{ij} + \varepsilon \delta_{ij}$ . Pour  $\varepsilon > 0$ , on a  $a_{ij}^\varepsilon \leq 0$  pour  $i \neq j$  et  $\sum_j a_{ij}^\varepsilon = \sum_j a_{ij} + \varepsilon \geq \varepsilon > 0$ . Par conséquent, d'après la proposition B.5, la matrice  $A + \varepsilon I$  est monotone.

- Montrons à présent que  $A$  est monotone par passage à la limite  $\varepsilon \rightarrow 0$ . La matrice  $A$  est inversible donc on peut écrire  $A + \varepsilon I = A(I + \varepsilon A^{-1})$ . De plus, on a  $\rho(\varepsilon A) = \varepsilon \rho(A) < 1$  pour  $\varepsilon$  suffisamment petit donc  $I + \varepsilon A^{-1}$  est inversible et  $(I + \varepsilon A^{-1})^{-1} = \sum_{k=0}^{\infty} (-\varepsilon A^{-1})^k$ . Par ailleurs,  $A + \varepsilon I$  est inversible avec

$$(A + \varepsilon I)^{-1} = (I + \varepsilon A^{-1})^{-1} A^{-1} = \left( \sum_{k=0}^{\infty} (-\varepsilon)^k (A^{-1})^k \right) A^{-1} = \left( I + \sum_{k=1}^{\infty} (-\varepsilon)^k (A^{-1})^k \right) A^{-1}.$$

Ainsi,  $(A + \varepsilon I)^{-1} - A^{-1} = \sum_{k=1}^{\infty} (-\varepsilon)^k (A^{-1})^k A^{-1}$  et donc

$$\left\| (A + \varepsilon I)^{-1} - A^{-1} \right\| \leq \left( \sum_{k=1}^{\infty} \varepsilon^k \|A^{-1}\|^k \right) \|A^{-1}\| = \frac{\varepsilon \|A^{-1}\|}{1 - \varepsilon \|A^{-1}\|} \|A^{-1}\| \rightarrow 0, \text{ quand } \varepsilon \rightarrow 0,$$

pour n'importe quelle norme matricielle  $\|\cdot\|$  subordonnée. Par conséquent chaque coefficient de  $(A + \varepsilon I)^{-1}$  qui est positif ou nul, tend vers le coefficient correspondant de  $A^{-1}$  qui est donc également positif ou nul. Ainsi  $A^{-1} \geq 0$  et  $A$  est une M-matrice.  $\square$

## B.4 Localisation des valeurs propres

**Théorème B.1 (Gershgorin-Hadamard)** Soit une matrice  $A$  carrée de taille  $N \times N$  à coefficients réels ou complexes.

1. Soit  $\lambda \in \mathbb{C}$  une valeur propre de  $A$ . Alors

$$\lambda \in \bigcup_{k=1}^N \mathcal{D}_k, \tag{B.1}$$

où  $\mathcal{D}_k$  sont les disques de Gershgorin définis par

$$\mathcal{D}_k = \left\{ z \in \mathbb{C}, |z - A_{k,k}| \leq \sum_{\substack{j=1 \\ j \neq k}}^N |A_{k,j}| \right\}.$$

2. On suppose que la matrice  $A$  est irréductible. Si une valeur propre  $\lambda$  est située sur la frontière de la réunion des disques de Gershgorin  $\mathcal{D}_k$  alors tous les cercles  $\partial \mathcal{D}_k$  passent par  $\lambda$ . Plus formellement, on a

$$\lambda \in \partial \left( \bigcup_{k=1}^N \mathcal{D}_k \right) \Rightarrow \lambda \in \bigcap_{k=1}^N \partial \mathcal{D}_k. \tag{B.2}$$

La figure B.1 illustre le cas où d'une valeur propre située sur le bord de l'union des disques de Gershgorin.

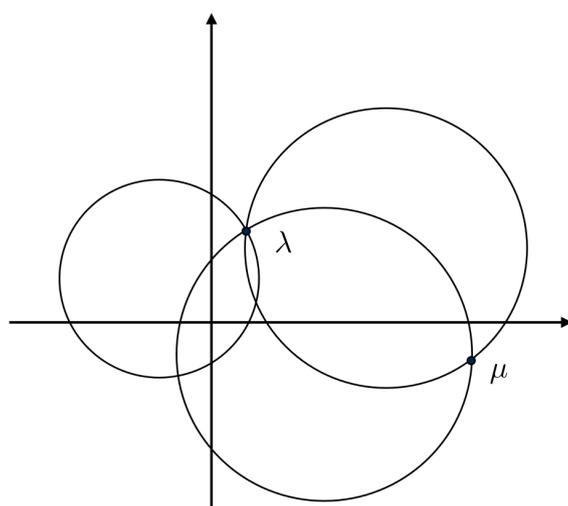


FIGURE B.1 – Disque de Gershgorin :  $\lambda$  peut être valeur propre mais  $\mu$  n'est pas valeur propre.



## Annexe C

### Quelques inégalités.

Pour les différentes inégalités classiques présentées ici, voir [14].

#### C.1 Inégalité de Cauchy-Schwarz

On note  $\langle \cdot, \cdot \rangle$  le produit scalaire et  $\| \cdot \|$  la norme euclidienne sur  $\mathbb{R}^n$ . Pour  $\mathbf{u} = (u_1, \dots, u_n)^\top$  et  $\mathbf{v} = (v_1, \dots, v_n)^\top$  deux vecteurs de  $\mathbb{R}^n$ , on rappelle que :

$$\langle \mathbf{u}, \mathbf{v} \rangle = \sum_{i=1}^n u_i v_i, \quad \|\mathbf{u}\| = \langle \mathbf{u}, \mathbf{u} \rangle^{1/2} = \left( \sum_{i=1}^n u_i^2 \right)^{1/2}.$$

*Inégalité de Cauchy-Schwarz* : pour tout  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ ,

$$\langle \mathbf{u}, \mathbf{v} \rangle \leq \|\mathbf{u}\| \|\mathbf{v}\|.$$

#### C.2 Inégalité de Young

Soient  $a$  et  $b$  deux nombres réels. Alors, pour tout  $\varepsilon > 0$ , on a

$$ab \leq \varepsilon a^2 + \frac{b^2}{4\varepsilon}$$

#### C.3 Inégalité de Gronwall

Soient  $u$  et  $v$  deux fonctions positives, définies et continues sur un intervalle  $I = [t_0, t_1] \subset \mathbb{R}$  telle que

$$\frac{du}{dt}(t) \leq v(t)u(t) \quad \text{pour tout } t \in I.$$

Alors, pour tout  $t \in I$ ,

$$u(t) \leq u(t_0) \exp \left( \int_{t_0}^{t_1} v(s) ds \right)$$

#### C.4 Inégalité de Gronwall discrète

1. Soient  $(a_n)_{n \in \mathbb{N}}$  et  $(b_n)_{n \in \mathbb{N}}$  deux suites positives telles que

$$a_n \leq \rho + \sum_{j=0}^{n-1} a_j b_j \quad \text{pour } \rho > 0, \forall n \in \mathbb{N}.$$

Alors, pour tout  $n \in \mathbb{N}$ ,

$$a_n \leq \rho \exp \left( \sum_{j=0}^{n-1} b_j \right). \tag{C.1}$$

2. Soit  $(a_n)_{n \in \mathbb{N}}$  une suite positive telle qu'il existe deux suites positives  $(\alpha_n)_{n \in \mathbb{N}}$ ,  $(\beta_n)_{n \in \mathbb{N}}$  pour lesquelles

$$a_{n+1} \leq (1 + \alpha_n)a_n + \beta_n, \quad \forall n \in \mathbb{N}.$$

Alors , pour tout  $n \in \mathbb{N}$ ,

$$a_n \leq \left( a_0 + \sum_{k=0}^{n-1} \beta_k \right) \exp \left( \sum_{k=0}^{n-1} \alpha_k \right). \quad (\text{C.2})$$

# Bibliographie

- [1] P. Blanc, R. Eymard, R. Herbin, *A staggered finite volume scheme on general meshes for the generalized Stokes problem in two space dimensions*, Int. J. Finite Volumes, 2 (2005), n°1, 31 pp.
- [2] H. Brézis, *Analyse fonctionnelle ; théorie et applications*, Masson, 1987.
- [3] P.G. Ciarlet, *Introduction à l'analyse matricielle et à l'optimisation*, Dunod, 5ème édition, 2007.
- [4] A. Ern, J.L. Guermond, *Éléments finis : théorie, applications, mise en oeuvre*, Mathématiques et Applications, Springer Berlin Heidelberg, 2002.
- [5] L.C. Evans, *Partial differential equations*, 1998.
- [6] R. Eymard, T. Gallouët, R. Herbin, *The finite volume method*, Handbook for Numerical Analysis, Ph. Ciarlet J.L. Lions eds, North Holland, 2000, 715-1022.
- [7] R. Eymard, R. Herbin, *A staggered finite volume scheme on general meshes for the Navier-Stokes equations in two space dimensions*, Int. J. Finite Volumes, 2 (2005), n°1, 19 pp.
- [8] R. A. Horn, Ch. R. Johnson, *Topics in Matrix Analysis*. Cambridge University Press, New York, 1991.
- [9] P. Lascaux, R. Théodor, *Analyse numérique matricielle appliquée à l'art de l'ingénieur*, Dunod, 2004.
- [10] R.J. Leveque, *Finite volume methods for hyperbolic problems*, 2002.
- [11] I. Mishev, *Finite Volume Methods on Voronoï Meshes*, Num. Meth. P.D.E, vol. 14, p.193-212,1998.
- [12] Protter, Weinberger, *Maximum principle in differential equations*, 1967.
- [13] A. Quarteroni, R. Sacco, F. Saleri, *Méthodes Numériques : Algorithmes, analyse et applications*, Springer Milan, 2008.
- [14] A. Quarteroni, A. Valli, *Numerical Approximation of Partial Differential Equations*, Springer Berlin Heidelberg, 2009.
- [15] M. Renardy, R.C. Rogers, *An introduction to partial differential equations*, 1993.