

## STAGE DE FIN D'ÉTUDE DE MASTER 1 CSMI

### SUJET: GALERKINE DISCONTINU EN 1D

Auteur :  
Adama DIENG

Superviseur :  
Joubine Aghili

Date : August 20, 2024

# Contents

<b>1</b>	<b>Context</b>	<b>3</b>
1.1	Introduction . . . . .	3
1.1.1	Informations Préliminaires . . . . .	3
1.1.2	Contexte du stage . . . . .	3
1.1.3	Approximation discontinue par méthode Galerkin Discontinue . . . . .	4
<b>2</b>	<b>Rappels sur la méthode des éléments finis MEF</b>	<b>5</b>
2.1	La méthode des éléments finis (MEF) . . . . .	5
2.1.1	Principe Fondamental de la MEF en 1D . . . . .	5
2.1.2	Discrétisation du Domaine . . . . .	5
2.1.3	Assemblage des Équations . . . . .	6
2.1.4	Extension en 2D . . . . .	6
<b>3</b>	<b>Méthode Galerkin Discontinue (DG)</b>	<b>7</b>
3.1	Equation différentielle ordinaire (EDO) en 1D . . . . .	7
3.1.1	Probleme abstrait et formulation faible . . . . .	7
3.1.2	Trace numérique $\hat{u}_h$ . . . . .	8
3.1.3	Assemblages locales des matrices . . . . .	11
3.1.4	Implementation du schema DG en 1D . . . . .	15
3.1.5	Résultats Numériques . . . . .	16
3.2	L'équation de transport 2D . . . . .	17
3.2.1	Formulation faible du problème . . . . .	17
3.2.2	Assemblage des matrices locales . . . . .	20
3.2.3	Assemblage d'un systeme lineaire: . . . . .	21
<b>4</b>	<b>Implémentation de la méthode DG pour l'équation de transport 2D</b>	<b>25</b>
<b>5</b>	<b>Conclusion et perspectives</b>	<b>29</b>

# Chapter 1

## Context

Pour réaliser ce rapport, on va suivre l'article de Cockburn [2]

### 1.1 Introduction

#### 1.1.1 Informations Préliminaires

##### Remerciements

Je tiens à exprimer ma gratitude à mon encadrant, M. Joubine Aghili, pour son accompagnement et ses enseignements durant ce stage. Je suis également très reconnaissant pour sa patience et sa disponibilité lorsque j'avais besoin d'aide, ainsi que pour m'avoir offert l'opportunité de travailler sur un sujet qui répond à mes intérêts académiques.

##### A propos du stage et du laboratoire de recherche

Ce stage fait partie du cursus du Master CSMI de l'Université de Strasbourg. Il s'agit d'un stage de 2 mois, qui a débuté le 3 juin et s'est terminé le 3 août 2024.

Il a eu lieu au sein du laboratoire **IRMA**, qui est une unité de recherche de CNRS (Centre national de la recherche scientifique) et de l'Université de Strasbourg. Le centre est situé au campus de l'université de Strasbourg, à Strasbourg, en France.

#### 1.1.2 Contexte du stage

Les méthodes numériques sont indispensables pour l'analyse et la modélisation des systèmes physiques complexes, en particulier dans le cadre de la résolution des équations différentielles partielles (EDP).

Ces dernières apparaissent fréquemment dans divers domaines scientifiques et techniques, allant de la mécanique des fluides à la propagation des ondes.

Parmi les méthodes numériques disponibles, les méthodes de Galerkin ont longtemps été privilégiées pour leur capacité à fournir des solutions précises avec un coût computationnel raisonnable. Toutefois, les méthodes classiques de Galerkin, comme les éléments finis continus, présentent des limitations lorsqu'elles sont confrontées à des solutions avec des discontinuités ou des gradients abrupts.

La méthode de Galerkin discontinue (DG) est une variante qui se distingue par sa capacité à traiter efficacement de tels défis de discontinuité. Contrairement aux méthodes continues, la méthode DG permet l'utilisation de fonctions de base discontinues, offrant ainsi une flexibilité accrue dans la représentation des solutions. Cette caractéristique est particulièrement utile dans les contextes où les solutions présentent des discontinuités naturelles, telles que les chocs ou les interfaces de phase, comme c'est souvent le cas dans les équations de conservation hyperboliques.

Ce rapport vise à explorer en détail la méthode de Galerkin discontinue dans le contexte unidimensionnel (1D). L'objectif principal est de fournir une compréhension claire et complète des principes sous-jacents de cette méthode, de démontrer son implémentation pratique, et d'examiner ses avantages par rapport aux autres méthodes numériques. Nous aborderons également les défis spécifiques associés à la mise en œuvre de la méthode DG, notamment le traitement des flux numériques et la stabilisation des solutions.

En première partie, nous introduirons les fondements théoriques de la méthode de Galerkin discontinue, en expliquant les concepts clés et en mettant en lumière ses avantages distinctifs. Ensuite, la deuxième partie sera consacrée à l'implémentation pratique de cette méthode en 1D, avec une description détaillée des étapes de discrétisation et des exemples d'application. Enfin, nous conclurons par une discussion sur les résultats obtenus, les limitations rencontrées, et les perspectives d'amélioration et d'extension de la méthode DG.

Cette exploration détaillée nous permettra de mieux comprendre comment la méthode de Galerkin discontinue peut être utilisée pour résoudre efficacement une large gamme de problèmes d'EDP, en mettant en lumière son potentiel et ses domaines d'application future.

### 1.1.3 Approximation discontinue par méthode Galerkin Discontinue

L'idée fondamentale de la méthode Galerkin Discontinu (DG) consiste, comme pour une méthode des éléments finis, à mettre la problème sous forme variationnelle, et à réaliser une approximation polynomiale de la solution du système.

Toutefois, contrairement aux éléments finis, où la solution est continue sur tout le domaine, l'approximation DG est une approximation locale dans chaque cellule du maillage, la valeur de la solution aux interfaces est à priori **multi-valuée** et la solution globale est donc discontinue sur les faces des éléments volumiques. De ce fait, des flux numériques sont calculés aux interfaces de chaque cellule, tel que cela est réaliser dans les méthodes de type volumes finis.

Par exemple en 1D, l'idée est de diviser le domaine de calcul (un intervalle) en plusieurs elements  $[x^{k-1}, x^k], [x^k, x^{k+1}], etc$  sur lesquels on va construire des  $u^k$  qui sont des approximations de la solution dans chaque élément. **figure1.1.** Contrairement aux méthodes d'éléments finis classiques, la méthode DG autorise des discontinuités aux frontières des éléments. Cela signifie que la valeur de la solution peut changer brusquement d'un élément à l'autre. Les points  $x^{k-1}, x^k, x^{k+1} etc$  sont les frontières (**les interfaces**) entre les éléments et les solution  $u^k$  sont des polynomes de degré  $k$  sur chaque élément.

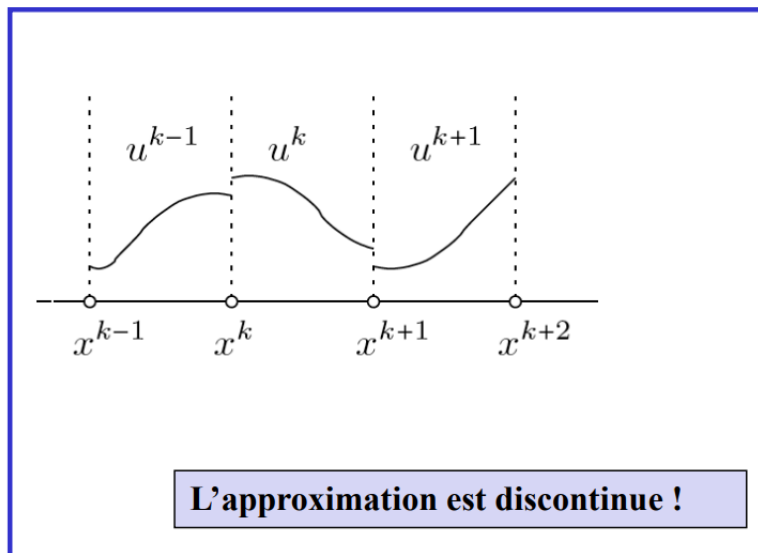


Figure 1.1: approximation discontinue en 1D [1]

## Chapter 2

# Rappels sur la méthode des éléments finis MEF

### 2.1 La méthode des éléments finis (MEF)

La méthode des éléments finis (MEF) est une technique numérique utilisée pour trouver des solutions approximatives aux équations différentielles partielles (EDP) et intégrales. Elle est largement employée dans les domaines de la mécanique, la physique, l'ingénierie, et plus particulièrement dans la résolution de problèmes liés aux déformations, à la diffusion de la chaleur, à la dynamique des fluides, et autres phénomènes physiques modélisables par des EDP.

#### 2.1.1 Principe Fondamental de la MEF en 1D

En une dimension (1D), la MEF consiste à diviser le domaine de calcul, qui est souvent un intervalle sur l'axe des abscisses, en un certain nombre de sous-domaines appelés éléments finis. Sur chaque élément, la solution est approximée par une fonction de forme simple, généralement un polynôme.

Soit une équation différentielle ordinaire (EDO) linéaire de la forme :

$$-\frac{d}{dx} \left( k(x) \frac{du}{dx} \right) + c(x)u = f(x), \quad x \in (a, b)$$

avec les conditions aux limites appropriées. L'idée principale est de multiplier cette équation par une fonction test  $v(x)$ , puis d'intégrer sur le domaine  $[a, b]$ . Cela conduit à la formulation dite **faible** ou **variationnelle** du problème. Celle-ci est de la forme:

$$-\int_a^b \frac{d}{dx} \left( k(x) \frac{du}{dx} \right) v(x) dx + \int_a^b c(x)u(x)v(x) dx = \int_a^b f(x)v(x) dx$$

pour tout  $v(x)$  appartenant à un espace de fonctions tests approprié.

#### 2.1.2 Discrétisation du Domaine

Le domaine  $[a, b]$  est discrétisé en un nombre fini de sous-intervalles  $[x_i, x_{i+1}]$ , formant ainsi un maillage. La solution approchée  $u_h(x)$  est construite comme une combinaison linéaire de fonctions de forme associées à chaque nœud du maillage :

$$u_h(x) = \sum_{i=1}^n u_i \phi_i(x)$$

où  $\phi_i(x)$  sont les fonctions de forme locales, et  $u_i$  sont les coefficients inconnus (les valeurs de la solution approchée aux nœuds).

### 2.1.3 Assemblage des Équations

En substituant l'expression de  $u_h(x)$  dans la formulation faible et en choisissant les fonctions tests  $v(x)$  égales aux fonctions de forme  $\phi_j(x)$ , on obtient un système d'équations linéaires pour les coefficients  $u_i$ . Ce système est généralement écrit sous forme matricielle :

$$\mathbf{K}\mathbf{u} = \mathbf{F}$$

où  $\mathbf{K}$  est la matrice de raideur,  $\mathbf{u}$  est le vecteur des inconnues, et  $\mathbf{F}$  est le vecteur des forces.

### 2.1.4 Extension en 2D

La méthode des éléments finis peut être étendue en deux dimensions (2D) en considérant des éléments finis comme des triangles ou des quadrilatères couvrant le domaine de calcul. La formulation et les concepts restent similaires, mais nécessitent une intégration numérique sur **des surfaces plutôt que sur des segments**.

## Chapter 3

# Méthode Galerkin Discontinue (DG)

### 3.1 Equation différentielle ordinaire (EDO) en 1D

#### 3.1.1 Probleme abstrait et formulation faible

Considérons une EDO en 1D de la forme:

$$\begin{cases} \frac{d}{dt}u(t) = f(t)u(t), & t \in (0, T), \\ u(0) = u_0, \end{cases} \quad (3.1)$$

On suppose que nous voulons déterminer une approximation  $u_h$  de la solution  $u$  de l'EDO en utilisant une méthode de Galerkin.

Pour cela, on va faire une discrétisation de l'intervalle  $(0, T)$  en petits intervalles  $I^n = (t_n, t_{n+1})$ , pour tout  $n = 0, 1, \dots, N-1$ .

On fait une sorte de reformulation faible de l'EDO en multipliant l'équation par une fonction  $v$  (polynôme de degré au plus égal à  $k^n$ ), car on cherche une fonction approximative  $u_h$  qui est de degré au plus égal à  $k^n$  sur chaque intervalle  $I^n$ . Ensuite on intègre sur chaque intervalle  $I^n$ .

On a donc:

$$\int_{I^n} \frac{d}{dt}u(t)v(t)dt = \int_{I^n} f(t)u(t)v(t)dt \quad (3.2)$$

En utilisant la formule d'intégration par partie, on a:

$$-\int_{I^n} u(t)\frac{d}{dt}v(t)dt + [u(t)v(t)]_{t_n}^{t_{n+1}} = \int_{I^n} f(t)u(t)v(t)dt \quad (3.3)$$

On cherche à approcher  $u$  par  $u_h$  tel que:

$$u_h = \sum_{n=0}^N u_i \phi_i(t) \quad (3.4)$$

Où les fonctions  $\tilde{\phi}_i$  sont des polynômes par morceaux de degré au plus égal à  $k^n$ , connue localement sur chaque intervalle  $I^n$  mais elles sont discontinues globalement sur  $(0, T)$ .

Dans la discrétisation de l'intervalle  $(0, T)$ , il est important aussi d'observer que la frontière de chaque élément  $I^n$  est partagée par deux intervalles.

**Par conséquent**, il est nécessaire de prendre en compte ce qui se passe aux interfaces de ces intervalles  $I^n$ .

Sur la base de cette partition, les méthodes DG permettent de gérer les valeurs des solutions aux interfaces des intervalles  $I^n$  en introduisant le terme de **trace numérique**  $\hat{u}_h$ .

On cherche alors  $u_h$  qui, sur l'intervalle  $I^n$ , est de degré au plus égal à  $k^n$  déterminée en imposant que:

$$-\int_{I^n} u_h(s)\frac{d}{ds}v(s)ds + \hat{u}_h v|_{t_n}^{t_{n+1}} = \int_{I^n} f(s)u_h(s)v(s)ds \quad (3.5)$$

pour tout polynome  $v$  de degré au plus égale à  $k^n$ . Pour completer cette definition de la méthode DG, il est necessaire de determiner la trace numérique  $\hat{u}_h$ . Notons que la méthode etablit un lien entre les valeurs de  $u_h$  dans les differents intervalles uniquement par le biais de la trace numérique  $\hat{u}_h$ .

### 3.1.2 Trace numérique $\hat{u}_h$

Pour une EDO, l'information voyage " du passé vers le futur", on peut donc determiner  $\hat{u}_h$  comme suit:

$$\hat{u}_h(t_n) = \begin{cases} u_0, & \text{si } t_n = 0, \\ \lim_{\epsilon \rightarrow 0} u_h(t_n - \epsilon), & \text{sinon.} \end{cases} \quad (3.6)$$

Cela achève la definition de la méthode DG. Dans cet exemple, on peut deja identifier les elements clés de la méthode :

- L'utilisation d'approximations discontinues pour la solution  $u_h$ ;
- L'application de l'équation différentielle sur chaque intervalle via une formulation faible de Galerkin;
- L'introduction et la definition adéquate de ce qu'on appelle la trace numérique  $\hat{u}_h$ .

La sélection de la trace numérique est probablement l'aspect le plus sensible et crucial de la méthode DG, car elle peut influencer sa cohérence, sa stabilité, et même sa précision. Le choix simple que nous avons effectué ici est bien adapté à ce cas et conduit à une méthode efficace.

On aborde d'abord la question de **la consistance** de la méthode de DG. Comme c'est habituel pour la majorité des méthodes des éléments finis, une méthode est dite consistante si la solution exacte  $u$  de l'EDO peut etre remplacée par la solution approximative  $u_h$  dans la formulation faible de Galerkin. Nous constatons immédiatement que cela est vérifié si et seulement si  $\hat{u} = u$ .

Quant à **la stabilité** de la méthode, un aspect plus delicat, l'approche consiste à etabli d'abord une propriété de stabilité pour l'équation (3.1) avant de l'étendre à la méthode DG (3.5) en definissant de maniere appropriée la trace numérique  $\hat{u}_h$ . Si on multiplions l'EDO par  $u$  et integrant sur  $(0, T)$ , on obtient:

$$\frac{1}{2}u^2(T) - \frac{1}{2}u_0^2 = \int_0^T f(s)u^2(s)ds \quad (3.7)$$

où la stabilité- $L^\infty$  de l'EDO est garantie par le fait que le terme de droite est positif. Pour la méthode DG, il est necessaire d'obtenir un égalité similaire. Pour y parvenir, il suffit de poser  $v = u_h$  dans l'équation (3.5), integrant par partie et en sommant sur  $n$ . On obtient alors:

$$\sum_{n=0}^{N-1} \left( -\frac{1}{2}u_h^2 + \hat{u}_h u_h \right) \Big|_{t_n}^{t_{n+1}} = \frac{1}{2}u_h^2(T^-) - \frac{1}{2}u_h^2(T^-) + \sum_{n=0}^{N-1} \left( -\frac{1}{2}u_h^2 + \hat{u}_h u_h \right) \Big|_{t_n}^{t_{n+1}} + \frac{1}{2}u_0^2 - \frac{1}{2}u_0^2 = \int_0^T f(s)u_h^2(s)ds \quad (3.8)$$

Autrement dit :

$$\sum_{n=0}^{N-1} \left( -\frac{1}{2}u_h^2 + \hat{u}_h u_h \right) \Big|_{t_n}^{t_{n+1}} = \frac{1}{2}u_h^2(T^-) + \Theta_h(T) - \frac{1}{2}u_0^2 = \int_0^T f(s)u_h^2(s)ds \quad (3.9)$$

où

$$\Theta_h(T) = -\frac{1}{2}u_h^2(T^-) + \sum_{n=0}^{N-1} \left( -\frac{1}{2}u_h^2 + \hat{u}_h u_h \right) \Big|_{t_n}^{t_{n+1}} + \frac{1}{2}u_0^2$$

**Quel critere sur  $\Theta_h(t)$  pour determiner cet trace numérique  $\hat{u}_h$ ?**

Il est important de noter que si  $\Theta_h(T)$  etait une quantité non-négative, l'égalité ci-dessus impliquerait la stabilité de la méthode DG. Autrement dit, la stabilité de la méthode DG est garantie si nous pouvons definir la trace numérique  $\hat{u}_h$  de sorte que  $\Theta_h(T) \geq 0$ , [2].



**Proposition 3.1: (L<sup>2</sup>–stabilité)[6]**

Nous avons:

$$\frac{1}{2} \|u_h\|_{L^2(0,T)}^2 + \Theta_h(T) \leq \frac{1}{2} \|u_0\|_{L^2(0,T)}^2$$

où  $\Theta_h(T) \geq 0$ .

En supposant que:

$$u_h(t) = u_0, \quad \text{si } t < 0,$$

et **en notant**:

$$\{u_h\} = \frac{1}{2}(u_h^+ + u_h^-), \quad [u_h] = u_h^- - u_h^+, \quad u_h^\pm = \lim_{\epsilon \rightarrow 0} u_h(t \pm \epsilon) \quad (3.10)$$

On peut ainsi réécrire  $\Theta_h(T)$  comme suit:

$$\begin{aligned} \Theta_h(T) = & -\frac{1}{2} u_h^2(T^-) \\ & + \left( -\frac{1}{2} u_h^2(T^-) + \hat{u}_h(T) u_h(T^-) \right) \\ & + \sum_{n=1}^{N-1} \left( -\frac{1}{2} [u_h]^2 + \hat{u}_h [u_h] \right) (t_n) \\ & - \left( -\frac{1}{2} u_h^2(0^+) + \hat{u}_h(0) u_h(0^+) \right) + \frac{1}{2} u_0^2 \end{aligned}$$

Avec l'identité suivante (obtenue avec les notations précédentes):

$$[u_h^2] = (u_h^+)^2 - (u_h^-)^2 = (u_h^+ - u_h^-)(u_h^+ + u_h^-) = 2 \{u_h\} [u_h]$$

On obtient:

$$\Theta_h(T) = (\hat{u}_h(T) - u_h(T^-)) u_h(T^-) + \sum_{n=1}^{N-1} ((\hat{u}_h - \{u_h\}) [u_h]) (t_n) - (\hat{u}_h(0) - u_0) u_h(0^+) + \frac{1}{2} [u_h]^2(0) \quad (3.11)$$

Il est clair que [2] si nous prenons:

$$\hat{u}_h(t_n) = \begin{cases} u_0, & \text{si } t_n = 0, \\ (\{u_h\} + C^n [u_h]) (t_n), & \text{si } t_n \in ]0, T[, \\ u_h(T^-), & \text{si } t_n = T, \end{cases} \quad (3.12)$$

Où  $C^n \geq 0$ , on peut voir que:

$$\Theta_h(T) = \sum_{n=1}^{N-1} C^n [u_h]^2(t_n) + \frac{1}{2} [u_h]^2(0) = \sum_{n=0}^{N-1} C^n [u_h]^2(t_n) \geq 0$$

avec  $C^0 = \frac{1}{2}$ .

En effet, en substituant  $\hat{u}_h$  dans l'expression de  $\Theta_h(T)$ , on obtient:

- Pour  $t^n = T$ :

$$\hat{u}_h(T) = u_h(T^-)$$

donc,

$$(\hat{u}_h(T) - u_h(T^-)) u_h(T^-) = (u_h(T^-) - u_h(T^-)) u_h(T^-) = 0$$

- Pour  $t_n = 0$ :

$$\hat{u}_h(0) = u_0$$

donc,

$$(\hat{u}_h(0) - u_0)u_h(0^+) = (u_0 - u_0)u_h(0^+) = 0$$

- Pour  $t_n \in (0, T)$ :

$$\hat{u}_h(t_n) = (\{u_h\} + C^n[u_h])(t_n)$$

Alors,

$$(\hat{u}_h - \{u_h\})(t_n) = (\{u_h\} + C^n[u_h] - \{u_h\})(t_n) = C^n[u_h](t_n)$$

Et donc,

$$(\hat{u}_h - \{u_h\})[u_h](t_n) = C^n[u_h]^2(t_n)$$

Ainsi on a bien:

$$\Theta_h(T) = 0 + \sum_{n=1}^{N-1} C^n[u_h]^2(t_n) - 0 + \frac{1}{2}[u_h]^2(0) = \sum_{n=0}^{N-1} C^n[u_h]^2(t_n) \geq 0$$

Le choix des  $C^n$  reste important pour une meilleure convergence de la méthode DG. Avec quelques tests numériques suivant les valeurs de  $C^n$ , on va voir que la solution est proche de la solution exacte selon les valeurs de  $C^n$ .

En effet, il est possible de définir la trace numérique  $\hat{u}_h$  afin de garantir la stabilité de la méthode. On note que le choix de  $C^n = \frac{1}{2}$  correspond à la trace numérique choisie initialement **(3.6)**, soit  $\hat{u}_h(t) = \hat{u}_h(t^-)$ .

Les méthodes DG qui en résulte sont non seulement stables mais aussi cohérentes pour toutes les valeurs de  $C^n \geq 0$ , car la condition  $\hat{u} = u$  est satisfaite.

De plus, en cherchant à assurer la stabilité, il a été trouvé naturellement que la trace numérique  $\hat{u}_h$  ne peut dépendre que des deux traces de  $\hat{u}$  à  $t$ , plus précisément de  $\hat{u}(t^-)$  et  $\hat{u}(t^+)$ .

Parlons maintenant de l'importance du choix des coefficients  $C^n$  sur la précision de la méthode. Il a été démontré que si  $C^n = \frac{1}{2}$ , l'ordre de la méthode aux points  $t_n$  est de  $2k + 1$ . Cependant, si  $C^n = 0$ , l'ordre est de  $2k + 2$  [1]. Il est important de noter dans ce cas, on perd la possibilité de résoudre le flux  $\hat{u}_h$  intervalle par intervalle, car si  $C^n \neq \frac{1}{2}$ , il devient nécessaire de résoudre  $\hat{u}_h$  sur l'ensemble du domaine de calcul  $(0, T)$  en une seule fois.

Par conséquent, si on souhaite une méthode DG qui soit cohérente, stable et capable de traiter la solution approximative intervalle par intervalle, il est nécessaire de fixer  $C^n = \frac{1}{2}$ .

Pour finir, on souhaite souligner trois propriétés essentielles des méthodes DG qui s'appliquent également aux cas multidimensionnels et à divers types de problèmes. Premièrement, la solution approchée des méthodes DG n'est pas contrainte par la continuité entre les éléments. Cela permet à la méthode d'être hautement parallélisable, particulièrement pour les problèmes hyperboliques dépendant du temps, et facilite l'utilisation de différentes approximations dans chaque élément, ce qui est idéal pour l'adaptativité hp.

Deuxièmement, les méthodes DG sont localement conservatrices. Cette propriété découle du fait que la méthode applique l'équation élément par élément et utilise une trace numérique. Dans notre cadre simple, cette propriété se lit comme suit:

$$\hat{u}|_{t_n}^{t_{n+1}} = \int_{I^n} f(s)u_h(s)ds,$$

et est obtenu en posant  $v = 1$  dans l'équation **(3.5)**.

Enfin, la troisième qui n'a pas été abordée suffisamment dans la littérature disponible sur les méthodes DG, concerne la relation étroite entre les résidus de  $u_h$  à l'intérieur des intervalles et ses discontinuités aux interfaces. Pour mettre en évidence cette relation, on intègre par partie le premier terme de l'équation **(3.5)** afin d'obtenir:

$$\int_{I^n} u_h(s) \frac{d}{ds} v(s) ds - u_h v|_{t_n}^{t_{n+1}} + \hat{u}_h v|_{t_n}^{t_{n+1}} = \int_{I^n} f(s)u_h(s)v(s)ds$$

$\Rightarrow$

$$\int_{I^n} v(s) \frac{d}{ds} u_h(s) ds - \int_{I^n} f(s) u_h(s) v(s) ds = (u_h - \hat{u}_h) v|_{t_n}^{t_{n+1}}$$

$\Rightarrow$

$$\int_{I^n} R_h(s) v(s) ds = (u_h - \hat{u}_h) v|_{t_n}^{t_{n+1}}$$

où  $R_h = (\frac{d}{ds} u_h - f(s) u_h)$  est le résidu de  $u_h$  sur l'intervalle  $I^n$ . Si on prends  $v = 1$  et on utilise la definition de la trace numérique  $\hat{u}_h$ , on obtient:

$$\int_{I^n} R_h(s) ds = \llbracket u_h \rrbracket(t_n)$$

En d'autres termes, l'intégrale du résidu de  $u_h$  sur l'intervalle  $I^n$  est égale au saut de  $u_h$ , noté  $\llbracket u_h \rrbracket$ , en  $t_n$ . Cela signifie que si l'EDO a été bien approximée dans l'intervalle  $I^n$ , alors le saut  $\llbracket u_h \rrbracket$  sera très faible. En revanche, si ce n'est pas le cas, c'est-à-dire si l'EDO n'a pas été bien approximée dans cette intervalle, le saut  $\llbracket u_h \rrbracket$  sera grand, et la méthode DG deviendra plus dissipative, comme on peut bien le constater directement à partir de l'identité de stabilité:

$$\frac{1}{2} u_h^2(T^-) - \frac{1}{2} \sum_{n=0}^{N-1} \llbracket u_h \rrbracket^2(t_n) - \frac{1}{2} u_0^2 = \int_0^T f(s) u_h^2(s) ds$$

Pour simplifier, on peut dire que les discontinuités  $\llbracket u_h \rrbracket$  fonctionnent comme des amortisseurs, stabilisant la méthode DG lorsque l'intégration devient difficile et que l'EDO n'est pas correctement approchée. Cela assure une grande stabilité à la méthode tout en maintenant sa précision. En passant aux dimensions supérieures, le résidu devient une fonction linéaire plus complexe des discontinuités, mais le mécanisme dissipatif décrit reste principalement le même. Ce mécanisme est aussi présent dans les méthodes d'éléments finis stabilisées, couramment utilisées pour les problèmes de convection et de second ordre.

### 3.1.3 Assemblages locales des matrices

On discrétise l'intervalle  $(0, T)$  en  $N - 1$  intervalles  $I^n$ :

$$0 = t_0 < t_1 < \dots < t_{N-1} = T$$

où les elements sont donnés par:

$$I^n = (t_n, t_{n+1}), \quad n \in [0, N-2]$$

On a alors:

$$(0, T) = \bigcup_{n=0}^{N-2} I^n$$

L'equation (3.5) peut se réécrire comme suit:

$$\int_{I^n} u_h(s) \frac{d}{ds} v(s) ds + \int_{I^n} f(s) u_h(s) v(s) ds - \hat{u}_h v|_{t_n}^{t_{n+1}} = 0 \quad (3.13)$$

En substituant (3.4) dans (3.13), on obtient:

$$\int_{I^n} \left( \sum_{i=0}^p u_i^n \phi_i^n(s) \right) \frac{d}{ds} \phi_j^n(s) ds + \int_{I^n} f(s) \left( \sum_{i=0}^p u_i^n \phi_i^n(s) \right) \phi_j^n(s) ds - \hat{u}_h v|_{t_n}^{t_{n+1}} = 0$$

qu'on peut réécrire encore:

$$\sum_{i=0}^p \int_{I^n} u_i^n \left( \frac{d}{ds} \phi_j^n(s) + f(s) \phi_j^n(s) \right) \phi_i^n(s) ds + \hat{u}_h(t_n) \phi_j^n(t_n) - \hat{u}_h(t_{n+1}) \phi_j^n(t_{n+1}) = 0 \quad (3.14)$$

Sur chaque intervalle de la partition, on on approche "localement" la solution:

- Sur le premier intervalle  $I^0 = (t_0, t_1)$ :

$$u_h = \sum_{i=0}^p u_i^0 \phi_i^0, \quad v = \phi_j^0$$

l'équation (3.14) devient:

$$\sum_{i=0}^p u_i^0 \int_{I^0} u_i^0 \left( f(s) \phi_j^0(s) + \frac{d}{ds} \phi_j^0(s) \right) ds + \hat{u}_h(t_0) \phi_j^0(t_0) - \hat{u}_h(t_1) \phi_j^0(t_1) = 0 \quad (3.15)$$

or d'après (3.12), on a:

$$\hat{u}_h(t_0) = u_0$$

et

$$\hat{u}_h(t_1) \phi_j^0(t_1^-) = (\{u_h\} + C^1 [u_h]) \phi_j^0(t_1^-)$$

alors,

$$\sum_{i=0}^p u_i^0 \int_{I^0} u_i^0 \left( f(s) \phi_j^0(s) + \frac{d}{ds} \phi_j^0(s) \right) ds + u_0 \phi_j^0(t_0) - \hat{u}_h(t_1) \phi_j^0(t_1) = 0$$

Si on aborde les notations dans (3.10), on obtient:

$$\begin{aligned} -\hat{u}_h(t_1) \phi_j^0(t_1^-) &= -(\{u_h\} + C^1 [u_h]) \phi_j^0(t_1^-) \\ &= -\left( \frac{1}{2} (u_h(t_1^-) + u_h(t_1^+)) + C^1 (u_h(t_1^-) - u_h(t_1^+)) \right) \phi_j^0(t_1^-) \\ &= \left[ \left( -\frac{1}{2} - C^1 \right) u_h(t_1^-) + \left( -\frac{1}{2} + C^1 \right) u_h(t_1^+) \right] \phi_j^0(t_1^-) \\ &= \left[ \left( C^1 - \frac{1}{2} \right) u_h(t_1^+) + \left( -\frac{1}{2} - C^1 \right) u_h(t_1^-) \right] \phi_j^0(t_1^-) \end{aligned}$$

Si on remplace  $u_h$  par son expression, on obtient:

$$-\hat{u}_h(t_1) \phi_j^0(t_1^-) = \left( C^1 - \frac{1}{2} \right) \sum_{i=0}^p \mathbf{u}_i^1 \phi_i^1(t_1^+) \phi_j^0(t_1^-) + \left( -\frac{1}{2} - C^1 \right) \sum_{i=0}^p \mathbf{u}_i^0 \phi_i^0(t_1^-) \phi_j^0(t_1^-)$$

Par suite, on peut reecire l'équation (3.15) sur l'intervalle  $I^0$  comme suit:

$$\begin{aligned} \sum_{i=0}^p \left\{ \int_{I^0} \left( f(s) \phi_i^0(s) + \frac{d}{ds} \phi_i^0(s) \right) \phi_j^0(s) ds \right\} \mathbf{u}_i^0 + \left( -\frac{1}{2} - C^1 \right) \sum_{i=0}^p \phi_i^0(t_1^-) \phi_j^0(t_1^-) \mathbf{u}_i^0 \\ + \left( C^1 - \frac{1}{2} \right) \sum_{i=0}^p \phi_i^1(t_1^+) \phi_j^0(t_1^-) \mathbf{u}_i^1 = -u_0 \phi_j^0(t_0^+) \end{aligned}$$

$\Rightarrow$

$$\begin{aligned} &\boxed{\sum_{i=0}^p \left\{ \int_{I^0} \left( f(s) \phi_i^0(s) + \frac{d}{ds} \phi_i^0(s) \right) \phi_j^0(s) ds + \left( -\frac{1}{2} - C^1 \right) \phi_i^0(t_1^-) \phi_j^0(t_1^-) \right\} \mathbf{u}_i^0} \\ &+ \boxed{\left( C^1 - \frac{1}{2} \right) \sum_{i=0}^p \phi_i^1(t_1^+) \phi_j^0(t_1^-) \mathbf{u}_i^1} \\ &= \boxed{-u_0 \phi_j^0(t_0^+)} \end{aligned}$$

Sous forme matricielle, on obtient:

$$\begin{bmatrix} Q_1 & Q_2 & 0 & \dots & 0 \\ 0 & 0 & 0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \end{bmatrix} \begin{bmatrix} \mathbf{u}^0 \\ \mathbf{u}^1 \\ \vdots \end{bmatrix} = \begin{bmatrix} \mathbf{b} \\ 0 \\ \vdots \end{bmatrix}$$

où

$$(Q_1)_{ij} = \left\{ \int_{I_0} \left( f(s) \phi_i^0(s) + \frac{d}{ds} \phi_i^0(s) \right) \phi_j^0(s) ds \right\} + \left( -\frac{1}{2} - C^1 \right) \phi_i^0(t_1^-) \phi_j^0(t_1^-)$$

$$(Q_2)_{ij} = \left( C^1 - \frac{1}{2} \right) \phi_i^1(t_1^+) \phi_j^0(t_1^-)$$

et

$$(\mathbf{b})_j = -u_0 \phi_j^0(t_0^+)$$

**Remarque:** Il faudra noter que  $Q_1$  et  $Q_2$  sont des blocs de matrices  $(p \times p)$ , et  $\mathbf{b}$  est un vecteur de taille  $p$  et que  $U^0$  et  $U^1$  sont aussi des vecteurs de taille  $p$ .

- Sur les intervalles interieurs  $I^n = (t_n, t_{n+1})$  avec  $n \in [1, N-3]$ :

$$u_h = \sum_{i=0}^p u_i^n \phi_i^n, \quad v = \phi_j^n$$

L'équation (3.14) devient:

$$\sum_{i=0}^p u_i^n \int_{I^n} \left( f(s) \phi_j^n(s) + \frac{d}{ds} \phi_j^n(s) \right) \phi_i^n(s) ds + \hat{u}_h(t_n) \phi_j^n(t_n) - \hat{u}_h(t_{n+1}) \phi_j^n(t_{n+1}) = 0 \quad (3.16)$$

or, toujours d'après (3.12), on a:

$$\begin{aligned} \hat{u}_h(t_n) &= \{u_h\} + C^n [u_h] \\ &= \frac{1}{2} (u_h(t_n^-) + u_h(t_n^+)) + C^n (u_h(t_n^-) - u_h(t_n^+)) \\ &= \left( \frac{1}{2} + C^n \right) u_h(t_n^-) + \left( \frac{1}{2} - C^n \right) u_h(t_n^+) \end{aligned}$$

Donc,

$$\begin{aligned} \hat{u}_h(t_n) \phi_j^n(t_n) &= \left[ \left( \frac{1}{2} + C^n \right) u_h(t_n^-) + \left( \frac{1}{2} - C^n \right) u_h(t_n^+) \right] \phi_j^n(t_n^+) \\ &= \left[ \left( \frac{1}{2} + C^n \right) \sum_{i=0}^p \mathbf{u}_i^{n-1} \phi_i^{n-1}(t_n^-) + \left( \frac{1}{2} - C^n \right) \sum_{i=0}^p \mathbf{u}_i^n \phi_i^n(t_n^+) \right] \phi_j^n(t_n^+) \\ &= \left( \frac{1}{2} + C^n \right) \sum_{i=0}^p \mathbf{u}_i^{n-1} \phi_i^{n-1}(t_n^-) \phi_j^n(t_n^+) + \left( \frac{1}{2} - C^n \right) \sum_{i=0}^p \mathbf{u}_i^n \phi_i^n(t_n^+) \phi_j^n(t_n^+) \end{aligned}$$

et de même,

$$\begin{aligned} \hat{u}_h(t_{n+1}) \phi_j^n(t_{n+1}) &= \left[ \left( \frac{1}{2} + C^{n+1} \right) u_h(t_{n+1}^-) + \left( \frac{1}{2} - C^{n+1} \right) u_h(t_{n+1}^+) \right] \phi_j^n(t_{n+1}^+) \\ &= \left( \frac{1}{2} + C^{n+1} \right) \sum_{i=0}^p \mathbf{u}_i^n \phi_i^n(t_{n+1}^-) \phi_j^n(t_{n+1}^+) + \left( \frac{1}{2} - C^{n+1} \right) \sum_{i=0}^p \mathbf{u}_i^{n+1} \phi_i^{n+1}(t_{n+1}^+) \phi_j^n(t_{n+1}^+) \end{aligned}$$

En substituant ces deux expressions dans l'équation (3.16), on obtient:

$$\begin{aligned} & \left[ \sum_{i=0}^p \left\{ \left( \frac{1}{2} + C^n \right) \phi_i^{n-1}(t_n^-) \phi_j^n(t_n^+) \right\} \mathbf{u}_i^{n-1} \right] + \\ & \left[ \sum_{i=0}^p \left\{ \int_{I^n} \phi_i^n(s) \left( f(s) \phi_j^n(s) + \frac{d}{ds} \phi_j^n(s) \right) ds + \left( \frac{1}{2} - C^n \right) \phi_i^n(t_n^+) \phi_j^n(t_n^+) - \left( \frac{1}{2} + C^{n+1} \right) \phi_i^n(t_{n+1}^-) \phi_j^{n+1}(t_{n+1}^+) \right\} \mathbf{u}_i^n \right] \\ & + \left[ \sum_{i=0}^p \left\{ - \left( \frac{1}{2} - C^{n+1} \right) \phi_i^{n+1}(t_{n+1}^+) \phi_j^{n+1}(t_{n+1}^+) \right\} \mathbf{u}_i^{n+1} \right] = 0 \end{aligned}$$

Sous forme matricielle, on obtient:

$$\begin{bmatrix} 0 & \dots & \dots & \dots & \dots & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & Q1 & Q2 & Q3 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & \dots & \dots & \dots & \dots & 0 \end{bmatrix} \begin{bmatrix} \vdots \\ \mathbf{u}^{n-1} \\ \mathbf{u}^n \\ \mathbf{u}^{n+1} \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix}$$

Avec:

$$\begin{aligned} (Q1)_{ij} &= \left\{ \left( \frac{1}{2} + C^n \right) \phi_i^{n-1}(t_n^-) \phi_j^n(t_n^+) \right\} \\ (Q2)_{ij} &= \left\{ \int_{I^n} \phi_i^n(s) \left( f(s) \phi_j^n(s) + \frac{d}{ds} \phi_j^n(s) \right) ds + \left( \frac{1}{2} - C^n \right) \phi_i^n(t_n^+) \phi_j^n(t_n^+) - \left( \frac{1}{2} + C^{n+1} \right) \phi_i^n(t_{n+1}^-) \phi_j^{n+1}(t_{n+1}^+) \right\} \\ (Q3)_{ij} &= - \left( \frac{1}{2} - C^{n+1} \right) \phi_i^{n+1}(t_{n+1}^+) \phi_j^{n+1}(t_{n+1}^+) \end{aligned}$$

**Remarque:** Même remarque que pour l'assemblage sur l'intervalle  $I^0$ :  $Q1$ ,  $Q2$  et  $Q3$  sont des blocs de matrices  $(p \times p)$ , et les solutions  $U^{n-1}$ ,  $U^n$  et  $U^{n+1}$  sont des vecteurs de taille  $p$ .

- Sur le dernier intervalle  $I^{N-2} = (t_{N-2}, t_{N-1}) = (t_{N-2}, T)$ :

$$u_h = \sum_{i=0}^p u_i^{N-2} \phi_i^{N-2}, \quad v = \phi_j^{N-2}$$

L'équation (3.14) devient:

$$\sum_{i=0}^p \mathbf{u}_i^{N-2} \int_{I^{N-2}} \left( f(s) \phi_j^{N-2}(s) + \frac{d}{ds} \phi_j^{N-2}(s) \right) \phi_i^{N-2}(s) ds + \hat{u}_h(t_{N-2}) \phi_j^{N-2}(t_{N-2}) - \hat{u}_h(T) \phi_j^{N-2}(T) = 0 \quad (3.17)$$

or

$$\hat{u}_h(T) = u_h(T^-) \quad \text{d'apres (3.12)}$$

donc,

$$\hat{u}_h(T) \phi_j^{N-2}(T) = \sum_{i=0}^p \mathbf{u}_i^{N-2} \phi_i^{N-2}(T^-) \phi_j^{N-2}(T^-)$$

Par ailleurs,

$$\begin{aligned} \hat{u}_h(t^{N-2}) &= \{u_h\} + C^{N-2} [u_h] \\ &= \frac{1}{2} (u_h(t_{N-2}^-) + u_h(t_{N-2}^+)) + C^{N-2} (u_h(t_{N-2}^-) - u_h(t_{N-2}^+)) \end{aligned}$$

$$= \left(\frac{1}{2} + C^{N-2}\right)u_h(t_{N-2}^-) + \left(\frac{1}{2} - C^{N-2}\right)u_h(t_{N-2}^+)$$

donc,

$$\begin{aligned}\hat{u}_h(t_{N-2})\phi_j^{N-2}(t_{N-2}^+) &= \left[\left(\frac{1}{2} + C^{N-2}\right)u_h(t_{N-2}^-) + \left(\frac{1}{2} - C^{N-2}\right)u_h(t_{N-2}^+)\right]\phi_j^{N-2}(t_{N-2}^+) \\ &= \left[\left(\frac{1}{2} + C^{N-2}\right)\sum_{i=0}^p \mathbf{u}_i^{N-3}\phi_i^{N-3}(t_{N-2}^-) + \left(\frac{1}{2} - C^{N-2}\right)\sum_{i=0}^p \mathbf{u}_i^{N-2}\phi_i^{N-2}(t_{N-2}^+)\right]\phi_j^{N-2}(t_{N-2}^+) \\ &= \left(\frac{1}{2} + C^{N-2}\right)\sum_{i=0}^p \mathbf{u}_i^{N-3}\phi_i^{N-3}(t_{N-2}^-)\phi_j^{N-2}(t_{N-2}^+) + \left(\frac{1}{2} - C^{N-2}\right)\sum_{i=0}^p \mathbf{u}_i^{N-2}\phi_i^{N-2}(t_{N-2}^+)\phi_j^{N-2}(t_{N-2}^+)\end{aligned}$$

En remplaçant ces expressions dans l'équation (3.17), on obtient:

$$\begin{aligned}&\sum_{i=0}^p \left\{ \left(\frac{1}{2} + C^{N-2}\right)\phi_i^{N-3}(t_{N-2}^-)\phi_j^{N-2}(t_{N-2}^+)\right\} \mathbf{u}_i^{N-3} + \\ &\sum_{i=0}^p \left\{ \int_{I^{N-2}} \phi_i^{N-2}(s) \left( f(s)\phi_j^{N-2}(s) + \frac{d}{ds}\phi_j^{N-2}(s) \right) ds + \left(\frac{1}{2} - C^{N-2}\right)\phi_i^{N-2}(t_{N-2}^+)\phi_j^{N-2}(t_{N-2}^+) - \phi_i^{N-2}(T^-)\phi_j^{N-2}(T^-) \right\} \\ &= 0\end{aligned}$$

Sous forme matricielle, on obtient:

$$\begin{bmatrix} 0 & \cdots & \cdots & 0 \\ 0 & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & Q1 & Q2 \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \mathbf{u}^{N-3} \\ \mathbf{u}^{N-2} \end{bmatrix}$$

Avec:

$$\begin{aligned}(Q1)_{ij} &= \left\{ \left(\frac{1}{2} + C^{N-2}\right)\phi_i^{N-3}(t_{N-2}^-)\phi_j^{N-2}(t_{N-2}^+) \right\} \\ (Q2)_{ij} &= \int_{I^{N-2}} \phi_i^{N-2}(s) \left( f(s)\phi_j^{N-2}(s) + \frac{d}{ds}\phi_j^{N-2}(s) \right) ds \\ &\quad + \left(\frac{1}{2} - C^{N-2}\right)\phi_i^{N-2}(t_{N-2}^+)\phi_j^{N-2}(t_{N-2}^+) \\ &\quad - \phi_i^{N-2}(T^-)\phi_j^{N-2}(T^-)\end{aligned}$$

### 3.1.4 Implementation du schema DG en 1D

Le schéma a été implémenté en Python et disponible à l'adresse suivante:

[https://gitlab.math.unistra.fr/aghili/dg/-/blob/dieng-main-patch-68158/dg1d.py?ref\\_type=heads](https://gitlab.math.unistra.fr/aghili/dg/-/blob/dieng-main-patch-68158/dg1d.py?ref_type=heads)

Dans cette section, nous présentons une implémentation de la méthode de Galerkin discontinue (DG) en une dimension. Le code utilise des fonctions de base discontinues, la formulation faible, et des flux numériques pour résoudre l'EDO (3.1). Voici les différentes étapes de l'implémentation:

- **Etape 1:** On commence par initialiser la classe, définir les paramètres du problème et les fonctions de base.
- **Etape 2:** On construit les matrices locales  $Q1$ ,  $Q2$  et  $Q3$  selon chaque intervalle  $I^n$ .
- **Etape 3:** On assemble les matrices locales pour obtenir la matrice globale  $A$ .
- **Etape 4:** On résout le système linéaire  $A\mathbf{u} = \mathbf{b}$  et estime les erreurs.

### 3.1.5 Résultats Numériques

On reprends le problème (3.1) avec  $f(t) = -t^2$ ,  $u_0 = 1$  et  $T = 5$ .

Dans la figure (3.1), on affiche l'erreur  $\|u_h - u\|_{L^2}^2$  en fonction du pas d'espace  $h$  allant de  $h = \frac{1}{2^2}$  jusqu'à  $h = \frac{1}{2^8}$  à échelles logarithmiques pour  $k = 0, 1, \dots, 5$  et pour  $C = 0$ .

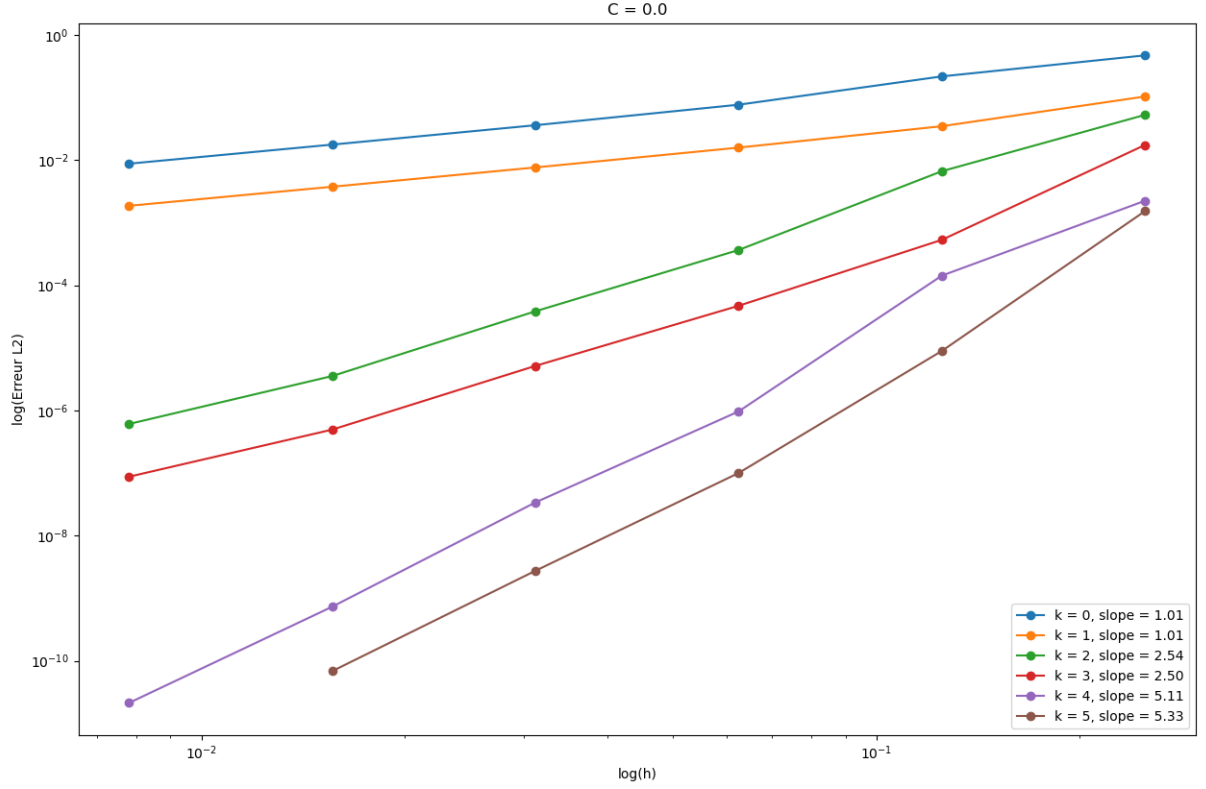


Figure 3.1: Erreurs et Pentés pour  $C = 0$  pour chaque ordre polynomial  $k$

De même dans la figure (3.2), on affiche les mêmes variations pour  $C = 0.5$ .



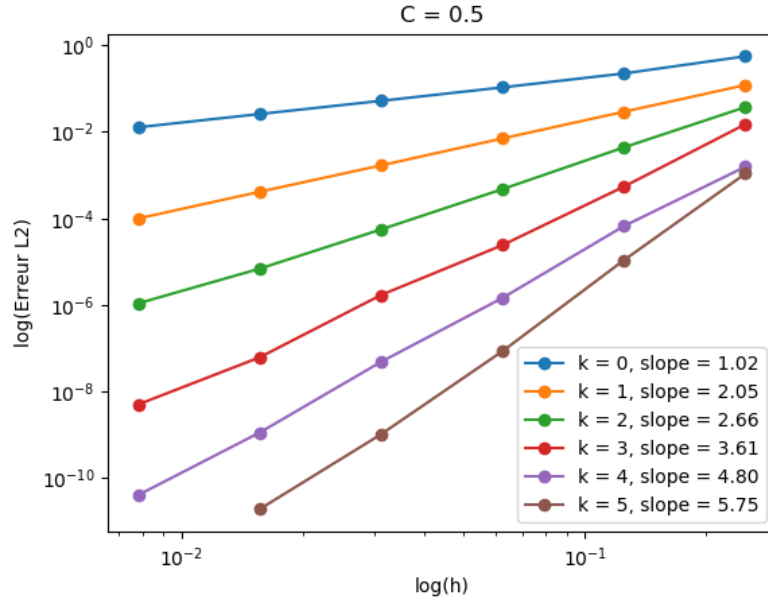


Figure 3.2: Erreurs et Pentes pour  $C = 0.5$  pour chaque ordre polynomial  $k$

**Interpretations:** Dans un premier temps, on peut remarquer que les ordres d'erreurs diminuent progressivement lorsque le degré  $k$  des polynômes de base augmente, ce qui confirme que les méthodes DG sont efficaces pour des polynômes de degré élevé.

Ensuite, on peut constater que le choix de  $C = 0.5$  conduit à des résultats **sous-optimaux**, dans le sens où les erreurs sont plus petites et les ordres d'erreurs se rapprochent de  $2k + 2$ . Cette sous-optimalité peut être justifiée par le choix (non optimal) des fonctions de base, qui pourrait être amélioré pour atteindre une meilleure précision.

## 3.2 L'équation de transport 2D

### 3.2.1 Formulation faible du problème

Dans cette section on considère les méthodes de DG pour la résolution de l'équation de transport 2D:

$$\begin{cases} \frac{\partial u}{\partial t} + \nabla \cdot (\mathbf{a}u) = 0, & \text{dans } \Omega \times (0, T) \\ u(t=0) = u_0, & \text{dans } \Omega \end{cases} \quad (3.18)$$

où  $\mathbf{a} = (a_x, a_y)$  est un vecteur fixé,  $\Omega = [0, 1] \times [0, 1]$  et

$$\nabla = \frac{\partial}{\partial x} \hat{\mathbf{i}} + \frac{\partial}{\partial y} \hat{\mathbf{j}}$$

avec  $\hat{\mathbf{i}}$  et  $\hat{\mathbf{j}}$  les vecteurs unitaires dans les directions  $x$  et  $y$  respectivement.

On se concentre sur trois propriétés clés de cette approche:

Premièrement, les méthodes DG sont étroitement liées aux méthodes classiques de volumes finis.

Deuxièmement, l'utilisation de polynômes de haut degré permet aux méthodes DG d'atteindre une grande précision tout en restant hautement parallélisables.

Troisièmement, la viscosité artificielle de la méthode est déterminée par l'amplitude des discontinuités, qui sont elles-mêmes liées aux résidus à l'intérieur des éléments.

Ainsi, lorsque le degré polynomial de la solution approchée augmente, la viscosité artificielle diminue, même en présence de discontinuités.

Pour discrétiser l'équation de transport dans l'espace à l'aide d'une méthode DG, on commence par subdiviser le domaine  $\Omega$  en un maillage triangulaire  $\mathcal{T}_h$ . On cherche ensuite une solution approchée discontinue  $u_h$ , qui appartient à un espace  $V(K)$  pour chaque élément  $K$  du maillage. Cet espace  $V(K)$  n'est pas contraint, bien que le choix le plus courant soit l'espace des polynômes de degré au plus  $k$ , noté  $P_k(K)$ .

Comme cas la section précédente, on commence par une formulation faible donnée par:

$$\int_K (u_h)_t v + \int_K \nabla \cdot (\mathbf{a} u_h) v = 0 \quad (3.19)$$

où  $v$  est une fonction test appartenant à un espace de fonctions  $V(K)$ .

Par une integration par partie de l'équation ci-dessus, on obtient:

$$\int_K (u_h)_t v - \int_K \mathbf{a} u_h \cdot \nabla v + \int_{\partial K} \widehat{\mathbf{a} u_h} \cdot \hat{\mathbf{n}}^k v = 0 \quad (3.20)$$

où  $\hat{\mathbf{n}}^k = (\hat{n}_x, \hat{n}_y)$  est la normale unitaire sortante de la frontière  $\partial_k$  de l'élément  $K$ . En 2D, le rôle des vecteurs normaux entre en jeu. Les différents vecteurs normaux d'un élément triangulaire ainsi que d'un de ses éléments voisins, comme illustré dans la **figure 3.4**.

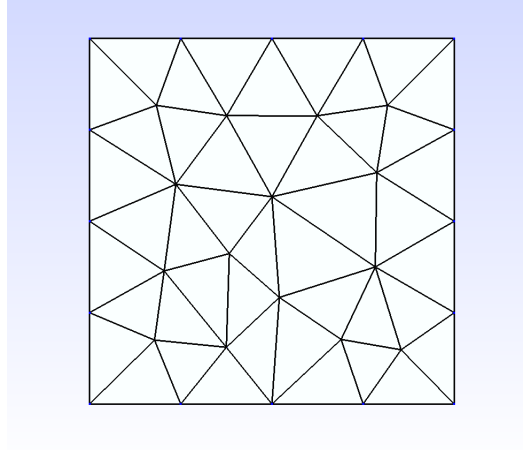


Figure 3.3: Maillage triangulaire du domaine  $\Omega$

Pour compléter la définition de la méthode DG, il est nécessaire de définir le flux numérique  $\widehat{\mathbf{a} u_h}$  sur les bords de chaque élément  $K$ . Pour cela, on procède comme dans la section précédente en obtenant un résultat de stabilité.

Pour ce faire, on multiplie par  $u$  et on intègre sur l'espace et le temps pour obtenir:

$$\frac{1}{2} \int_{\mathbf{R}^2} u^2 + \frac{1}{2} \int_0^T \int_{\mathbf{R}^2} \nabla \cdot \mathbf{a} u^2 = \frac{1}{2} \int_{\mathbf{R}^2} u_0^2$$

Maintenant, imitons la procédure précédente pour la méthode DG considérée. En prenant  $v = u_h$  dans l'équation (3.20), et en additionnant sur les éléments  $K$ , on obtient:

$$\frac{1}{2} \int_K u_h^2(x, y, T) + \frac{1}{2} \int_0^T \int_K \nabla \cdot \mathbf{a} u_h^2(x, y, t) dx dy dt = \frac{1}{2} \int_K u_0^2(x, y) dx dy$$

où

$$\Theta_h(t) = \sum_{K \in \mathcal{T}_h} \left( -\frac{1}{2} \int_K \nabla \cdot (\mathbf{a} u_h)(x, y, t) dx dy + \int_{\partial K} \widehat{\mathbf{a} u_h} \cdot \mathbf{n} u_h(x, y, t) ds \right)$$

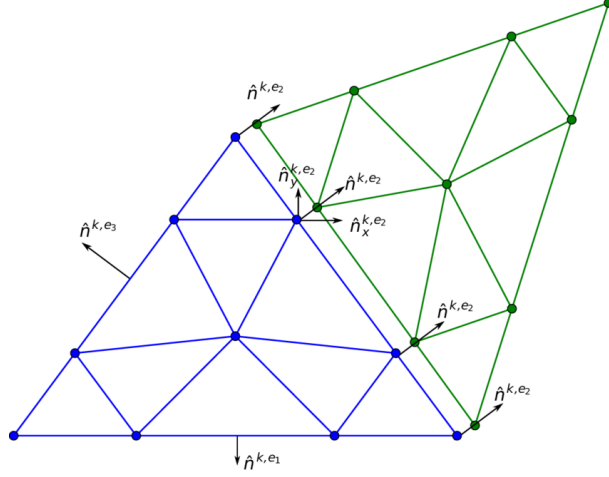


Figure 3.4: Vecteurs normaux du deuxième bord (e2) du  $\mathbf{K}$ -ième élément (bleu) avec son élément voisin (vert) [1]

Ensuite, on définit le flux numérique  $\widehat{\mathbf{a}u_h}$  de sorte que cette quantité  $\Theta_h(t)$  soit positive.

Considérons un point  $(x, y)$  situé sur l'ensemble  $e = \overline{\partial K^+} \cap \overline{\partial K^-}$ . Soient  $n^+$  et  $n^-$  les vecteurs normaux sortants de  $K^+$  et  $K^-$  respectivement en  $(x, y)$ . La valeur de  $u_h^\pm(x, y)$  est définie comme  $\lim_{\epsilon \rightarrow 0} u_h(x - \epsilon n^+, y - \epsilon n^+)$ . On définit le flux comme suit:

$$\widehat{\mathbf{a}u_h^n} = \mathbf{a} \{u_h\} + C \llbracket u_h \rrbracket \quad \text{avec} \quad \llbracket u_h \rrbracket = u_h^- \mathbf{n}^- + u_h^+ \mathbf{n}^+ \quad \text{et} \quad \{u_h\} = \frac{1}{2}(u_h^+ + u_h^-)$$

Enfin soit  $\mathcal{E}_h$  l'ensemble des ensembles  $e = \overline{\partial K^+} \cap \overline{\partial K^-}$  pour tout  $K^+$  et  $K^-$  dans  $\mathcal{T}_h$ . On peut alors écrire:

$$\begin{aligned} \Theta_h(t) &= \sum_{K \in \mathcal{T}_h} \left( -\frac{1}{2} \int_K \nabla \cdot (\mathbf{a}u_h)(x, y, t) dx dy + \int_{\partial K} \widehat{\mathbf{a}u_h} \cdot \mathbf{n} u_h(x, y, t) ds \right) \\ &= \sum_{K \in \mathcal{T}_h} \int_{\partial K} \left( \widehat{\mathbf{a}u_h} \cdot \mathbf{n} u_h - \frac{1}{2} \mathbf{a} u_h^2 \cdot \mathbf{n} \right) ds \\ &= \sum_{e \in \mathcal{E}_h} \int_e \left( \widehat{\mathbf{a}u_h} u_h - \frac{1}{2} \mathbf{a} u_h^2 \right) ds \\ &= \sum_{e \in \mathcal{E}_h} \int_e \left( \widehat{\mathbf{a}u_h} \llbracket u_h \rrbracket - \frac{1}{2} \llbracket \mathbf{a} u_h^2 \rrbracket \right) ds \end{aligned}$$

Comme

$$\frac{1}{2} \llbracket u_h^2 \rrbracket = \llbracket u_h \rrbracket \{u_h\} \quad \text{et} \quad \llbracket \mathbf{a} u_h^2 \rrbracket = \llbracket u_h \rrbracket \{\mathbf{a} u_h\}$$

alors, on obtient:

$$\begin{aligned} \Theta_h(t) &= \sum_{e \in \mathcal{E}_h} \int_e \left( \widehat{\mathbf{a}u_h} \llbracket u_h \rrbracket - \frac{1}{2} \llbracket \mathbf{a} u_h^2 \rrbracket \right) ds \\ &= \sum_{e \in \mathcal{E}_h} \int_e (\widehat{\mathbf{a}u_h} \llbracket u_h \rrbracket - \llbracket u_h \rrbracket \{\mathbf{a} u_h\}) ds \\ &= \sum_{e \in \mathcal{E}_h} \int_e (\widehat{\mathbf{a}u_h} - \mathbf{a} \{u_h\}) \llbracket u_h \rrbracket ds \end{aligned}$$

En prenant

$$\widehat{\mathbf{a}u_h} = \mathbf{a} \{u_h\} + C \llbracket u_h \rrbracket \tag{3.21}$$

on obtient que

$$\Theta_h(t) = \sum_{e \in \mathcal{E}_h} \int_e C \llbracket u_h \rrbracket^2 ds \geq 0$$

pour  $C \geq 0$ .

Ainsi, le choix du flux numérique dans (3.21) garantit la stabilité et la consistance de la méthode DG.

### 3.2.2 Assemblage des matrices locales

On cherche une solution approchée  $u_h^n(x, y) \approx u(x, y, t^n)$  telle que restreinte à l'élément  $K$ , elle s'écrit:

$$u_h^n(x, y) \Big|_K = \alpha^n + \beta^n x + \gamma^n y$$

On fait l'approximation en temps suivante:

$$(u_h)_t = \frac{u_h^{n+1} - u_h^n}{\delta t}$$

On réécrit alors l'équation (3.20) comme suit:

$$\int_K \frac{u_h^{n+1} - u_h^n}{\delta t} v - \int_K \mathbf{a} u_h^n \cdot \nabla v + \int_{\partial K} \widehat{\mathbf{a} u_h^n} \cdot \hat{\mathbf{n}}^k v = 0 \quad \forall v \in V = P_k(K)$$

Pour tout  $v \in V = P_k(K)$ , on a:

$$\int_K \frac{(\alpha^{n+1} - \alpha^n) + (\beta^{n+1} - \beta^n)x + (\gamma^{n+1} - \gamma^n)y}{\delta t} \cdot v - \int_K \mathbf{a} (\alpha^n + \beta^n x + \gamma^n y) \cdot \nabla v + \int_{\partial K} \widehat{\mathbf{a} u_h^n} \cdot \hat{\mathbf{n}}^k \cdot v = 0 \quad (3.22)$$

Ce qui équivaut à:

$$\begin{aligned} \int_K \frac{(\alpha^{n+1} - \alpha^n) + (\beta^{n+1} - \beta^n)x + (\gamma^{n+1} - \gamma^n)y}{\delta t} \cdot v \\ - \int_K a_x (\alpha^n + \beta^n x + \gamma^n y) \frac{\partial v}{\partial x} \\ - \int_K a_y (\alpha^n + \beta^n x + \gamma^n y) \frac{\partial v}{\partial y} \\ + \sum_{\text{Arete} \in \partial K} \int_{\text{Arete}} \left( \widehat{\mathbf{a} u_h^n} \cdot \hat{\mathbf{n}}^k \right) \cdot v = 0 \end{aligned}$$

Supposons un élément  $K$  avec les arêtes  $e_1, e_2, e_3$  qui sont en commun avec les éléments voisins L, M, et N respectivement et posons  $\mathbf{n} = \hat{\mathbf{n}}^k$  pour simplifier la notation. On a alors:

$$\begin{aligned} \int_{e_1} \left( \widehat{\mathbf{a} u_h^n} \right) \cdot \mathbf{n} \cdot v &= \int_{e_1} (\mathbf{a} \{u_h\} + C \llbracket u_h \rrbracket) \cdot \mathbf{n} \cdot v \\ &= \int_{e_1} (\mathbf{a} \{u_h\}) \cdot \mathbf{n} \cdot v + C \int_{e_1} \llbracket u_h \rrbracket \cdot \mathbf{n} \cdot v \\ &= \int_{e_1} \frac{1}{2} \mathbf{a} (u_h^+ + u_h^-) \cdot \mathbf{n} \cdot v + C \int_{e_1} (u_h^- \mathbf{n}^- + u_h^+ \mathbf{n}^+) \cdot \mathbf{n} \cdot v \\ &= \int_{e_1} \frac{1}{2} \mathbf{a} \left( u_h^n \Big|_K + u_h^n \Big|_L \right) \cdot \mathbf{n} \cdot v + C \int_{e_1} \left( u_h^n \Big|_K \mathbf{n}^- + u_h^n \Big|_L \mathbf{n}^+ \right) \cdot \mathbf{n} \cdot v \\ &= \int_{e_1} \frac{1}{2} \mathbf{a} ((\alpha^{K,n} + \beta^{K,n}x + \gamma^{K,n}y) + (\alpha^{L,n} + \beta^{L,n}x + \gamma^{L,n}y)) \cdot \mathbf{n} \cdot v \\ &\quad + C \int_{e_1} ((\alpha^{K,n} + \beta^{K,n}x + \gamma^{K,n}y)\mathbf{n}^- + (\alpha^{L,n} + \beta^{L,n}x + \gamma^{L,n}y)\mathbf{n}^+) \cdot \mathbf{n} \cdot v \\ &= \int_{e_1} \frac{1}{2} \mathbf{a} ((\alpha^{K,n} + \alpha^{L,n}) + (\beta^{K,n} + \beta^{L,n})x + (\gamma^{K,n} + \gamma^{L,n})y) \cdot \mathbf{n} \cdot v \\ &\quad + C \int_{e_1} ((\alpha^{K,n}\mathbf{n}^- + \alpha^{L,n}\mathbf{n}^+) + (\beta^{K,n}\mathbf{n}^- + \beta^{L,n}\mathbf{n}^+)x + (\gamma^{K,n}\mathbf{n}^- + \gamma^{L,n}\mathbf{n}^+)y) \cdot \mathbf{n} \cdot v \end{aligned}$$

Par analogie sur les aretes  $e_2$  (commune avec M) et  $e_3$  (commune avec N), on obtient respectivement:

$$\begin{aligned}\int_{e_2} \left( \widehat{\mathbf{a}u_h^n} \right) \cdot \mathbf{n} \cdot v &= \int_{e_2} \frac{1}{2} \mathbf{a} \left( (\alpha^{K,n} + \alpha^{M,n}) + (\beta^{K,n} + \beta^{M,n})x + (\gamma^{K,n} + \gamma^{M,n})y \right) \cdot \mathbf{n} \cdot v \\ &+ C \int_{e_2} \left( (\alpha^{K,n} \mathbf{n}^- + \alpha^{M,n} \mathbf{n}^+) + (\beta^{K,n} \mathbf{n}^- + \beta^{M,n} \mathbf{n}^+)x + (\gamma^{K,n} \mathbf{n}^- + \gamma^{M,n} \mathbf{n}^+)y \right) \cdot \mathbf{n} \cdot v \\ \int_{e_3} \left( \widehat{\mathbf{a}u_h^n} \right) \cdot \mathbf{n} \cdot v &= \int_{e_3} \frac{1}{2} \mathbf{a} \left( (\alpha^{K,n} + \alpha^{N,n}) + (\beta^{K,n} + \beta^{N,n})x + (\gamma^{K,n} + \gamma^{N,n})y \right) \cdot \mathbf{n} \cdot v \\ &+ C \int_{e_3} \left( (\alpha^{K,n} \mathbf{n}^- + \alpha^{N,n} \mathbf{n}^+) + (\beta^{K,n} \mathbf{n}^- + \beta^{N,n} \mathbf{n}^+)x + (\gamma^{K,n} \mathbf{n}^- + \gamma^{N,n} \mathbf{n}^+)y \right) \cdot \mathbf{n} \cdot v\end{aligned}$$

### 3.2.3 Assemblage d'un systeme lineaire:

Nous considérons une solution approchée  $u_h^n(x, y)$  dans l'élément  $K$  à un instant  $t^n$  donnée par:

$$u_h^n(x, y) = \alpha^{K,n} + \beta^{K,n}x + \gamma^{K,n}y = \sum_{i=1}^3 \alpha_i^{K,n} \phi_i^K(x, y)$$

où  $\phi_i^K(x, y)$  sont les fonctions de base associées à l'élément  $K$  et  $\alpha_i^{K,n}$  sont les coefficients de la solution approchée sur l'élément  $K$  à l'instant  $t^n$ .

Soit  $v = \phi_j^K(x, y)$  pour  $j = 1, 2, 3$  une fonction test dans l'élément  $K$ . En substituant  $u_h^n$  dans l'équation (3.22), on obtient:

$$\int_K \frac{u_h^{n+1} - u_h^n}{\delta t} \phi_j^K(x, y) - \int_K \mathbf{a} \cdot \nabla u_h^n \phi_j^K(x, y) + \sum_{Arete \in \partial K} \int_{Arete} \left( \widehat{\mathbf{a}u_h^n} \cdot \hat{\mathbf{n}}^k \right) \cdot \phi_j^K(x, y) = 0$$

En remplaçant les expressions de  $u_h^n$  dans l'équation ci-dessus, on obtient:

$$\sum_{i=1}^3 \alpha_i^{K,n+1} \int_K \phi_i^K(x, y) \phi_j^K(x, y) - \sum_{i=1}^3 \alpha_i^{K,n} \int_K \mathbf{a} \cdot \nabla \phi_i^K(x, y) \phi_j^K(x, y) + \sum_{Arete \in \partial K} \int_{Arete} \left( \widehat{\mathbf{a}u_h^n} \cdot \hat{\mathbf{n}}^k \right) \cdot \phi_j^K(x, y) = 0$$

#### Matrice de Masse:

La matrice de masse locale sur un élément triangulaire  $K$  est définie par:

$$M = (M_{ij}) \quad \text{où} \quad M_{ij} = \int_K \phi_i^K(x, y) \phi_j^K(x, y) dx dy$$

#### Matrice d'advection

La matrice d'advection est donnée par:

$$S = (S_{ij}) \quad \text{où} \quad S_{ij} = \int_K \mathbf{a} \cdot \nabla \phi_i^K(x, y) \phi_j^K(x, y) dx dy$$

#### Matrice de Flux aux interfaces:

Pour chaque bord  $e_i$  ( $i = 1, 2, 3$ ), on a:

$$F = (F_{e_i}) \quad \text{où} \quad F_{e_i} = \int_{e_i} \left( \widehat{\mathbf{a}u_h^n} \cdot \hat{\mathbf{n}}^k \right) \cdot \phi_j^K(x, y) ds$$

Pour une arete  $e_i$ , ( $i = 1, 2, 3$ ) partagée par  $K$  et  $L$ , On decompose cette expression en deux parties:

**La moyenne des flux:**

$$\int_{e_i} \mathbf{a} \left( u_h^n \Big|_K + u_h^n \Big|_L \right) \cdot \mathbf{n} \cdot \phi_j^K(x, y) ds = \int_{e_i} \mathbf{a} u_h^n \Big|_K \cdot \mathbf{n} \cdot \phi_j^K(x, y) ds + \int_{e_i} \mathbf{a} u_h^n \Big|_L \cdot \mathbf{n} \cdot \phi_j^K(x, y) ds$$

En remplaçant les expressions de  $u_h^n$  sur les éléments  $K$  et  $L$ , on obtient:

$$\int_{e_i} \frac{1}{2} \mathbf{a} \left( \sum_{i=1}^3 \alpha_i^{K,n} \phi_i^K(x, y) + \sum_{i=1}^3 \alpha_i^{L,n} \phi_i^L(x, y) \right) \cdot \mathbf{n} \cdot \phi_j^K(x, y) ds$$

En termes matriciels, cela se traduit par une multiplication vecteurs/matrices:

$$\mathbf{F}_1 = \frac{1}{2} \mathbf{n}^T \mathbf{a} (\mathbf{M}_{KL} \alpha^{K,n} + \mathbf{M}_{LK} \alpha^{L,n})$$

où:

- $\mathbf{M}_{KL}$  est la matrice des intégrales de  $\phi_i^K(x, y) \phi_j^L(x, y)$  sur l'arête  $e_i$ :

$$(\mathbf{M}_{KL})_{ij} = \int_{e_i} \phi_i^K(x, y) \phi_j^L(x, y) ds$$

- $\alpha^{K,n}$  et  $\alpha^{L,n}$  sont les vecteurs des coefficients de  $u_h^n$  sur les éléments  $K$  et  $L$  respectivement
- $\mathbf{n}^T \mathbf{a}$  est le produit scalaire entre le vecteur normal  $\mathbf{n}$  et le vecteur  $\mathbf{a}$ .

**Et le terme de saut:**

$$\int_{e_i} C(u_h^- \mathbf{n}^- + u_h^+ \mathbf{n}^+) \cdot \mathbf{n} \cdot \phi_j^K(x, y) ds$$

Encore une fois, en remplaçant les expressions de  $u_h^-$  et  $u_h^+$ , on obtient:

$$\int_{e_i} C \left( \sum_{i=1}^3 \alpha_i^{K,n} \phi_i^K(x, y) \mathbf{n}^- + \sum_{i=1}^3 \alpha_i^{L,n} \phi_i^L(x, y) \mathbf{n}^+ \right) \cdot \mathbf{n} \cdot \phi_j^K(x, y) ds$$

En termes matriciels:

$$\mathbf{F}_2 = C (\mathbf{n}^- \mathbf{M}_K \alpha^{K,n} + \mathbf{n}^+ \mathbf{M}_L \alpha^{L,n})$$

où  $\mathbf{M}_K$  et  $\mathbf{M}_L$  sont les matrices associées aux éléments  $K$  et  $L$  respectivement:

$$(\mathbf{M}_K)_{ij} = \int_{e_i} \phi_i^K(x, y) \phi_j^K(x, y) ds, \quad (\mathbf{M}_L)_{ij} = \int_{e_i} \phi_i^L(x, y) \phi_j^L(x, y) ds$$

### **Triangle de référence:**

Pour calculer ces matrices  $M$  et  $S$ , on se place tout d'abord dans un triangle «simple», appelé **triangle de référence**  $\hat{K}$ , dont les sommets sont  $\hat{s}_0 = (0, 0)$ ,  $\hat{s}_1 = (1, 0)$ ,  $\hat{s}_2 = (0, 1)$ , ordonné dans le sens trigonométrique. Pour différencier ce triangle d'un triangle du maillage, on lui adjoint un repère  $(\xi, \eta)$ , dit **repère paramétrique**.

Notons  $\hat{\phi}_i \in P^1(\hat{K})$  les trois fonctions de bases associées aux sommets  $\hat{s}_i$ , pour  $i = 0, 1, 2$ , définies par  $\hat{\phi}_i(\hat{s}_j) = \delta_{ij}$ . Ces fonctions  $\hat{\phi}_i$  étant des polynômes de degré 1, elles sont données par:

$$\begin{cases} \hat{\phi}_1(\xi, \eta) = 1 - \xi - \eta \\ \hat{\phi}_2(\xi, \eta) = \xi \\ \hat{\phi}_3(\xi, \eta) = \eta \end{cases} \quad (3.23)$$

**Transformation affine entre triangle de référence et élément du Maillage** Pour passer du triangle de référence à un élément  $K$  du maillage (élément physique), on utilise une transformation affine. Cette transformation permet de cartographier chaque point du triangle de référence  $\hat{K}$  vers un point de l'élément  $K$ .

On considère un élément  $K$  du maillage avec les sommets  $s_0 = (x_0, y_0)$ ,  $s_1 = (x_1, y_1)$ ,  $s_2 = (x_2, y_2)$ . la transformation affine entre  $\hat{K}$  et  $K$  est donnée par:

$$\begin{pmatrix} x \\ y \end{pmatrix} = T \begin{pmatrix} \xi \\ \eta \end{pmatrix} + \begin{pmatrix} x_0 \\ y_0 \end{pmatrix}$$

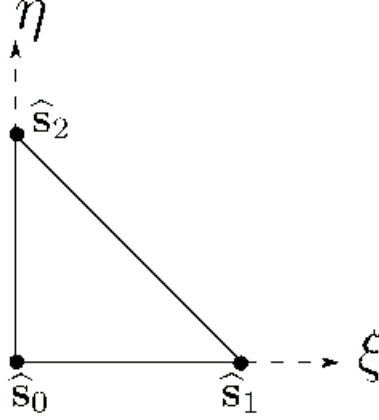


Figure 3.5: Triangle de référence  $\hat{K}$

où  $T$  est la matrice Jacobienne de la transformation, définie par:

$$T = \begin{pmatrix} x_1 - x_0 & x_2 - x_0 \\ y_1 - y_0 & y_2 - y_0 \end{pmatrix}$$

et  $\begin{pmatrix} x_0 \\ y_0 \end{pmatrix}$  est le vecteur de translation.

**Calcul du Jacobien de la transformation et de son déterminant** Le Jacobien de la transformation est important pour calculer les intégrales. Le déterminant  $|J_K| = \det(T)$  mesure le rapport entre l'élément physique  $K$  et le triangle de référence  $\hat{K}$ .

En particulier lorsqu'on intègre une fonction  $f(x, y)$  sur l'élément physique  $K$ , on a:

$$\int_K f(x, y) dx dy = \int_{\hat{K}} f(T(\xi, \eta) + (x_0, y_0)) |J_K| d\xi d\eta$$

**Transformation du gradient des fonctions de base** Lorsqu'on calcule le gradient d'une fonction de base sur l'élément réel, il faut le transformer depuis le triangle de référence. Ainsi, si  $\nabla_{\xi, \eta} \hat{\phi}_j$  est le gradient de la fonction de base  $\hat{\phi}_j$  dans l'élément de référence, alors le gradient dans l'élément réel est donnée par:

$$\nabla_{x, y} \phi_j = T^{-1} \nabla_{\xi, \eta} \hat{\phi}_j$$

Cette relation sera utilisée pour calculer la matrice d'advection.

#### **Intégrations sur le triangle de référence**

Dans le triangle de référence, les matrices de masse local et d'advection sont données respectivement par  $\hat{M} = (\hat{M}_{ij})_{1 \leq i, j \leq 3}$  et  $\hat{S} = (\hat{S}_{ij})_{1 \leq i, j \leq 3}$ , où:

$$\hat{M}_{ij} = |J_K| \int_{\hat{K}} \hat{\phi}_i(\xi, \eta) \hat{\phi}_j(\xi, \eta) d\xi d\eta$$

$$\hat{S}_{ij} = |J_K| \int_{\hat{K}} \mathbf{a} \cdot (T^{-1} \nabla_{\xi, \eta} \hat{\phi}_j(\xi, \eta)) \hat{\phi}_i(\xi, \eta) d\xi d\eta$$

La matrice de masse élémentaire  $\hat{M}$  est donnée par:

$$\hat{M} = \frac{1}{24} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix}$$

en effet, pour  $i = j = 2$ , soit  $\hat{\phi}_i(\xi, \eta) = \hat{\phi}_j(\xi, \eta) = \xi$ . Dans ce cas on a:

$$\hat{M}_{22} = |J_K| \int_{\hat{K}} \xi^2 d\xi d\eta = \int_0^1 \int_0^{1-\xi} \xi^2 d\eta d\xi = \left[ \frac{\xi^3}{3} - \frac{\xi^4}{4} \right]_0^1 = \frac{1}{3} - \frac{1}{4} = \frac{1}{12}$$

Les calculs sont similaires pour  $i = 1$ , et  $i = 3$ . Prenons maintenant le cas  $i \neq j$ , par exemple  $i = 3$  et  $j = 2$ :

$$\hat{M}_{32} = |J_K| \int_{\hat{K}} \hat{\phi}_3(\xi, \eta) \hat{\phi}_2(\xi, \eta) d\xi d\eta = \int_0^1 \left( \int_0^{1-\xi} \eta d\eta \right) \xi d\xi = \frac{1}{2} \int_0^1 \xi(1-\xi)^2 d\xi = \frac{1}{2} \left[ \frac{1}{2} - \frac{2}{3} + \frac{1}{4} \right] = \frac{1}{24}$$

Les calculs sont similaires pour les autres coefficients  $M_{11}, M_{12}, M_{13}, M_{21}, M_{23}$  et  $M_{31}$ .

Système linéaire:  $A\mathbf{U}^{n+1} = F$

En assemblant les matrices élémentaires, on obtient le système linéaire suivant:

$$A\mathbf{U}^{n+1} = F$$

où  $\mathbf{U}^{n+1} = [\alpha^{K,n+1}, \beta^{K,n+1}, \gamma^{K,n+1}]^T$  est le vecteur des coefficients de la solution approchée à l'instant  $t^{n+1}$ ,  $A$  est la matrice de taille  $3N \times 3N$  et  $F$  est le vecteur de taille  $3N$  donné par:

$$F = \begin{pmatrix} F_1 \\ F_2 \\ F_3 \end{pmatrix}$$

où  $F_i$  est le vecteur de taille  $N$  associé à l'élément  $i$ :

$$F_i = \begin{pmatrix} F_{e_1} \\ F_{e_2} \\ F_{e_3} \end{pmatrix}$$

et  $N$  est le nombre total d'éléments dans le maillage.

La matrice  $A$  est donnée par:

$$A = \begin{pmatrix} M_1 + S_1 & -F_1 & 0 \\ -F_1 & M_2 + S_2 & -F_2 \\ 0 & -F_2 & M_3 + S_3 \end{pmatrix}$$

où  $M_i$  et  $S_i$  sont les matrices de masse et d'advection associées à l'élément  $i$  respectivement.



## Chapter 4

# Implémentation de la méthode DG pour l'équation de transport 2D

Pour illustrer la méthode DG pour l'équation de transport 2D, voici quelques extraits de codes que nous avons eu à faire pour la résolution du problème de transports (3.18)

```
1      class dG2DimTransport(AbstractProblem):
2          def __init__(self, T=1.0, a=(1., 1.), k=1, N=20):
3              self.steady = False
4              super().__init__([(0., 1.), (0., 1.)], N)
5              self.a = a
6
7              self.phi_hat = [
8                  lambda xi, eta: 1 - xi - eta,
9                  lambda xi, eta: xi,
10                 lambda xi, eta: eta
11             ]
12
13             self.dphi_xi = [
14                 lambda xi, eta: -1,
15                 lambda xi, eta: 1,
16                 lambda xi, eta: 0
17             ]
18
19             self.dphi_eta = [
20                 lambda xi, eta: -1,
21                 lambda xi, eta: 0,
22                 lambda xi, eta: 1
23             ]
24
25             self.p = len(self.phi_hat)
26             self.dofs = self.mesh.ne * self.p
27             self.T = T # Temps final
28
29             self.u = np.ones((self.mesh.ne, self.p)) #
30             self.u0 = lambda x, y: 1 # Condition initiale
31             self.ue = lambda t, x, y: self.u0(x - t * self.a[0], y - t * self.a[1])
32
```

Ceci permet d'initialiser la classe dG2DimTransport qui hérite de la classe AbstractProblem.

Ensuite les quelques méthodes importantes pour assembler les matrices élémentaires M, S, et F.

```
1      def eval_at_point(self, t, x, y):
2          n = self.mesh.get_index((x, y))
3
4          x0, y0 = self.mesh.get_vertex(n, 0)
5          x1, y1 = self.mesh.get_vertex(n, 1)
6          x2, y2 = self.mesh.get_vertex(n, 2)
7
```

```

8         T, B = self.transformer((x0, y0), (x1, y1), (x2, y2))
9         xi_eta = np.linalg.inv(T) @ (np.array([x, y]) - B)
10
11         phi_values = np.array([phi(xi_eta[0], xi_eta[1]) for phi in self.phi_hat])
12         u_t = self.u[n, :]
13
14         ut = u_t @ phi_values
15         return ut
16
17     def transformer(self, X0, X1, X2):
18         x0, y0 = X0
19         x1, y1 = X1
20         x2, y2 = X2
21
22         T = np.array([[x1 - x0, x2 - x0], [y1 - y0, y2 - y0]])
23         B = np.array([x0, y0])
24
25         return T, B
26
27     def transform_to_global(self, T, B, xi, eta):
28         global_coords = T @ np.array([xi, eta]) + B
29         return global_coords[0], global_coords[1]
30
31     def integrand_M(self, i, j, T, B, xi, eta):
32         return self.phi_hat[i](xi, eta) * self.phi_hat[j](xi, eta)
33
34     def assemble_Sa_local(self, i, j, T, a, xi, eta):
35         T_inv_T = np.linalg.inv(T).T
36         integrand_Sa = (a @ (T_inv_T @ np.array([self.dphi_xi[j](xi, eta), self.dphi_eta[j](xi,
37 eta)]))) * self.phi_hat[i](xi, eta)
38         return integrand_Sa
39
40     def assembly_local(self, C=0.5, element=0):
41         x0, y0 = self.mesh.get_vertex(element, 0)
42         x1, y1 = self.mesh.get_vertex(element, 1)
43         x2, y2 = self.mesh.get_vertex(element, 2)
44
45         M_local = np.zeros((self.p, self.p))
46         Sa_local = np.zeros((self.p, self.p))
47         F_e = np.zeros(self.p)
48
49         X0, X1, X2 = (x0, y0), (x1, y1), (x2, y2)
50         T, B = self.transformer(X0, X1, X2)
51         detT = np.abs(np.linalg.det(T))
52
53         for i in range(self.p):
54             for j in range(self.p):
55                 M_local[i, j] = integrate.dblquad(lambda xi, eta: self.integrand_M(i, j, T, B,
56 xi, eta) * detT,
57                                                     0, 1, lambda xi: 0, lambda xi: 1 - xi)[0]
58                 Sa_local[i, j] = integrate.dblquad(lambda xi, eta: self.assemble_Sa_local(i, j,
59 T, self.a, xi, eta) * detT,
60                                                     0, 1, lambda xi: 0, lambda xi: 1 - xi)[0]
61
62         voisins = self.mesh.voisins[element]
63         for index, voisin in enumerate(voisins):
64             if voisin == -1:
65                 continue
66             else:
67                 x0_v, y0_v = self.mesh.get_vertex(voisin, 0)
68                 x1_v, y1_v = self.mesh.get_vertex(voisin, 1)
69                 x2_v, y2_v = self.mesh.get_vertex(voisin, 2)
70
71                 T_v, B_v = self.transformer((x0_v, y0_v), (x1_v, y1_v), (x2_v, y2_v))
72
73                 # Normalisation de la normale
74                 nx, ny = y1 - y0, x0 - x1
75                 norme = math.sqrt(nx**2 + ny**2)

```

```

73         nx, ny = nx / norme, ny / norme
74
75         for i in range(self.p):
76             integrand_flux = lambda s: (0.5 * (self.a[0] * nx + self.a[1] * ny) * self.
phi_hat[i](s, 0) +
77                                     C * (nx**2 + ny**2) * self.phi_hat[i](s, 0))
78             F_e[i] += integrate.quad(integrand_flux, 0, 1, limit=100)[0]
79
80         return M_local, Sa_local, F_e
81
82

```

Une assemblage des matrices élémentaires pour obtenir la matrice globale  $A$  et le vecteur  $F$ .

```

1     def assemblage_global(self, C=0.5):
2         n_dofs = self.mesh.ne * self.p
3         A = np.zeros((n_dofs, n_dofs))
4         S = np.zeros((n_dofs, n_dofs))
5         F = np.zeros(n_dofs)
6
7         for K in range(self.mesh.ne):
8             MK, SK, Fe = self.assembly_local(C, K)
9             for i in range(self.p):
10                 dof_i = K * self.p + i
11                 F[dof_i] += Fe[i]
12                 for j in range(self.p):
13                     dof_j = K * self.p + j
14                     A[dof_i, dof_j] += MK[i, j]
15                     S[dof_i, dof_j] += SK[i, j]
16
17         return A, S, F
18
19
20

```

Voir le code complet à l'adresse suivante:

[https://gitlab.math.unistra.fr/aghili/dg/-/blob/dieng-main-patch-68158/EquationTransport.py?ref\\_type=heads](https://gitlab.math.unistra.fr/aghili/dg/-/blob/dieng-main-patch-68158/EquationTransport.py?ref_type=heads)

et le repo gitlab complet se trouve à l'adresse suivante:

[https://gitlab.math.unistra.fr/aghili/dg/-/tree/dieng-main-patch-68158?ref\\_type=heads](https://gitlab.math.unistra.fr/aghili/dg/-/tree/dieng-main-patch-68158?ref_type=heads)

### Quelques tests unitaires pour vérifier les méthodes:

Afin de vérifier ces méthodes implémentées, nous avons fait quelques tests unitaires avec **pytest** pour vérifier les assemblages des matrices élémentaires.

```

1     import pytest
2     import numpy as np
3     from scipy import integrate
4     from EquationTransport import dG2DimTransport
5
6     expected_matrix = (1/24) * np.array([[2,1, 1],
7                                           [1, 2, 1],
8                                           [1, 1, 2]])
9
10    def setup_test_problem():
11        T = 1.0
12        a = (0.25, 0.25)
13        k = 1
14        N = 20
15        problem = dG2DimTransport(T=T, a=a, k=k, N=N)
16        return problem
17
18    def test_integrand_M():
19        problem = setup_test_problem()
20        element = 0
21        x0, y0 = problem.mesh.get_vertex(element, 0)

```

```

22     x1, y1 = problem.mesh.get_vertex(element, 1)
23     x2, y2 = problem.mesh.get_vertex(element, 2)
24     X0, X1, X2 = (x0, y0), (x1, y1), (x2, y2)
25     A, B = problem.transformer(X0, X1, X2)
26
27     for i in range(problem.p):
28         for j in range(problem.p):
29             result = integrate.dblquad(lambda x, y: problem.integrand_M(i, j, A, B, x, y),
30             0, 1, lambda x: 0, lambda x: 1 - x)[0]
31             assert np.isclose(result, expected_matrix[i, j]), f"Integrand M({i},{j}) = {
32             result}, expected {expected_matrix[i, j]}"
33
34     if __name__ == "__main__":
35         pytest.main()

```

Dans la sortie de la commande **pytest**, on obtient le résultat suivant qui indique bien que cette assemblage de la matrice de masse est correcte.

```

1      ===== test session starts =====
2      platform linux -- Python 3.10.9, pytest-7.1.2, pluggy-1.0.0
3      rootdir: /home/mai-pret/Bureau/stage/StageDG/dg
4      plugins: anyio-3.5.0
5      collected 1 item
6
7      test_dG2DimTransport_integrals.py . [100%]
8
9      ===== 1 passed in 0.80s =====
10

```

## Chapter 5

# Conclusion et perspectives

### Conclusion

Ce travail a permis d'explorer en profondeur les aspects théoriques et pratiques liés à l'application de la méthode de Galerkin discontinue (DG) pour la résolution d'équations différentielles ordinaires et d'équations de transport en 2D.

Ce travail a permis d'explorer en profondeur les aspects théoriques et pratiques liés à l'application de la méthode de Galerkin discontinue (DG) pour la résolution

Du point de vue de l'implémentation, nous avons d'abord traité le cas en 1D, ce qui nous a permis de reproduire les ordres d'erreurs théoriques attendus. Cette étape a été déterminante pour valider notre approche et ajuster les paramètres de la méthode DG. Toutefois, il a été constaté que le choix des fonctions de base n'était pas optimal, ce qui a conduit à des résultats sous-optimaux. Une amélioration de ce choix pourrait potentiellement augmenter la précision de la méthode.

Le passage au cas 2D a introduit de nouvelles difficultés, notamment liées à la complexité accrue des flux aux interfaces et à l'assemblage des matrices. Ces défis ont nécessité une adaptation des techniques utilisées en 1D pour les rendre compatibles avec les exigences d'une discrétisation en deux dimensions.

### Perspectives

Les perspectives de ce travail sont multiples:

- **Amélioration en 1D:** Un des premiers objectifs serait d'optimiser le choix des fonctions de base afin de réduire les erreurs et d'atteindre des ordres de convergence plus proches des prévisions théoriques.
- **Finalisation en 2D:** L'application de la méthode DG à un cas simple en 2D, tel que l'advection scalaire constante, constitue une prochaine étape logique. Ce test permettra de vérifier la robustesse de l'approche développée et d'identifier les ajustements nécessaires avant de passer à des problèmes plus complexes.

# Bibliography

- [1] Bharat Tripathi.  
Discontinuous Galerkin Method for Propagation of Acoustical Shock Waves in Complex Geometry. Acoustics [physics.class-ph]. Université Pierre et Marie Curie - Paris VI, 2015. English. ffNNT : 2015PA066344ff. fftel-01297050f
- [2] Bernardo Cockburn.  
Discontinuous Galerkin methods. Journal of Applied Mathematics and Mechanics / Zeitschrift für Angewandte Mathematik und Mechanik, 2003, 83 (11), pp.731-754. 10.1002/zamm.200310088. hal-01352444
- [3] [https://fr.wikipedia.org/wiki/M%C3%A9thode\\_de\\_Galerkine\\_discontinue](https://fr.wikipedia.org/wiki/M%C3%A9thode_de_Galerkine_discontinue)
- [4] Delfour, M., W. Hager, and F. Trochu.  
“Discontinuous Galerkin Methods for Ordinary Differential Equations. ” Mathematics of Computation 36, no. 154 (1981): 455–73. <https://doi.org/10.2307/2007652>
- [5] Jean-Baptiste Clément, Mehmet Ersoy, Frederic Golay, Damien Sous.  
Discontinuous Galerkin Method for steady-state Richards Equation. Topical Problems of Fluid Mechanics 2019, Feb 2019, Prague, Czech Republic. pp.53-62, 10.14311/TPFM.2019.008. hal-02075109
- [6] Bernardo Cockburn.  
An Introduction to the Discontinuous Galerkin Method for Convection-Dominated Problems School of Mathematics, University of Minnesota,