

Erreur de troncature

1) PROBLÉMATIQUE

• Nous cherchons à connaître les solutions $t \mapsto u(t)$ de l'équation différentielle

$$(1) \quad \frac{du}{dt} = f(u), \quad t > 0$$

au moins de manière approchée aux instants

$$(2) \quad t^k = k \Delta t, \quad k \in \mathbb{N}, \quad \Delta t > 0 \text{ fixé.}$$

Pour cela, nous avons à notre disposition les quatre schémas proposés au chapitre précédent (et il en existe aussi beaucoup d'autres !) qui fournissent une valeur approchée u^k à l'instant t^k . Nous notons $u_{\Delta t}^k$ la valeur numérique approchée avec le schéma numérique pour un pas de temps Δt . Nous cherchons à comparer cette valeur approchée $u_{\Delta t}^k$ à la valeur exacte $u(t^k)$.

• Plus précisément, fixons $T > 0$ relativement “grand” devant le temps caractéristique τ de variation des solutions de l'équation (1). L'écart entre la solution exacte et cette valeur approchée est notée $\epsilon_{\Delta t}^k$:

$$(3) \quad \epsilon_{\Delta t}^k = u_{\Delta t}^k - u(t^k), \quad k \Delta t \leq T.$$

Nous voulons que cet écart reste “petit”, ce pour **tous** les instants discrets $k \Delta t$ jusqu'au temps T . On introduit donc l'**erreur** $\delta(\Delta t)$:

$$(4) \quad \delta(\Delta t) = \sup_{0 \leq k \Delta t \leq T} |\epsilon_{\Delta t}^k|$$

et on souhaite que $\delta(\Delta t)$ soit “petit” pour Δt “assez petit”. En d'autres termes, on souhaite que $\delta(\Delta t)$ **tende vers zéro** lorsque le pas de temps Δt tend vers zéro.

• Comment aborder l'étude de ce problème ? En effet, si on sait *a priori* calculer $u_{\Delta t}^k$ avec l'algorithme associé au schéma numérique (et sa mise en œuvre sur ordinateur), on ignore tout de $u(k\Delta t)$, valeur de la solution exacte de l'équation (1) à l'instant discret t^k . L'astuce consiste à **renverser les rôles**, ce qui conduit à la notion d'erreur de troncature. Nous verrons en fin de chapitre que si le schéma est **stable**, l'erreur de troncature donne une bonne estimation de l'erreur $\delta(\Delta t)$.

2) SCHÉMA D'EULER EXPLICITE

- Nous écrivons ce schéma

$$(5) \quad \frac{1}{\Delta t} (u^{k+1} - u^k) - f(u^k) = 0.$$

A partir d'une valeur u^k à l'instant discret $k \Delta t$, le schéma d'Euler explicite reconstitue une valeur u^{k+1} qui veut être une valeur approchée de $u((k+1)\Delta t)$. Imaginons qu'on parte de la valeur exacte $u(k\Delta t)$ *i.e.* qu'on suppose $u^k = u(k\Delta t)$. Le schéma numérique (5) calcule une valeur u^{k+1} **différente** de $u((k+1)\Delta t)$, puisque le schéma numérique (5) n'est qu'une **approximation** de la relation exacte

$$(6) \quad \frac{1}{\Delta t} (u((k+1)\Delta t) - u(k\Delta t)) - \frac{1}{\Delta t} \int_{k\Delta t}^{(k+1)\Delta t} f(u(t)) dt = 0.$$

- Nous injectons $u^k = u(t^k)$, valeur exacte dans le schéma (5). Nous avons alors $u^{k+1} = u(k\Delta t) + \Delta t f(u(k\Delta t)) \neq 0$, soit en d'autres termes

$$(7) \quad \frac{1}{\Delta t} (u((k+1)\Delta t) - u(k\Delta t)) - \frac{1}{\Delta t} f(u(k\Delta t)) \neq 0.$$

La solution exacte de l'équation (1) n'est en général **pas** solution du schéma numérique (5). C'est cet écart qu'on appelle "erreur de troncature". Nous définissons cette erreur de troncature \mathcal{T} autour d'un temps $t > 0$ arbitraire, d'un pas de temps de temps $\Delta t > 0$ quelconque également et pour une **solution** $u(\bullet)$ de l'équation différentielle (1). Nous posons pour le schéma d'Euler explicite :

$$(8) \quad \mathcal{T}(\Delta t, t; u(\bullet)) \equiv \frac{1}{\Delta t} (u(t + \Delta t) - u(t)) - f(u(t)).$$

Cette erreur de troncature mesure "en quoi le schéma est mal vérifié pour une solution exacte de l'équation à résoudre". Si elle est grande, on a peu d'espoir. Si elle est "petite", très petite pour Δt assez petit, on imagine que le schéma "simule bien" l'équation (1) et qu'en conséquence c'est l'**erreur** $\delta(\Delta t)$ qui sera petite !

- Même si on ne connaît pas la solution $u(\bullet)$ de l'équation (1), on peut faire le **développement limité** de l'erreur de troncature $\mathcal{T}(\Delta t, t; u(\bullet))$ lorsque le pas de temps Δt tend vers zéro. Grâce à la formule de Taylor

$$(9) \quad u(t + \Delta t) = u(t) + \Delta t \frac{du}{dt}(t) + \frac{1}{2} \Delta t^2 \frac{d^2u}{dt^2}(t) + O(\Delta t^3).$$

Si $u(\bullet)$ est assez régulière, nous tirons de (8) :

$$(10) \quad \mathcal{T}(\Delta t, t; u(\bullet)) = \left[\frac{du}{dt}(t) - f(u(t)) \right] + \frac{\Delta t}{2} \frac{d^2u}{dt^2}(t) + O(\Delta t^2),$$

ce qui constitue un développement limité de l'erreur de troncature. Puisque $u(\bullet)$ est solution du système dynamique (1), le premier terme du membre de droite de la relation (10) est **nul**. Nous en déduisons, puisqu'*a priori* d^2u/dt^2 est non nul :

$$(11) \quad \mathcal{T}(\Delta t, t; u(\bullet)) = O(\Delta t).$$

- Nous avons mis en évidence l'ordre asymptotique de convergence de l'erreur de troncature du schéma d'Euler explicite. Il est de la forme $O(\Delta t^1)$, avec la valeur "unité" comme exposant de Δt . Pour cette raison, on dit que le schéma d'Euler explicite est **d'ordre 1**.

3) SCHÉMA D'EULER IMPLICITE

- Nous procédons pour le schéma d'Euler implicite

$$(12) \quad \frac{1}{\Delta t} (u^{k+1} - u^k) - f(u^{k+1}) = 0$$

comme pour le schéma d'Euler explicite. Nous introduisons une **solution** $u(\bullet)$ de l'équation (1), et injectons les valeurs $u(t)$ et $u(t + \Delta t)$ dans l'expression (12) du schéma, remplaçant le temps discret $t^k = k \Delta t$ par le temps continu t :

$$(13) \quad \mathcal{T}(\Delta t, t; u(\bullet)) = \frac{1}{\Delta t} (u(t + \Delta t) - u(t)) - f(u(t + \Delta t)).$$

Nous définissons ainsi l'erreur de troncature du schéma d'Euler implicite. Afin de connaître son comportement asymptotique pour Δt tendant vers zéro, nous avons besoin de développer $f(u(t + \Delta t))$. Nous l'effectuons au troisième ordre de précision.

Lemme 1. Développement limité.

Pour $t \mapsto u(t)$ régulière et $u \mapsto f(u)$ régulière, nous avons

$$(14) \quad \begin{cases} f(u(t + \Delta t)) = f(u(t)) + \Delta t \frac{du}{dt} f'(u(t)) + \\ + \frac{\Delta t^2}{2} \frac{d^2u}{dt^2} f'(u(t)) + \frac{1}{2} \left(\Delta t \frac{du}{dt} \right)^2 f''(u(t)) + O(\Delta t^3). \end{cases}$$

- Preuve du lemme 1.

Nous écrivons la formule de Taylor pour $f(u(t) + v)$, pour un infiniment petit v *a priori* arbitraire :

$$(15) \quad f(u(t) + v) = f(u(t)) + v f'(u(t)) + \frac{v^2}{2} f''(u(t)) + O(v^3)$$

puis nous particularisons v compte tenu du développement donné en (9) :

$$(16) \quad v = \Delta t \frac{du}{dt}(t) + \frac{1}{2} \Delta t^2 \frac{d^2u}{dt^2}(t) + O(\Delta t^3).$$

On a donc

$$(17) \quad \frac{1}{2} v^2 = \frac{1}{2} \left(\Delta t \frac{du}{dt} \right)^2 + O(\Delta t^3)$$

$$(18) \quad O(v^3) = O(\Delta t^3).$$

On injecte les relations (16) à (18) au sein du développement (15). On remarque que $f(u(t + \Delta t)) = f(u(t) + v + O(\Delta t^3)) = f(u(t) + v) + O(\Delta t^3)$, donc

$$\begin{aligned} f(u(t + \Delta t)) &= f(u(t)) + \left(\Delta t \frac{du}{dt} + \frac{1}{2} \Delta t^2 \frac{d^2u}{dt^2} \right) f'(u(t)) + \\ &\quad + \frac{1}{2} \left(\Delta t \frac{du}{dt} \right)^2 f''(u(t)) + O(\Delta t^3), \end{aligned}$$

ce qui constitue exactement le développement (14) annoncé. \square

• Avec la définition (13) de l'erreur de troncature et le développement (14), on a

$$\begin{aligned} \mathcal{T}(\Delta t, t; u(\bullet)) &= \frac{du}{dt} + \frac{\Delta t}{2} \frac{d^2u}{dt^2} - \left[f(u(t)) + \Delta t \frac{du}{dt} f'(u(t)) \right] + O(\Delta t^2) \\ &= \left(\frac{du}{dt} - f(u(t)) \right) + \Delta t \left(\frac{1}{2} \frac{d^2u}{dt^2} - \frac{du}{dt} f'(u(t)) \right) + O(\Delta t^2). \end{aligned}$$

Si $u(\bullet)$ est solution de l'équation différentielle (1), on a $\frac{du}{dt} = f(u(t))$ et par dérivation par rapport au temps de cette identité :

$$\frac{d^2u}{dt^2} = \frac{d}{dt} \left(f(u(t)) \right) = f'(u(t)) \bullet \frac{du}{dt}$$

donc le développement de l'erreur de troncature s'écrit

$$(19) \quad \mathcal{T}(\Delta t, t; u(\bullet)) = -\frac{\Delta t}{2} \frac{d^2u}{dt^2} + O(\Delta t^2)$$

ce qui montre que **le schéma d'Euler rétrograde est d'ordre 1.**

4) SCHÉMA DE CRANK-NICOLSON

- C'est en quelque sorte la "moyenne" entre les deux schémas d'Euler (5) et (12) :

$$(20) \quad \frac{1}{\Delta t} (u^{k+1} - u^k) - \frac{1}{2} [f(u^k) + f(u^{k+1})] = 0.$$

De manière analogue aux deux autres schémas, on introduit une solution de l'équation différentielle (1), on remplace u^k par $u(t)$ et u^{k+1} par $u(t + \Delta t)$ dans l'expression (20) du schéma, et le résultat obtenu **définit** l'erreur de troncature :

$$(21) \quad \mathcal{T}(\Delta t, t; u(\bullet)) \equiv \frac{u(t + \Delta t) - u(t)}{\Delta t} - \frac{1}{2} [f(u(t)) + f(u(t + \Delta t))].$$

- Le développement limité de l'erreur de troncature (21) du schéma de Crank-Nicolson s'obtient en rapprochant les développements (9) et (14). Nous obtenons :

$$(22) \quad \begin{aligned} \mathcal{T}(\Delta t, t; u(\bullet)) &= \frac{du}{dt} + \frac{\Delta t}{2} \frac{d^2u}{dt^2} + \frac{\Delta t^2}{6} \frac{d^3u}{dt^3} + O(\Delta t^3) \\ &\quad - \frac{1}{2} \left[f(u(t)) + f(u(t)) + \Delta t \frac{du}{dt} f'(u(t)) + \frac{\Delta t^2}{2} \frac{d^2u}{dt^2} f'(u(t)) \right. \\ &\quad \left. + \frac{1}{2} \left(\Delta t \frac{du}{dt} \right)^2 f''(u(t)) \right] + O(\Delta t^3) \\ \left\{ \begin{aligned} \mathcal{T}(\Delta t, t; u(\bullet)) &= \left[\frac{du}{dt} - f(u(t)) \right] + \frac{\Delta t}{2} \left[\frac{d^2u}{dt^2} - \frac{du}{dt} f'(u(t)) \right] + \\ &+ \Delta t^2 \left[\frac{1}{6} \frac{d^3u}{dt^3} - \frac{1}{4} \frac{d^2u}{dt^2} f'(u(t)) - \frac{1}{4} \left(\frac{du}{dt} \right)^2 f''(u(t)) \right] + O(\Delta t^3). \end{aligned} \right. \end{aligned}$$

- Nous constatons que le terme constant du développement (22) est nul car $u(\bullet)$ est solution de l'équation (1). Quand on dérive une fois cette relation, nous avons :

$$(23) \quad \frac{d^2u}{dt^2} = f'(u(t)) \bullet \frac{du}{dt},$$

ce qui montre que le coefficient du terme en Δt dans le développement (22) est nul, donc que le schéma de Crank-Nicolson est au moins d'ordre deux. Par dérivation en temps de la relation (23), nous avons

$$(24) \quad \frac{d^3u}{dt^3} = f''(u(t)) \left(\frac{du}{dt} \right)^2 + f'(u(t)) \frac{d^2u}{dt^2}$$

donc le coefficient de Δt^2 dans le développement (22) vaut $\left(\frac{1}{6} - \frac{1}{4} \right) \frac{d^3u}{dt^3} = -\frac{1}{12} \frac{d^3u}{dt^3}$; il est en général **non nul**. Nous retenons

$$(25) \quad \mathcal{T}(\Delta t, t; u(\bullet)) = -\frac{1}{12} \frac{d^3 u}{dt^3} \Delta t^2 + O(\Delta t^3), \quad \text{Crank Nicolson}$$

et l'erreur de troncature du schéma de Crank-Nicolson tend vers zéro comme $O(\Delta t^2)$. On dit pour cette raison qu'il est d'ordre deux.

5) SCHÉMA DE HEUN

- On rappelle que ce schéma prend la forme

$$(26) \quad \frac{1}{\Delta t} (u^{k+1} - u^k) - \frac{1}{2} \left[f(u^k) + f(u^k + \Delta t f(u^k)) \right] = 0.$$

L'erreur de troncature se définit comme dans les cas précédents

$$(27) \quad \left\{ \begin{array}{l} \mathcal{T}(\Delta t, t; u(\bullet)) \equiv \frac{1}{\Delta t} (u(t^{k+1}) - u(t^k)) \\ \quad - \frac{1}{2} \left[f(u(t^k)) + f(u(t^k) + \Delta t f(u(t^k))) \right]. \end{array} \right.$$

- Pour déterminer l'ordre de ce schéma, on développe l'erreur de troncature (27), sans oublier les relations (1), (23) et (24) qui lui sont dérivées. De manière analogue au lemme 1, on a la relation (15), à appliquer avec $v = \Delta t f(u(t))$ cette fois, pour lequel nous avons $\frac{v^2}{2} = \frac{\Delta t^2}{2} (f(u(t)))^2$ et $O(v^3) = O(\Delta t^3)$. Il vient

$$(28) \quad \left\{ \begin{array}{l} f(u(t^k) + \Delta t f(u(t^k))) = f(u(t)) + \Delta t f'(u(t)) \bullet f(u(t)) \\ \quad + \frac{1}{2} \Delta t^2 f''(u(t)) \bullet f(u(t))^2 + O(\Delta t^3). \end{array} \right.$$

On reporte cette expression au second membre de la relation (27) et on tient compte du développement (9). Il vient

$$(29) \quad \left\{ \begin{array}{l} \mathcal{T}(\Delta t, t; u(\bullet)) = \frac{du}{dt} + \frac{\Delta t}{2} \frac{d^2 u}{dt^2} + O(\Delta t^3) - f'(u(t)) \\ \quad - \frac{1}{2} \left[\Delta t f'(u(t)) f(u(t)) + \frac{\Delta t^2}{2} f''(u(t)) f(u(t))^2 \right] + O(\Delta t^3) \\ \mathcal{T}(\Delta t, t; u(\bullet)) = \frac{du}{dt} - f(u(t)) + \frac{\Delta t}{2} \left[\frac{d^2 u}{dt^2} - f'(u(t)) f(u(t)) \right] \\ \quad + \Delta t^2 \left(\frac{1}{6} \frac{d^3 u}{dt^3} - \frac{1}{4} f''(u(t)) f(u(t))^2 \right) + O(\Delta t^3). \end{array} \right.$$

Le terme constant dans le développement (29) est identiquement nul compte tenu de la relation (1). Le terme en Δt l'est également compte tenu des relations (23) et (1). On a aussi suite à (24) :

$$(30) \quad \frac{d^3u}{dt^3} = f''(u(t)) (f(u(t)))^2 + (f'(u(t)))^2 f(u(t))$$

donc

$$(31) \quad \left\{ \begin{array}{l} \mathcal{T}(\Delta t, t; u(\bullet)) = \left[-\frac{1}{12} f''(u(t)) f^2(u(t)) + \right. \\ \left. \frac{1}{6} (f'(u(t)))^2 f(u(t)) \right] \Delta t^2 + O(\Delta t^3) \end{array} \right.$$

et le coefficient du terme d'ordre deux dans le développement (31) est en général **non nul**. Le schéma de Heun est d'ordre **deux**.

6) DÉFINITION GÉNÉRALE

• Dans le cas d'un schéma général qui s'écrit par exemple sous la forme

$$(32) \quad \frac{1}{\Delta t} (u^{k+1} - u^k) - \Phi(u^k, u^{k+1}, f(u^k), f(u^{k+1})) = 0,$$

nous définissons l'erreur de troncature $\mathcal{T}(\Delta t, t; u(\bullet))$ par la relation

$$(33) \quad \left\{ \begin{array}{l} \mathcal{T}(\Delta t, t; u(\bullet)) \equiv \frac{1}{\Delta t} (u(t + \Delta t) - u(t)) \\ - \Phi(u(t), u(t + \Delta t), f(u(t)), f(u(t + \Delta t))) \end{array} \right.$$

pour une solution $u(\bullet)$ de l'équation (1).

• On dit que le schéma (32) est **consistant** avec l'équation (1) lorsque l'erreur de troncature \mathcal{T} définie à la relation (33) **tend vers zéro** si Δt tend vers zéro. On peut vérifier sans peine que c'est le cas si et seulement si la fonction Φ est continue et satisfait à

$$(34) \quad \Phi(u, u, f(u), f(u)) = f(u), \quad \forall u.$$

On dit que le schéma (32) est **d'ordre p** si l'erreur de troncature (33) admet le développement

$$(35) \quad \mathcal{T}(\Delta t, t; u(\bullet)) = O(\Delta t^p)$$

lorsque Δt tend vers zéro.

7) ERREUR ET ERREUR DE TRONCATURE

• Nous illustrons sur un exemple un résultat général qui énonce qu'un schéma **stable** et **consistant** est alors **convergent**. Nous fixons $\lambda > 0$ et étudions l'équation modèle

$$(36) \quad \frac{du}{dt} + \lambda u = 0, \quad t > 0.$$

Nous la discrétisons avec un schéma d'Euler explicite pour un pas de temps Δt :

$$(37) \quad \frac{1}{\Delta t} (u_{\Delta t}^{k+1} - u_{\Delta t}^k) + \lambda u_{\Delta t}^k = 0.$$

L'erreur de troncature est définie par

$$(38) \quad \mathcal{T}(\Delta t, t; u(\bullet)) \equiv \frac{1}{\Delta t} (u((k+1)\Delta t) - u(k\Delta t)) - \lambda u(k\Delta t)$$

et on a vu à la relation (10) qu'on a

$$(39) \quad \mathcal{T}(\Delta t, t; u(\bullet)) = \frac{\Delta t}{2} \frac{d^2 u}{dt^2}(t) + O(\Delta t^2).$$

• Comme à la relation (3), nous introduisons l'erreur $\epsilon_{\Delta t}^k$ par

$$(40) \quad \epsilon_{\Delta t}^k = u_{\Delta t}^k - u(k\Delta t).$$

On soustrait la relation (38) de (37). Il vient

$$(41) \quad \frac{1}{\Delta t} (\epsilon_{\Delta t}^{k+1} - \epsilon_{\Delta t}^k) + \lambda \epsilon_{\Delta t}^k = \mathcal{T}(\Delta t, k\Delta t; u(\bullet))$$

et l'erreur $\epsilon_{\Delta t}^k$ vérifie une équation analogue à celle du schéma numérique, avec comme **source** (au membre de droite) l'erreur de troncature.

• On suppose $0 \leq k\Delta t \leq T$, avec T fixé. On peut donc majorer uniformément la dérivée seconde $\frac{d^2 u}{dt^2}(t)$ sur cet intervalle :

$$(42) \quad \left| \frac{d^2 u}{dt^2}(k\Delta t) \right| \leq C, \quad \forall k \in \mathbb{N} \text{ tel que } k\Delta t \leq T.$$

On tire alors de (39) et (42)

$$(43) \quad |\mathcal{T}(\Delta t, k\Delta t; u(\bullet))| \leq C\Delta t, \quad k\Delta t \leq T.$$

• On suppose de plus que le schéma d'Euler est **stable**, propriété qui est vérifiée si le pas de temps n'est pas choisi trop grand :

$$(44) \quad 0 < \lambda\Delta t \leq 1.$$

On peut alors écrire la relation (41) sous la forme

$$(45) \quad \epsilon_{\Delta t}^{k+1} = (1 - \lambda \Delta t) \epsilon_{\Delta t}^k - \Delta t | \mathcal{T}(\Delta t, k \Delta t; u(\bullet)) | .$$

On tire alors de la relation (44) : $0 < 1 - \lambda \Delta t < 1$ et en prenant les valeurs absolues de part et d'autre de (45) :

$$| \epsilon_{\Delta t}^{k+1} | \leq | \epsilon_{\Delta t}^k | + \Delta t | \mathcal{T}(\Delta t, k \Delta t; u(\bullet)) | ,$$

donc en tenant compte de la relation (43) :

$$(46) \quad | \epsilon_{\Delta t}^{k+1} | \leq | \epsilon_{\Delta t}^k | + C \Delta t^2 , \quad k \Delta t \leq T .$$

• On écrit la chaîne d'inégalités (46), depuis $\epsilon_{\Delta t}^0 \equiv 0$ jusqu'à $\epsilon_{\Delta t}^k$. Il vient

$$| \epsilon_{\Delta t}^k | \leq | \epsilon_{\Delta t}^{k-1} | + C \Delta t^2$$

$$| \epsilon_{\Delta t}^{k-1} | \leq | \epsilon_{\Delta t}^{k-2} | + C \Delta t^2$$

...

$$| \epsilon_{\Delta t}^1 | \leq | \epsilon_{\Delta t}^0 | + C \Delta t^2 .$$

Puis on ajoute toutes ces inégalités, ce qui est possible grâce à la stabilité (44). On en déduit, puisque $k \Delta t \leq T$:

$$(47) \quad | \epsilon_{\Delta t}^k | \leq C k \Delta t^2 \leq C T \Delta t .$$

Nous venons d'établir que l'**erreur** est majoré par une constante multipliée par Δt . Elle est d'ordre un en Δt tout comme l'erreur de troncature. Si nous introduisons l'erreur $\delta(\Delta t)$ comme à la relation (4), on tire de la relation (47) l'estimation

$$(48) \quad \delta(\Delta t) \leq C T \Delta t ,$$

ce qui confirme bien le résultat établi.