

## Chapitre VI : Méthodes de résolution approchée des équations différentielles ordinaires

### Introduction

Nous nous intéressons à la résolution approchée d'une équation différentielle ordinaire du premier ordre :

$$y'(t) = f(t, y(t)).$$

Afin de simplifier la présentation, nous nous limitons au cas où la fonction  $f$  est définie sur un sous-ensemble de  $\mathbb{R}^2$  de la forme  $I \times U$ , où  $I$  est un intervalle de  $\mathbb{R}$  et  $U$  désigne un sous-ensemble ouvert de  $\mathbb{R}$ , et  $\bar{a}$  valeurs réelles.

Les différentes méthodes de résolution approchée que nous décrivons dans la suite s'étendent sans difficulté majeure au cas où la fonction  $f$  est définie sur un sous-ensemble de  $\mathbb{R} \times \mathbb{R}^n$  de la forme  $I \times U$ , où  $U$  désigne désormais un sous-ensemble ouvert de  $\mathbb{R}^n$ , et  $\bar{a}$  valeurs vectorielles. De plus, elles s'étendent aussi à des équations différentielles ordinaires d'ordre strictement supérieur à un.

Nous nous plaçons de plus dans un cadre où le problème de Cauchy associé à l'équation différentielle  $y'(t) = f(t, y(t))$  est bien posé (au moins localement). Étant donné un nombre  $t_0 \in I$ , ceci signifie qu'il existe une unique solution  $y \in \mathcal{C}^1(I_0, U)$  sur un intervalle  $I_0 \subset I$  contenant  $t_0$  pour la donnée initiale  $y(t_0) = y_0 \in U$ , et que cette solution dépend de manière continue de la valeur  $y_0$  de la donnée initiale. Rappelons que le théorème de Cauchy-Lipschitz garantit le caractère bien posé de l'équation différentielle  $y'(t) = f(t, y(t))$ .

Théorème de Cauchy-Lipschitz: Soit  $I$ , un intervalle ouvert de  $\mathbb{R}$  et  $U$ , un sous-

ensemble ouvert de  $\mathbb{R}$ . Considérons une fonction  $f \in \mathcal{C}^0(I \times U, \mathbb{R})$  qui est localement lipschitzienne par rapport à  $y \in U$ , soit telle que:

$$\forall (t_0, y_0) \in I \times U, \exists (\delta_0, \rho_0) \in (\mathbb{R}_+^*)^2 \text{ t.q. } \exists K_0 \in \mathbb{R}_+ \text{ t.q. } \forall t \in I, \\ \forall (y, y') \in U^2, |t - t_0| < \delta_0 \text{ et } \max(|y - y_0|, |y' - y_0'|) < \rho_0 \Rightarrow |f(t, y) - f(t, y')| \leq K_0 |y - y'|.$$

Étant donné deux nombres  $t_0 \in I$  et  $y_0 \in U$ , il existe un nombre strictement positif  $\varepsilon_0$  tel que il existe une unique solution  $y \in \mathcal{C}^1([t_0 - \varepsilon_0, t_0 + \varepsilon_0], U)$  de l'équation différentielle:

$$\forall t \in (t_0 - \varepsilon_0, t_0 + \varepsilon_0), y'(t) = f(t, y(t)),$$

pour la donnée initiale  $y(t_0) = y_0$ . De plus, si deux suites  $(t_0^{(n)})_{n \in \mathbb{N}} \in I^{\mathbb{N}}$  et  $(y_0^{(n)})_{n \in \mathbb{N}} \in U^{\mathbb{N}}$  satisfont:

$$t_0^{(n)} \xrightarrow{n \rightarrow +\infty} t_0 \text{ et } y_0^{(n)} \xrightarrow{n \rightarrow +\infty} y_0,$$

alors, les solutions  $y^{(n)} \in \mathcal{C}^1([t_0 - \varepsilon^{(n)}, t_0 + \varepsilon^{(n)}], U)$  de l'équation différentielle pour les données initiales  $y^{(n)}(t_0^{(n)}) = y_0^{(n)}$  satisfont:

$$\lim_{n \rightarrow +\infty} \varepsilon^{(n)} \geq \varepsilon_0 \text{ et } \max_{t \in (t_0 - \varepsilon_0, t_0 + \varepsilon_0)} |y^{(n)}(t) - y(t)| \xrightarrow{n \rightarrow +\infty} 0.$$

Lorsque cela s'avèrera nécessaire, nous nous placerons de plus dans le cas où les solutions de l'équation différentielle sont globales, soit définies sur tout l'intervalle. Ceci n'inclut pas de traiter des équations différentielles pour lesquelles l'unicité des solutions n'est pas garantie, mais cela facilite l'analyse de la convergence des méthodes numériques proposées puisque leur limite est alors définie de manière unique (et si nécessaire, globalement).

Les méthodes de résolution approchée auxquelles nous nous intéressons reposent sur un processus de discrétisation. Étant donné une solution  $y \in \mathcal{C}^1([t_0, t_0 + T], U)$  du problème de Cauchy pour une donnée initiale  $y(t_0) = y_0$ , nous commencerons par introduire une subdivision  $t_0 < t_1 < \dots < t_N = t_0 + T$  de son intervalle de définition  $(t_0, t_0 + T)$ , et nous cherchons à déterminer des

valeurs approchées  $(y_n)_{0 \leq n \leq N}$  des valeurs exactes  $(y(t_n))_{0 \leq n \leq N}$  prises par la solution  $y$  aux temps  $(t_n)_{0 \leq n \leq N}$ .

Les méthodes de résolution approchées que nous discuterons seront des méthodes à un pas. Le calcul de la valeur approchée  $y_{n+1}$  ne dépendra que de la valeur approchée  $y_n$ . Les méthodes d'Euler explicite et implicite, et du point milieu sont des exemples de méthodes à un pas. Nous chercherons à déterminer sous quelles conditions ces méthodes sont convergentes, c'est-à-dire celles que l'approximation ainsi obtenue converge vers la solution exacte lorsque le pas de la subdivision:

$$h = \max_{0 \leq n \leq N-1} |y_{n+1} - y_n|,$$

tend vers 0. En particulier, nous introduisons les notions de consistance et de stabilité qui garantissent la convergence de ces méthodes. Nous analyserons également l'ordre de ces méthodes afin de déterminer leur précision, et l'influence des erreurs d'arrondis.

Nous détaillerons ensuite le cas des méthodes de Runge-Kutta, avant de conclure par une présentation succincte des difficultés à surmonter lorsque les problèmes de Cauchy étudiés ne sont plus bien posés mathématiquement ou numériquement, ou encore bien conditionnés.

Nous renvoyons aux ouvrages de Michel Crouzeix et Alain Hignat intitulé "Analyse numérique des équations différentielles", et de Jean Pierre Demailly intitulé "Analyse numérique et équations différentielles" pour de plus amples détails au sujet des diverses méthodes de résolution approchées des équations différentielles ordinaires.

## 2. Étude générale des méthodes à un pas

Soit  $I$ , un intervalle ouvert de  $\mathbb{R}$ . Nous considérons une fonction  $f \in C^0(I)$

$\mathbb{R}, \mathbb{R}$ ), et nous nous intéressons à la résolution approchée du problème de Cauchy :

$$\begin{cases} \forall t \in [t_0, t_0 + T], & y'(t) = f(t, y(t)), \\ y(t_0) = y_0, \end{cases}$$

où les nombres  $t_0$  et  $t_0 + T$  appartiennent à l'intervalle  $I$ , et la valeur  $y_0$  de la donnée initiale est réelle.

Afin de simplifier la présentation qui suit, nous supposons que la fonction  $f$  est globalement lipschitzienne par rapport à  $y \in \mathbb{R}$ , à savoir qu'il existe un nombre positif  $K$  tel que :

$$\forall t \in I, \forall (y, y') \in \mathbb{R}^2, |f(t, y') - f(t, y)| \leq K |y' - y|.$$

Ceci garantit qu'il existe une unique solution  $y \in C^1([t_0, t_0 + T], \mathbb{R})$  du problème de Cauchy précédent, qui dépend de plus de manière continue de la donnée initiale  $y_0$ . Cette hypothèse simplifie la preuve sur la forme de certaines des estimations d'erreur que nous allons établir. Notons cependant que la plupart de ces estimations demeurent valables sous des hypothèses plus faibles telles que la caractéristique localement lipschitzienne par rapport à  $y \in \mathbb{R}$  de la fonction  $f$ .

Afin de déterminer une valeur approchée de la solution  $y$ , nous introduisons une subdivision :

$$t_0 < t_1 < \dots < t_N = t_0 + T,$$

du segment  $[t_0, t_0 + T]$ . Nous notons :

$$\forall n \in [0, N-1], h_n = t_{n+1} - t_n,$$

les pas entre les temps  $t_n$  et  $t_{n+1}$ , et :

$$h_{\max} = \max_{0 \leq n \leq N-1} h_n,$$

le pas maximal de la subdivision. Notre objectif est de déterminer

des valeurs approchées  $(y_n)_{0 \leq n \leq N}$  de la solution  $(y(t_n))_{0 \leq n \leq N}$

aux temps  $(t_n)_{0 \leq n \leq N}$ . Dans ce but, nous introduisons les méthodes

numériques à un pas, qui sont définies de la façon suivante.

## 1. Définition et exemples

Définition: Une méthode numérique à un pas pour la résolution approchée du problème de Cauchy:

$$\left\{ \begin{array}{l} \forall t \in [t_0, t_0 + T], y'(t) = f(t, y(t)), \\ y(t_0) = y_0, \end{array} \right.$$

est définie par la donnée d'une fonction  $\mathbb{F} \in \mathcal{C}^0([t_0, t_0 + T] \times \mathbb{R}^2, \mathbb{R})$  et la formule de récurrence pour les valeurs approchées  $(y_n)_{0 \leq n \leq N}$  aux temps  $(t_n)_{0 \leq n \leq N}$ :

$$\forall n \in [0, N-1], y_{n+1} = y_n + h_n \mathbb{F}(t_n, y_n, h_n).$$

La terminologie de méthodes numériques à un pas s'explique par le fait que le calcul de la valeur approchée  $y_{n+1}$  au temps  $t_{n+1}$  ne dépend que de la valeur approchée  $y_n$  au temps  $t_n$  (et du pas  $h_n$ ).

Il est possible d'affaiblir les hypothèses sur la fonction  $\mathbb{F}$  qui peut ne pas être définie qu'en un voisinage du triplet  $(t_0, y_0, 0)$ . Il est néanmoins nécessaire que son ensemble de définition recouvre les valeurs prises par les couples  $(t, y)$  pour  $t \in [t_0, t_0 + T]$ , où  $y$  désigne l'unique solution du problème de Cauchy considéré.

Exemples: (i) Méthode d'Euler explicite (ou progressive)

Cette méthode correspond à la fonction  $\mathbb{F}$  définie par:

$$\forall (t, y, h) \in [t_0, t_0 + T] \times \mathbb{R}^2, \mathbb{F}(t, y, h) = f(t, y).$$

La suite  $(y_n)_{0 \leq n \leq N}$  des valeurs approchées au temps  $(t_n)_{0 \leq n \leq N}$  est donnée par la formule de récurrence:

$$\forall n \in [0, N-1], y_{n+1} = y_n + h_n f(t_n, y_n).$$

Cette suite est toujours bien définie sous les hypothèses générales que nous avons imposés pour la fonction  $f$ .

(ii) Méthode d'Euler implicite (ou rétrograde)

Cette méthode correspond à la fonction  $\Phi$  définie par:

$$\forall (t, y, h) \in (t_0, t_0 + \tau] \times \mathbb{R}^2, \Phi(t, y, h) = f\left(t+h, \left[1-h \cdot h f(t, y)\right] \cdot y\right)$$

La suite  $(y_n)_{0 \leq n \leq N}$  des valeurs approchées aux temps  $(t_n)_{0 \leq n \leq N}$  est donnée par la formule de récurrence implicite:

$$\forall 0 \leq n \leq N-1, y_{n+1} = y_n + h_n f(t_{n+1}, y_{n+1}).$$

Il n'est pas clair que cette suite soit bien définie. Le calcul de la valeur  $y_{n+1}$  passe par la résolution d'une équation de type point fixe, ce qui ne s'avère possible que si la fonction  $f$  satisfait les hypothèses ad hoc. Mettons que ces hypothèses reviennent à supposer que l'application  $\xi \mapsto h \cdot f(t+h, \xi)$  est un homéomorphisme de  $\mathbb{R}$  sur  $\mathbb{R}$ , qui dépend de manière continue des paramètres  $t$  et  $h$  (auquel cas la fonction  $\Phi$  est bien définie et continue sur  $(t_0, t_0 + \tau] \times \mathbb{R}^2$ ).

### (iii) Méthode du point milieu (ou d'Euler modifiée)

Cette méthode correspond à la fonction  $\Phi$  définie par:

$$\forall (t, y, h) \in (t_0, t_0 + h] \times \mathbb{R}^2, \Phi(t, y, h) = f\left(t + \frac{h}{2}, y + \frac{h}{2} f(t, y)\right).$$

La suite  $(y_n)_{0 \leq n \leq N}$  des valeurs approchées aux temps  $(t_n)_{0 \leq n \leq N}$  est donnée par l'algorithme récursif:

$$\forall 0 \leq n \leq N-1, \begin{cases} y_{n+\frac{1}{2}} = y_n + \frac{h_n}{2} f(t_n, y_n), \\ y_{n+1} = f\left(t_n + \frac{h_n}{2}, y_{n+\frac{1}{2}}\right). \end{cases}$$

Cette suite est toujours bien définie sous les hypothèses générales que nous avons imposées pour la fonction  $f$ .

Nous cherchons dans la suite à établir sous quelles conditions une méthode numérique à un pas est convergente, à savoir qu'elle vérifie la propriété:

$$\max_{0 \leq n \leq N} |y(t_n) - y_n| \rightarrow 0,$$

lorsque le pas maximal  $h_{\max}$  de la subdivision tend vers 0. Ceci passe par l'introduction des notions de consistance et de stabilité.

## 2. Consistance d'une méthode numérique à un pas

Considérons la méthode numérique à un pas associée à la fonction  $\mathbb{F} \in \mathcal{C}^0([t_0, t_0 + T] \times \mathbb{R}^2, \mathbb{R})$ .

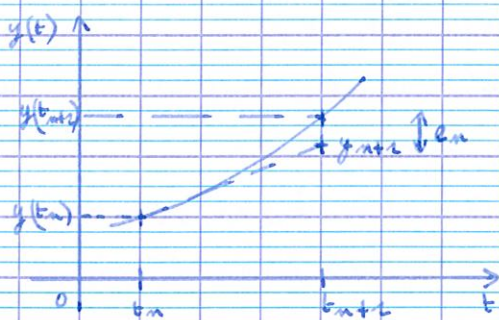
Definition: Soit  $0 \leq n \leq N-1$ . Considérons une solution  $y \in \mathcal{C}^1([t_0, t_0 + T], \mathbb{R}^2)$  de l'équation différentielle:

$$\forall t \in [t_0, t_0 + T], y'(t) = \mathbb{F}(t, y(t)).$$

L'erreur de consistance au temps  $t_n$  relative à la solution  $y$  est définie par:

$$e_{cn} = y(t_{n+1}) - y(t_n) - h_n \mathbb{F}(t_n, y(t_n), h_n).$$

L'erreur de consistance au temps  $t_n$  mesure donc l'écart entre la valeur exacte au temps  $t_{n+1}$  donnée par  $y(t_{n+1})$  et la valeur approchée fournie par la méthode numérique à un pas à partir de la valeur exacte  $y(t_n)$  au temps  $t_n$ . Autrement dit, elle mesure l'erreur faite au  $n$ <sup>ème</sup> pas lorsque l'équation différentielle est remplacée par le schéma numérique à un pas.



Cette erreur n'est pas égale à l'erreur finale entre la valeur exacte  $y(t_n)$  au temps  $t_n$  et sa valeur approchée  $y_n$ , laquelle est donnée par la différence  $y(t_n) - y_n$ . Cependant, il semble clair que cette erreur sera au moins de l'ordre de  $\sum_{k=0}^{n-1} \tau_k$  et  $h$ , quantité qu'il est donc important de pouvoir contrôler.

Definition: Une méthode numérique à un pas est consistante si et seulement si quelle que soit la solution  $y \in \mathcal{C}^1([t_0, t_0 + T] \times \mathbb{R}^2, \mathbb{R})$  de l'é-

équation différentielle:

$$\forall t \in (t_0, t_0 + T], y'(t) = f(t, y(t)),$$

les erreurs de consistance  $(e_n)_{0 \leq n \leq N-1}$  aux temps  $(t_n)_{0 \leq n \leq N-1}$  relatives à la solution  $y$  satisfont:

$$\sum_{n=0}^{N-1} |e_n| \xrightarrow{h_{\max} \rightarrow 0} 0.$$

Une méthode numérique à un pas est donc consistante lorsqu'il est possible de contrôler l'accumulation des erreurs de consistance quelle que soit la solution  $y$  considérée. Nous avons le critère de consistance suivant pour une méthode numérique à un pas.

Théorème: Une méthode numérique à un pas est consistante avec l'équation différentielle  $y'(t) = f(t, y(t))$  si et seulement si:

$$\forall (t, y) \in (t_0, t_0 + T], \Phi(t, y, 0) = f(t, y).$$

Preuve:

Soit  $y \in C^1((t_0, t_0 + T], \mathbb{R})$ , une solution de l'équation différentielle:

$$\forall t \in (t_0, t_0 + T], y'(t) = f(t, y(t)).$$

Pour  $0 \leq n \leq N-1$ , l'erreur de consistance  $e_n$  au temps  $t_n$  relative à la solution  $y$  est donnée par:

$$e_n = y(t_{n+1}) - y(t_n) - h_n \Phi(t_n, y(t_n), h_n).$$

Pour le théorème des accroissements finis, il existe un nombre  $\xi_n \in a_n \leq \xi_n \leq b_{n+1}$  tel que:

$$y(t_{n+1}) - y(t_n) = h_n y'(\xi_n) = h_n f(\xi_n, y(\xi_n)),$$

De sorte que:

$$e_n = h_n [f(\xi_n, y(\xi_n)) - \Phi(\xi_n, y(\xi_n), 0)] + h_n [\Phi(\xi_n, y(\xi_n), 0) - \Phi(t_n, y(t_n), h_n)].$$

À ce stade, nous pouvons supposer que:

$$h_{\max} \leq 1.$$

Notons alors que la fonction  $(t, h) \mapsto \Phi(t, y(t), h)$  est continue sur le compact  $[t_0, t_0 + T] \times (0, 1]$ , donc uniformément continue sur cet



ensemble. Étant donné un nombre strictement positif  $\varepsilon$ , il existe donc un nombre strictement positif  $\delta$  tel que :

$$\forall (t, t') \in [t_0, t_0 + T]^2, \forall (k, k') \in [0, 1]^2, |t - t'| \leq \delta \text{ et } |k - k'| \leq \delta \Rightarrow |F(t', y(t'), k') - F(t, y(t), k)| \leq \varepsilon.$$

Lorsque  $h_{\max} \leq \delta$ , nous concluons donc que :

$$\sum_{n=0}^{N-1} h_n |F(t_n, y(t_n), 0) - F(t_n, y(t_n), h_n)| \leq \sum_{n=0}^{N-1} h_n \varepsilon = T\varepsilon$$

Il nous reste alors que la fonction  $t \mapsto f(t, y(t)) - F(t, y(t), 0)$  est continue sur

$[t_0, t_0 + T]$ , de sorte que par le théorème des séries de Riemann :

$$\sum_{n=0}^{N-1} h_n |f(t_n, y(t_n)) - F(t_n, y(t_n), 0)| \xrightarrow[h_{\max} \rightarrow 0]{} \int_{t_0}^{t_0+T} dt |f(t, y(t)) - F(t, y(t), 0)|$$

En conclusion, nous avons établi que :

$$\sum_{n=0}^{N-1} |e_n| \xrightarrow[h_{\max} \rightarrow 0]{} \int_{t_0}^{t_0+T} |f(t, y(t)) - F(t, y(t), 0)| dt$$

La consistance de la méthode numérique est donc équivalente au fait que

$$\int_{t_0}^{t_0+T} |f(t, y(t)) - F(t, y(t), 0)| dt = 0,$$

quelle que soit la solution  $y$  considérée, soit au fait que :

$$\forall t \in [t_0, t_0 + T], F(t, y(t), 0) = f(t, y(t)).$$

Par le théorème de Cauchy-Lipschitz, quelle que soit  $(t, y) \in [t_0, t_0 + T] \times \mathbb{R}$ , il existe une solution  $y$  telle que :

$$y(t) = y,$$

ce qui suffit à établir l'équivalence recherchée.

Exemples : (i) Méthode d'Euler explicite

Cette méthode est consistante puisque :

$$\forall (t, y) \in [t_0, t_0 + T] \times \mathbb{R}, F(t, y, 0) = f(t, y).$$

(ii) Méthode d'Euler implicite

Cette méthode est consistante (lorsqu'elle est bien définie), puisque :

$$\forall (t, y) \in [t_0, t_0 + T] \times \mathbb{R}, F(t, y, 0) = f(t, \text{Id}^{-1}(y)) = f(t, y).$$

(iii) Méthode du point milieu

Cette méthode est consistante puisque :

$$\forall (t, y) \in [t_0, t_0 + T] \times \mathbb{R}, F(t, y, 0) = f(t, y).$$

### 3. Stabilité d'une méthode numérique à un pas

Considérons la méthode numérique à un pas associée à la fonction  $\mathbb{F} \in \mathcal{C}^0([a, b], \mathbb{R}^2, \mathbb{R})$ . En pratique, le calcul des valeurs approchées  $y_n$  est entaché de erreurs d'arrondis. Il est crucial de contrôler ces erreurs afin que les valeurs approchées  $y_n$  soient significatives.

Definition: Une méthode numérique à un pas est stable si et seulement s'il existe un nombre positif  $S$  telle que, quelles que soient les suites  $(y_n)_{0 \leq n \leq N} \in \mathbb{R}^{N+1}$ ,  $(y'_n)_{0 \leq n \leq N} \in \mathbb{R}^{N+1}$  et  $(h_n)_{0 \leq n \leq N-1} \in \mathbb{R}^N$  telles que:

$$\forall 0 \leq n \leq N-1, \begin{cases} y_{n+1} = y_n + h_n \mathbb{F}(t_n, y_n, h_n), \\ \text{et} \\ y'_{n+1} = y'_n + h_n \mathbb{F}'(t_n, y'_n, h_n) + h_n, \end{cases}$$

nous avons l'inégalité:

$$\max_{0 \leq n \leq N} |y'_n - y_n| \leq S \left[ |y'_0 - y_0| + \sum_{n=0}^{N-1} |h_n| \right].$$

Le nombre  $S$  est alors une constante de stabilité de la méthode.

Une notion de stabilité est intrinsèque à la méthode numérique consistant au sens où elle ne dépend pas de l'équation différentielle  $y'(t) = f(t, y(t))$  à résoudre. Elle garantit qu'une petite erreur sur la donnée initiale  $y_0$  et de petites erreurs d'arrondis à chaque étape de la méthode ne conduisent pas à une erreur finale incontrôlable. Nous pouvons établir la condition suivante pour la stabilité d'une méthode numérique.

Théorème: Si la fonction  $\mathbb{F}$  est globalement lipschitzienne par rapport à la variable  $y \in \mathbb{R}$ , c'est-à-dire s'il existe un nombre positif  $K$  tel que:

$$\forall t \in [a, b], \forall (y, y') \in \mathbb{R}^2, \forall h \in \mathbb{R}, |\mathbb{F}(t, y, h) - \mathbb{F}(t, y', h)| \leq K|y - y'|$$

alors, la méthode numérique à un pas est stable. De plus, le nombre  $S = e^{KT}$  est une constante de stabilité de la méthode.

Preuve:

Considérons des suites  $(y_n)_{0 \leq n \leq N} \in \mathbb{R}^{N+2}$ ,  $(y'_n)_{0 \leq n \leq N} \in \mathbb{R}^{N+2}$  et  $(\delta_n)_{0 \leq n \leq N-1}$  telles que:

$$\forall 0 \leq n \leq N-1, \begin{cases} y_{n+1} = y_n + h_n \Phi(t_n, y_n, h_n), \\ \text{et} \\ y'_{n+1} = y'_n + h_n \Psi(t_n, y'_n, h_n) + \delta_n. \end{cases}$$

Il nous vient alors:

$$\forall 0 \leq n \leq N-1, y'_{n+1} - y_{n+1} = y'_n - y_n + h_n [\Psi(t_n, y'_n, h_n) - \Phi(t_n, y_n, h_n)] + \delta_n$$

$$\begin{aligned} \Rightarrow |y'_{n+1} - y_{n+1}| &\leq (1 + h_n k) |y'_n - y_n| + |\delta_n| \\ &\leq e^{k h_n} |y'_n - y_n| + |\delta_n| \\ &= e^{k(t_{n+1} - t_n)} |y'_n - y_n| + |\delta_n| \end{aligned}$$

$$\Rightarrow e^{-k t_{n+1}} |y'_{n+1} - y_{n+1}| \leq e^{-k t_n} |y'_n - y_n| + e^{-k t_{n+1}} |\delta_n|$$

De suite que:

$$\sum_{l=0}^n (e^{-k t_{l+1}} |y'_{l+1} - y_{l+1}| - e^{-k t_l} |y'_l - y_l|) \leq \sum_{l=0}^n e^{-k t_{l+1}} |\delta_l|.$$

Il s'ensuit que:

$$e^{-k t_{n+1}} |y'_{n+1} - y_{n+1}| \leq e^{-k t_0} |y'_0 - y_0| + \sum_{l=0}^n e^{-k t_{l+1}} |\delta_l|.$$

Il vient alors:

$$\forall 0 \leq n \leq N+1, |y'_{n+1} - y_{n+1}| \leq e^{k T} (|y'_0 - y_0| + \sum_{l=0}^n |\delta_l|),$$

ce qui assure l'inégalité de stabilité pour une constante de stabilité  $S = e^{kT}$ .

Exemples: Rappelons que la fonction  $f$  est globalement lipschitzienne par rapport

à  $y \in \mathbb{R}$ , soit qu'il existe un nombre positif  $K$  tel que:

$$\forall t \in I, \forall (y, y') \in \mathbb{R}^2, |f(t, y') - f(t, y)| \leq K |y - y'|.$$

(i) Méthode d'Euler explicite

Comme

$$\forall t \in [t_0, t_0 + 1], \forall (y, y') \in \mathbb{R}^2, \forall h \in \mathbb{R}, |\Phi(t, y', h) - \Phi(t, y, h)| \leq K |y - y'|,$$

la méthode est stable pour une constante de stabilité  $S = e^{kT}$ .

(ii) Méthode d'Euler implicite

Rappelons que la fonction  $\mathbb{E}$  est définie par:

$$\forall (t, y, h) \in (t_0, t_0 + \tau) \times \mathbb{R}^2, \mathbb{I}(t, y, h) = f(t+h, (Id - h f(t+h, \cdot))^{-1}(y))$$

et supposons que, pour des nombres  $t \in (t_0, t_0 + \tau)$  et  $h \in \mathbb{R}$  fixés, l'application  $Id - h f(t+h, \cdot)$  soit un homéomorphisme de  $\mathbb{R}^2$  sur  $\mathbb{R}^2$ . Dans ce cas, la fonction  $y \mapsto \mathbb{I}(t, y, h)$  est bien définie. De plus, nous pouvons dire:

$$\forall (y, y') \in \mathbb{R}^2, |\mathbb{I}(t, y', h) - \mathbb{I}(t, y, h)| \leq K |(Id - h f(t+h, \cdot))^{-1}(y') - (Id - h f(t+h, \cdot))^{-1}(y)|$$

Soit alors  $z' = (Id - h f(t+h, \cdot))^{-1}(y')$  et  $z = (Id - h f(t+h, \cdot))^{-1}(y)$ . Il vient:

$$\begin{aligned} |y' - y| &= |z' - h f(t+h, z') - z + h f(t+h, z)| \\ &\geq |z' - z| - |h| |f(t+h, z') - f(t+h, z)| \\ &\geq |z' - z| - K |h| |z' - z|. \end{aligned}$$

Lorsque  $|h| < \frac{1}{K}$ , nous en déduisons que:

$$|z' - z| \leq \frac{|y' - y|}{1 - K|h|},$$

ce qui conduit à l'inégalité:

$$\forall (y, y') \in \mathbb{R}^2, |\mathbb{I}(t, y', h) - \mathbb{I}(t, y, h)| \leq \frac{K}{1 - K|h|} |y' - y|.$$

Notons qu'en pratique, nous ne nous intéressons qu'à des valeurs de  $h$  petites, de sorte que nous pouvons introduire une notion de stabilité où le pas maximal  $h_{\max}$  est borné par le nombre  $\frac{1}{K}$ . Pour cette notion de stabilité, la méthode d'Euler implicite s'avère stable (lorsqu'elle est bien définie) pour une constante de stabilité  $\frac{K\tau}{1 - K h_{\max}}$ , qui est du même ordre que celle de la méthode d'Euler explicite lorsque le pas maximal  $h_{\max}$  tend vers 0.

### (iii) Méthode du point milieu

Rappelons que la fonction  $\mathbb{I}$  est définie par:

$$\forall (t, y, h) \in (t_0, t_0 + \tau) \times \mathbb{R}^2, \mathbb{I}(t, y, h) = f\left(t + \frac{h}{2}, y + \frac{h}{2} f(t, y)\right)$$

Il s'ensuit que:

$$\begin{aligned} \forall (t, h) \in (t_0, t_0 + \tau) \times \mathbb{R}, \forall (y, y') \in \mathbb{R}^2, |\mathbb{I}(t, y', h) - \mathbb{I}(t, y, h)| \\ \leq K \left| y' + \frac{h}{2} f(t, y') - y - \frac{h}{2} f(t, y) \right| \end{aligned}$$

$$\leq K |y' - y| \left( 1 + K \frac{h_1}{2} \right)$$

Comme pour la méthode d'Euler implicite, nous pouvons introduire une notion de stabilité pour laquelle le pas maximal  $h_{\max}$  est borné, par exemple, par 1. La méthode du point milieu s'avère alors stable pour une constante de stabilité  $\frac{1}{2} K T (1 + K \frac{h_{\max}}{2})$ , qui est du même ordre que celle de la méthode d'Euler explicite lorsque le pas maximal  $h_{\max}$  tend vers 0.

#### 4. Convergence d'une méthode numérique à un pas

Considérons la méthode numérique à un pas associée à la fonction  $\Phi \in \mathcal{C}^0([t_0, t_0+T], \mathbb{R}^2, \mathbb{R})$ .

Définition: La méthode numérique à un pas est convergente si et seulement, quelle que soit une solution  $y \in \mathcal{C}^2([t_0, t_0+T], \mathbb{R})$  de l'équation différentielle:

$$\forall t \in [t_0, t_0+T], y'(t) = f(t, y(t)),$$

les suites  $(y_n)_{0 \leq n \leq N} \in \mathbb{R}^{N+2}$  telles que :

$$\forall 0 \leq n \leq N-1, y_{n+1} = y_n + h_n \Phi(t_n, y_n, h_n),$$

satisfont à la convergence :

$$\max_{0 \leq n \leq N} |y(t_n) - y_n| \rightarrow 0,$$

lorsque  $y(t_0) \rightarrow y_0$  et le pas maximal  $h_{\max} = \max_{0 \leq n \leq N-1} h_n \rightarrow 0$ .

Cette notion garantit que la méthode numérique à un pas considérée fournit une approximation valable d'une solution  $y \in \mathcal{C}^2([t_0, t_0+T], \mathbb{R})$  de l'équation différentielle considérée, puisque l'erreur globale définie par :

$$E_N = \max_{0 \leq n \leq N} |y(t_n) - y_n|,$$

est aussi petite que souhaitée lorsque le pas maximal  $h_{\max}$  est suffisamment petit, et les données initiales  $y(t_0)$  et  $y_0$ , suffisamment proches.

Les notions de consistance et de stabilité introduites précédemment permettent de garantir la convergence d'une méthode numérique à un pas, comme l'exprime le résultat suivant.

Théorème: Si la méthode numérique à un pas est consistante et stable, alors elle est convergente.

Preuve:

Considérons une solution  $y \in C^1([t_0, t_0+T], \mathbb{R})$  de l'équation différentielle

$$\forall t \in [t_0, t_0+T], y'(t) = f(t, y(t)),$$

et désignons par  $e_n$  l'erreur de consistance au temps  $t_n$  relative à la solution qui est définie par:

$$\forall 0 \leq n \leq N-1, e_n = y(t_{n+1}) - y(t_n) - h_n \mathbb{F}(t_n, y(t_n), h_n).$$

Étant donnée une suite  $(y_n)_{0 \leq n \leq N}$  telle que:

$$\forall 0 \leq n \leq N-1, y_{n+1} = y_n + h_n \mathbb{F}(t_n, y_n, h_n),$$

la stabilité de la méthode garantit l'existence d'un nombre positif  $s$  tel que :

$$\forall 0 \leq n \leq N, |y(t_n) - y_n| \leq s \left( |y(t_0) - y_0| + \sum_{k=0}^{n-1} |e_k| \right).$$

La consistance de la méthode garantit quant à elle que:

$$\sum_{k=0}^{N-1} |e_k| \xrightarrow{h_{\max} \rightarrow 0} 0,$$

ce qui suffit à achever la preuve.

Le théorème fondamental ramène la validité d'une méthode numérique à un pas aux seules notions de consistance et de stabilité dont nous avons pu évaluer qu'elles étaient faciles à caractériser. Nous en déduisons en particulier le corollaire suivant.

Corollaire: Supposons que:

$$(i) \forall (t, y) \in [t_0, t_0+T], \mathbb{F}(t, y, 0) = f(t, y)$$

(ii) La fonction  $\mathbb{F}$  est globalement lipschitzienne par rapport à la variable  $y \in \mathbb{R}$ , c'est-à-dire qu'il existe un nombre positif  $K$

tel que :

$\forall t \in (t_0, t_0+h], \forall (y, y') \in \mathbb{R}^2, \forall h \in \mathbb{R}, |\Phi(t, y', h) - \Phi(t, y, h)| \leq K|y'|$   
La méthode numérique à un pas est alors convergente.

Preuve :

Stabilité

Exemples : (i) Méthode d'Euler explicite

La méthode est convergente par la corollaire précédent.

(ii) Méthode d'Euler implicite

La méthode est convergente sous les conditions de stabilité du paragraphe précédent.

(iii) Méthode du point milieu

La méthode est convergente sous les conditions de stabilité du paragraphe précédent.

## 5. Ordre d'une méthode numérique à un pas

Considérons la méthode numérique à un pas associée à la fonction  $\Phi \in \mathcal{C}^p([t_0, t_0+T] \times \mathbb{R}^2, \mathbb{R})$ .

Définition : Soit  $p \in \mathbb{N}^*$ . La méthode numérique à un pas est d'ordre (au moins égal à)  $p$  si et seulement si il existe un nombre positif  $C$  tel que, quelle que soit la solution  $y \in \mathcal{C}^{p+1}([t_0, t_0+T], \mathbb{R})$  de l'équation différentielle :

$$\forall t \in [t_0, t_0+T], y'(t) = f(t, y(t)),$$

les erreurs de consistance  $(e_n)_{0 \leq n \leq N-1}$  aux temps  $(t_n)_{0 \leq n \leq N-1}$  relatives à la solution  $y$  satisfont :

$$\sum_{n=0}^{N-1} |e_n| = \sum_{n=0}^{N-1} |y(t_{n+1}) - y(t_n) - h_n \Phi(t_n, y(t_n), h_n)| \leq C h_{\max}^p$$

où le pas maximal  $h_{\max}$  est défini par:

$$h_{\max} = \max_{0 \leq n \leq N-1} h_n.$$

Remarque: La méthode numérique est (parfois) dite d'ordre exactement égal à  $p$  si et seulement si elle est d'ordre (au moins égal à)  $p$ , mais pas d'ordre (au moins égal à)  $p+1$ . Les deux définitions précédentes cohabitent dans la littérature.

La notion d'ordre d'une méthode numérique à un pas permet de caractériser la précision qu'elle est capable d'atteindre.

Lemme: Soit  $p \in \mathbb{N}^*$ . Si la méthode numérique à un pas est d'ordre (au moins égal à)  $p$  et stable, alors, il existe des nombres positifs  $C$  et  $S$  tels que, quelle que soit la solution  $y \in C^{p+1}([t_0, t_0+T], \mathbb{R})$  de l'équation différentielle:

$$\forall t \in (t_0, t_0+T), y'(t) = f(t, y(t)),$$

les valeurs approchées  $(y_n)_{0 \leq n \leq N}$  données par la méthode numérique satisfont:

$$\max_{0 \leq n \leq N} |y(t_n) - y_n| \leq S (|y(t_0) - y_0| + C h_{\max}^p),$$

$$\text{où } h_{\max} = \max_{0 \leq n \leq N-1} h_n.$$

Preuve:

La stabilité de la méthode garantit l'existence d'un nombre positif  $S$  tel que:

$$\max_{0 \leq n \leq N} |y(t_n) - y_n| \leq S (|y(t_0) - y_0| + \sum_{m=0}^{n-1} |e_m|).$$

Le lemme découle alors de l'estimation précédente de la somme des erreurs de consistance  $(e_m)_{0 \leq m \leq n-1}$ .

En pratique, l'erreur initiale  $|y(t_0) - y_0|$  est marginale. L'erreur de la méthode numérique est donc gouvernée par son ordre  $p$  à travers le terme  $C h_{\max}^p$ . itmilieren cet ordre permet d'augmenter la précision de



La méthode.

Afin de caractériser l'ordre d'une méthode numérique, nous introduisons les notions suivantes.

Définitions: Soit  $p \in \mathbb{N}^*$ . Supposons que la fonction  $f$  soit de classe  $\mathcal{C}^p$  sur  $I \times \mathbb{R}$ . Les fonctions  $(f^{(q)})_{0 \leq q \leq p}$  sont définies sur l'ensemble  $I \times \mathbb{R}$  par les formules de récurrence:

$$\forall (t, y) \in I \times \mathbb{R}, \begin{cases} f^{(0)}(t, y) = f(t, y) \\ \forall 0 \leq q \leq p-1, f^{(q+1)}(t, y) = \partial_t f^{(q)}(t, y) + \partial_y f^{(q)}(t, y) \end{cases}$$

Les fonctions  $(f^{(q)})_{0 \leq q \leq p}$  sont définies de manière à satisfaire aux identités suivantes.

Lemme: Soit  $p \in \mathbb{N}^*$ . Supposons que la fonction  $f$  soit de classe  $\mathcal{C}^p$  sur  $I \times \mathbb{R}$  et considérons une solution  $y \in \mathcal{C}^{p+1}([t_0, t_0+T], \mathbb{R})$  de l'équation différentielle:

$$\forall t \in [t_0, t_0+T], y'(t) = f(t, y(t)).$$

Les dérivées successives de la fonction  $y$  satisfont alors les formules:

$$\forall 0 \leq q \leq p, \forall t \in [t_0, t_0+T], y^{(q+1)}(t) = f^{(q)}(t, y(t)).$$

Preuve:

Immédiate par récurrence sur l'entier  $0 \leq q \leq p$ .

Nous pouvons alors établir la caractérisation suivante de l'ordre d'une méthode numérique à un pas.

Théorème: Soit  $p \in \mathbb{N}^*$ . Supposons que la fonction  $f$  est de classe  $\mathcal{C}^p$  sur  $I \times \mathbb{R}$  et que la fonction  $\Phi$  est de classe  $\mathcal{C}^p$  sur  $(t_0, t_0+T) \times \mathbb{R}^2$ . La méthode numérique à un pas est d'ordre (au moins égal à)  $p$  si et seulement si:

$$\forall 0 \leq q \leq p-1, \forall (t, y) \in (t_0, t_0+T) \times \mathbb{R}, \partial_{tt}^q \Phi(t, y, 0) = \frac{1}{h^2} f^{(q)}(t, y).$$

Preuve:

Soit  $y \in \mathcal{C}^{p+2}([t_0, t_0+T], \mathbb{R})$  une solution de l'équation différentielle:

$$\forall t \in [t_0, t_0+T], y'(t) = f(t, y(t)).$$

Pour  $0 \leq n \leq N-1$ , l'erreur de consistance  $e_n$  au temps  $t_n$  relative à la solution  $y$  est donnée par:

$$e_n = y(t_{n+1}) - y(t_n) - h_n \Phi(t_n, y(t_n), h_n).$$

Introduisons les fonctions  $(\theta_q)_{0 \leq q \leq p-1}$  définies par les formules:

$$\forall 0 \leq q \leq p-1, \forall (t, y) \in [t_0, t_0+T] \times \mathbb{R}, \theta_q(t, y) = \frac{f^{(q)}(t, y)}{(q+2)!} - \frac{\partial^q \Phi(t, y, 0)}{q!}.$$

Comme  $y(t_{n+1}) = y(t_n + h_n)$ , la formule de Taylor avec reste intégral permet d'établir que:

$$e_n = \sum_{q=0}^{p-1} h_n^{q+2} \left[ \frac{y^{(q+2)}(t_n)}{(q+2)!} - \frac{\partial^q \Phi(t_n, y(t_n), 0)}{q!} \right] + O_{h_n \rightarrow 0}(h_n^{p+2}),$$

De sorte que, par le lemme précédent,

$$e_n = \sum_{q=0}^{p-1} h_n^{q+2} \theta_q(t_n, y(t_n)) + O_{h_n \rightarrow 0}(h_n^{p+2}).$$

Les identités du théorème, nous savons que:

$$\forall 0 \leq q \leq p-1, \forall (t, y) \in [t_0, t_0+T] \times \mathbb{R}, \theta_q(t, y) = 0,$$

De sorte qu'il existe un nombre positif  $C$  tel que:

$$\forall 0 \leq n \leq N-1, |e_n| \leq C h_n^{p+2} \leq C h_n h_{\max}^p.$$

Il s'ensuit que:

$$\sum_{n=0}^{N-1} |e_n| \leq C T h_{\max}^p,$$

Et la méthode numérique est bien d'ordre (au moins égal à)  $p$ .

Réciproquement, il existe un nombre positif  $C$  tel que:

$$\sum_{n=0}^{N-1} \left| \sum_{q=0}^{p-1} h_n^{q+2} \theta_q(t_n, y(t_n)) + O_{h_n \rightarrow 0}(h_n^{p+2}) \right| \leq C h_{\max}^p.$$

La preuve des identités du théorème est alors par récurrence sur l'entier

$0 \leq q \leq p-1$ . Au rang  $q=0$ , l'inégalité précédente conduit au fait

que:

$$\sum_{n=0}^{N-1} |h_n \theta_0(t_n, y(t_n))| \leq C h_{\max}.$$

de la limite  $h_{\max} \rightarrow 0$ , le théorème sur les sommes de Riemann

garantit que:

$$\int_{t_0}^{t_0+T} |\theta_0(t, y(t))| dt = 0,$$

puis, par continuité de la fonction  $t \mapsto \theta_0(t, y(t))$ , il vient:

$$\forall t \in [t_0, t_0+T], \theta_0(t, y(t)) = 0.$$

Cette identité est vraie pour toute solution  $y$  de l'équation différentielle considérée. Le théorème de Cauchy-Lipschitz assure alors que:

$$\forall (t, y) \in (t_0, t_0 + T) \times \mathbb{R}, \varphi_0(t, y) = 0,$$

soit l'identité:

$$\forall (t, y) \in (t_0, t_0 + T) \times \mathbb{R}, \mathbb{I}(t, y, 0) = f^{(0)}(t, y).$$

Il s'ensuit alors que:

$$\sum_{n=0}^{p-2} \left| \sum_{q=2}^{p-2} h_n^{q+1} \varphi_q(t_n, y(t_n)) + O(h_n^{p+2}) \right| \leq C h_{\max}^p,$$

Et la même preuve s'applique aux fonctions  $\varphi_2$ , puis  $\varphi_3$ , et jusqu'à  $\varphi_p$ . Ceci achève la preuve du théorème.

Nous déduisons du théorème précédent la corollaire suivant.

Corollaire: Supposons que la fonction  $f$  est de classe  $C^2$  sur  $I \times \mathbb{R}$  et que la fonction  $\mathbb{I}$  est de classe  $C^2$  sur  $(t_0, t_0 + T) \times \mathbb{R}^2$ . La méthode numérique à un pas est d'ordre (au moins égal à) 1 si et seulement si elle est consistante.

Preuve:

Les deux propriétés sont en effet équivalentes au fait que:

$$\forall (t, y) \in (t_0, t_0 + T) \times \mathbb{R}, \mathbb{I}(t, y, 0) = f(t, y).$$

Exemples: (i) Méthode d'Euler explicite

Cette méthode est d'ordre exactement égal à 1 puisque:

$$\forall (t, y) \in (t_0, t_0 + T) \times \mathbb{R}, \mathbb{I}(t, y, 0) = f(t, y)$$

$$\text{et } \left\{ \begin{array}{l} \partial_h \mathbb{I}(t, y, 0) = 0 \neq \frac{1}{2} (\partial_t f(t, y) + \partial_y f(t, y) f(t, y)) \\ \text{(en général)} \end{array} \right.$$

(ii) Méthode d'Euler implicite

Cette méthode est d'ordre exactement égal à 1 puisque:

$$\forall (t, y) \in (t_0, t_0 + T) \times \mathbb{R}, \mathbb{I}(t, y, 0) = f(t, y)$$

$$\text{et } \left\{ \begin{array}{l} \partial_h \mathbb{I}(t, y, 0) = \partial_t f(t, y) + \partial_y f(t, y) f(t, y) \\ \neq \frac{1}{2} (\partial_t f(t, y) + \partial_y f(t, y) f(t, y)) \text{ (en général)} \end{array} \right.$$

### (iii) Méthode du point milieu

Cette méthode est d'ordre exactement égal à 2 puisque:

$$\forall (t, y) \in [t_0, t_0 + \tau] \times \mathbb{R}^d, \left\{ \begin{array}{l} \Phi(t, y, 0) = f(t, y) \\ \text{et} \left\{ \begin{array}{l} \partial_t \Phi(t, y, 0) = \frac{1}{2} \left[ \partial_t f(t, y) + f(t, y) \partial_y f(t, y) \right] \\ \partial_t^2 \Phi(t, y, 0) = \frac{1}{4} \left[ \partial_t^2 f(t, y) + 2 f(t, y) \partial_t^2 f(t, y) + \right. \\ \left. f(t, y)^2 \partial_y^2 f(t, y) \right] \\ \neq \frac{1}{2} \left[ \partial_t^2 f(t, y) + \partial_t f(t, y) \partial_y f(t, y) + 2 f(t, y) \right. \\ \left. \partial_t^2 f(t, y) + f(t, y) (\partial_y f(t, y))^2 + f(t, y) \partial_y^2 f(t, y) \right] \end{array} \right. \end{array} \right.$$

## 6. Influence des erreurs d'arrondis

Considérons la méthode numérique à un pas associée à la fonction  $\Phi \in \mathcal{C}^2([t_0, t_0 + \tau] \times \mathbb{R}^d; \mathbb{R}^d)$  et supposons qu'elle est d'ordre (au moins égal à)  $p \in \mathbb{N}^*$ .

Rappelons que cette méthode revient à calculer les valeurs approchées  $(y_m)_{0 \leq m \leq N}$  aux temps  $(t_m)_{0 \leq m \leq N}$  données par les formules:

$$\forall m \in [0, N-1], y_{m+1} = y_m + h_m \Phi(t_m, y_m, h_m).$$

En pratique, ce calcul induit des erreurs d'arrondis  $(\epsilon_m)_{0 \leq m \leq N-1}$  sur la valeur de  $(\Phi(t_m, y_m, h_m))_{m \in [0, N-1]}$ , et  $(\delta_m)_{0 \leq m \leq N}$  sur la valeur de  $(y_{m+1})_{0 \leq m \leq N-1}$  de sorte que les valeurs approchées  $(\tilde{y}_m)_{0 \leq m \leq N}$  satisfont:

$$\forall m \in [0, N-1], \tilde{y}_{m+1} = \tilde{y}_m + h_m \Phi(t_m, \tilde{y}_m, h_m) + h_m \epsilon_m + \delta_m.$$

Nous avons alors l'estimation suivante de l'erreur entre les valeurs approchées  $\tilde{y}_m$  et  $y_m$ .

Lemme: Supposons que la méthode numérique est stable et qu'il existe des nombres positifs  $\rho$  et  $\sigma$  tels que:

$$\forall 0 \leq m \leq N-1, |\epsilon_m| \leq \rho \text{ et } |\delta_m| \leq \sigma.$$

Il existe un nombre positif  $\xi$  tel que:

$$\max_{0 \leq m \leq N} |\tilde{y}_m - y_m| \leq \xi \left[ |\tilde{y}_0 - y_0| + T \rho + N \sigma \right].$$

Preuve:

La stabilité de la méthode fournit l'existence d'un nombre positif  $S$  tel que

$$\begin{aligned} \max_{0 \leq m \leq N} |\tilde{y}_m - y_m| &\leq S \left( |\tilde{y}_0 - y_0| + \sum_{n=0}^{m-1} (|h_n p_n| + |b_n|) \right) \\ &\leq S (|\tilde{y}_0 - y_0| + T p + N b) \end{aligned}$$

Lorsque la méthode numérique est stable et d'ordre (au moins égal à)  $p$ , l'erreur totale commise par rapport aux valeurs exactes  $(y(t_m))_{0 \leq m \leq N}$  de la solution est donc égale à :

$$\begin{aligned} \max_{0 \leq m \leq N} |\tilde{y}_m - y(t_m)| &\leq \max_{0 \leq m \leq N} |\tilde{y}_m - y_m| + \max_{0 \leq m \leq N} |y_m - y(t_m)| \\ &\leq S (|\tilde{y}_0 - y_0| + |y_0 - y(t_0)|) + T p + N b + C h_{\max}^p \end{aligned}$$

En supposant de plus que le pas est constant, c'est-à-dire que :

$$\forall 0 \leq m \leq N-1, h_m = \frac{T}{N} = h_{\max},$$

il vient :

$$\max_{0 \leq m \leq N} |\tilde{y}_m - y(t_m)| \leq S (|\tilde{y}_0 - y_0| + |y_0 - y(t_0)| + T p) + \frac{S T b}{h_{\max}} + S C h_{\max}^p$$

L'erreur commise est ainsi minimale lorsque :

$$h_{\max} = \left( \frac{T b}{p C} \right)^{\frac{1}{p+2}}$$

Elle est alors de l'ordre de :

$$\max_{0 \leq m \leq N} |\tilde{y}_m - y(t_m)| \leq S (|\tilde{y}_0 - y_0| + |y_0 - y(t_0)| + T p + C' (T b)^{\frac{p}{p+2}})$$

Cette limitation sur le choix du pas maximal  $h_{\max}$  est cohérente, d'une part, avec le fait que les erreurs d'arrondis limitent la précision possible, et d'autre part, avec la propriété que plus le pas maximal  $h_{\max}$  est petit, plus l'entier  $N$  est grand, et plus les erreurs d'arrondis s'additionnent.

En pratique, choisir un pas maximal  $h_{\max}$  très petit ne garantit donc pas une précision plus grande sur la solution approchée en raison des erreurs d'arrondis. Il existe des méthodes de contrôle de choix du pas afin d'optimiser à la fois le temps de calcul et la précision de la solution approchée. Nous renvoyons à l'ouvrage de Michel Crouzeix et Alain Dujardin "Analyse numérique des équations différentielles" pour de plus amples détails sur ce sujet.

## II Méthodes de Runge - Kutta

### 1. Principe et exemples

Considérons une solution  $y \in \mathcal{C}^1([t_0, t_0+T], \mathbb{R})$  de l'équation différentielle :

$$\forall t \in [t_0, t_0+T], \quad y'(t) = f(t, y(t)).$$

Étant donné une subdivision  $t_0 < t_1 < \dots < t_N = t_0+T$  de l'intervalle  $[t_0, t_0+T]$ , nous cherchons à déterminer des valeurs approchées  $(y_n)_{0 \leq n \leq N}$  des valeurs exactes  $(y(t_n))_{0 \leq n \leq N}$  prises par la solution  $y$  aux temps  $(t_n)_{0 \leq n \leq N}$ . Dans ce but, nous pouvons utiliser l'identité :

$$\forall 0 \leq n \leq N-1, \quad y(t_{n+1}) = y(t_n) + \int_{t_n}^{t_{n+1}} f(t, y(t)) dt.$$

Calculer les valeurs approchées  $(y_n)_{0 \leq n \leq N}$  se réduit donc à calculer de manière approchée des intégrales. Il est alors naturel d'introduire des formules de quadrature pour le calcul de ces intégrales.

Plus précisément, nous pouvons introduire une formule de quadrature élémentaire :

$$\int_0^1 g(x) dx \approx \sum_{j=1}^q b_j g(c_j),$$

dans laquelle les points  $(c_j)_{1 \leq j \leq q}$  appartiennent au segment  $[0, 1]$ .

Obtons alors :

$$\forall 0 \leq n \leq N-1, \quad \forall 1 \leq j \leq q, \quad t_{n,j} = t_n + c_j h_n.$$

Le changement de variables  $t = t_n + \theta h_n$  conduit à la formule de quadrature :

$$\int_{t_n}^{t_{n+1}} f(t, y(t)) dt \approx \sum_{j=1}^q b_j h_n f(t_{n,j}, y(t_{n,j})).$$

Il reste encore à déterminer la valeur des nombres  $(y(t_{n,j}))_{1 \leq j \leq q}$ .

Rappelons que ces nombres satisfont aux identités :

$$\forall 1 \leq j \leq q, \quad y(t_{n,j}) = \int_{t_n}^{t_{n,j}} f(t, y(t)) dt.$$

Nous pouvons donc utiliser de nouvelles formules de quadrature élémentaires pour le calcul de valeurs approchées  $(y_{n,j})_{1 \leq j \leq q}$  des valeurs exactes  $(y(t_{n,j}))_{1 \leq j \leq q}$ .

suite à renumérotage, nous pouvons supposer que:

$$0 \leq c_1 \leq c_2 \leq \dots \leq c_q \leq 1.$$

Dans ce cas, nous pouvons définir des méthodes de quadrature élémentaires, de sorte que:

$$\forall 1 \leq j \leq q, \int_0^{c_j} g(t) dt \approx \sum_{i=1}^j a_{ji} g(c_i),$$

chacune de ces méthodes pouvant être distincte des autres. Nous avons alors le calcul approché:

$$\int_{t_n}^{t_{n+1}} f(t, y(t)) dt \approx \sum_{i=1}^j h_n a_{ji} f(t_{n,i}, y(t_{n,i})).$$

Ceci permet de définir les valeurs approchées  $(y_{n,i})_{1 \leq i \leq j}$ , puis de déduire à chaque étape la valeur approchée  $y_{n+1}$  de la valeur approchée  $y_n$ .

Ce principe conduit à la définition suivante des méthodes de Runge-Kutta.

Definition: Soit  $(c_j)_{1 \leq j \leq q} \in [0, 1]^q$  tel que:  $c_1 \leq \dots \leq c_q$ . Étant donné des coefficients réels  $(b_j)_{1 \leq j \leq q}$  et  $(a_{ji})_{1 \leq i \leq j \leq q}$ , la méthode de Runge-Kutta de résolution de l'équation différentielle:

$$\forall t \in [t_0, t_0 + T], y'(t) = f(t, y(t))$$

est définie par les formules de récurrence suivantes pour le calcul des valeurs approchées  $(y_n)_{0 \leq n \leq N}$  aux temps  $(t_n)_{0 \leq n \leq N}$ :

$$\forall n \in [0, N-1], \begin{cases} \forall 1 \leq j \leq q, t_{n,j} = t_n + c_j h_n \\ \text{et } \begin{cases} y_{n,i} = y_n + h_n \sum_{j=1}^i a_{ji} f(t_{n,j}, y_{n,j}) \\ y_{n+1} = y_n + h_n \sum_{j=1}^q b_j f(t_{n,j}, y_{n,j}) \end{cases} \end{cases}$$

Notation: Soit  $(c_j)_{1 \leq j \leq q} \in [0, 1]^q$  tel que:  $c_1 \leq \dots \leq c_q$ . La méthode de Runge-Kutta associée aux coefficients réels  $(b_j)_{1 \leq j \leq q}$  et  $(a_{ji})_{1 \leq i \leq j \leq q}$  est représentée par le tableau:

$c_1$	$a_{11}$	0	...	0
$c_q$	$a_{q2}$	...	$a_{qq}$	0
	$b_1$	...	$b_q$	

Exemples: (i) Méthode d'Euler explicite

C'est la méthode de Runge - Kutta associée au tableau:

$$\begin{array}{c|c} 0 & 0 \\ \hline & 1 \end{array}$$

Cette méthode correspond à l'application de la méthode de quadrature des rectangles à gauche.

(ii) Méthode d'Euler implicite

C'est la méthode de Runge - Kutta associée au tableau:

$$\begin{array}{c|c} 1 & 1 \\ \hline & 1 \end{array}$$

Cette méthode correspond à l'application de la méthode de quadrature des rectangles à droite.

(iii) Méthode du point milieu

C'est la méthode de Runge - Kutta associée au tableau:

$$\begin{array}{c|cc} 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \hline & 0 & 1 \end{array}$$

Cette méthode correspond à l'application de la méthode de quadrature du point milieu.

(iv) Méthode de Heun

C'est la méthode de Runge - Kutta associée au tableau:

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 1 & 1 & 0 \\ \hline & \frac{1}{2} & \frac{1}{2} \end{array}$$

Cette méthode correspond à l'application de la méthode de quadrature des trapèzes.

(v) Méthode de Runge - Kutta classique



C'est la méthode de Runge - Kutta associée au tableau :

0	0 0 0 0
$\frac{1}{2}$	$\frac{1}{2}$ 0 0 0
$\frac{1}{2}$	0 $\frac{1}{2}$ 0 0
1	0 0 1 0
	1 2 2 1
	6 6 6 6

Cette méthode correspond à l'application de la méthode de quadrature de Simpson.

Mais pouvons remarquer que les méthodes de Runge - Kutta sont explicites dès lors que les coefficients  $(a_{ji})_{2 \leq i \leq q, 1 \leq j \leq i-1}$  satisfont :

$$\forall 1 \leq j \leq q, a_{jj} = 0.$$

Mais nous limiterons dans la suite de ce cours à ce cas particulier, même si certains des résultats qui suivent s'étendent aux méthodes de Runge - Kutta implicites.

## 2. Convergence et ordre

Les méthodes de Runge - Kutta sont des méthodes numériques à un pas.

Théorème : La méthode de Runge - Kutta associée aux points  $(c_j)_{2 \leq j \leq q}$ , et aux coefficients réels  $(b_j)_{2 \leq j \leq q}$  et  $(a_{ji})_{2 \leq i < j \leq q}$  est la méthode numérique à un pas associée à la fonction  $\Phi \in \mathcal{C}^0([t_0, t_0 + T] \times \mathbb{R}^L, \mathbb{R})$  donnée par la formule :

$$\forall (t, y, h) \in [t_0, t_0 + T] \times \mathbb{R}^L, \Phi(t, y, h) = \sum_{j=2}^q b_j f(t + c_j h, y_j),$$

où les points  $(y_j)_{2 \leq j \leq q}$  sont donnés par les formules :

$$y_2 = y \text{ et } \forall 2 \leq j \leq q, y_j = y + h \sum_{i=2}^{j-1} a_{ji} f(t + c_i h, y_i).$$

Démonstration :

Immédiat.

Il est donc possible de leur appliquer les résultats généraux quant à la consistance, la stabilité et la convergence de ces méthodes.

Lemme: La méthode de Runge - Kutta associée aux points  $(c_j)_{1 \leq j \leq q}$ , et aux coefficients réels  $(b_j)_{1 \leq j \leq q}$  et  $(a_{ji})_{1 \leq i < j \leq q}$  est consistante si et seulement si :

$$\sum_{j=1}^q b_j = 1.$$

Preuve:

En effet, nous savons que :

$$\forall (t, y) \in [t_0, t_0 + T] \times \mathbb{R}, \quad \Phi(t, y, 0) = \left( \sum_{j=1}^q b_j \right) f(t, y).$$

L'équivalence résulte alors de la condition de consistance des méthodes numériques à un pas.

En ce qui concerne la stabilité, nous pouvons établir le résultat suivant.

Lemme: Supposons que la fonction  $f$  est globalement lipschitzienne par rapport à  $y$  et soit qu'il existe un nombre positif  $K$  tel que :

$$\forall t \in [t_0, t_0 + T], \forall (y, y') \in \mathbb{R}^2, |f(t, y) - f(t, y')| \leq K |y - y'|$$

La méthode de Runge - Kutta associée aux points  $(c_j)_{1 \leq j \leq q}$  et aux coefficients réels  $(b_j)_{1 \leq j \leq q}$  et  $(a_{ji})_{1 \leq i < j \leq q}$  est stable dès lors que le pas maximal  $h_{\max}$  est borné par 1.

Preuve:

Il suffit de vérifier que la fonction  $\Phi$  est globalement lipschitzienne par rapport à la variable  $y \in \mathbb{R}$ . Soit  $(t, h) \in [t_0, t_0 + T] \times ]0, h_{\max}]$ . Alors pouvons calculer :

$$\forall (y, y') \in \mathbb{R}^2, |\Phi(t, y', h) - \Phi(t, y, h)| \leq K \sum_{j=1}^q |b_j| |y'_j - y_j|,$$

où

$$\forall 1 \leq j \leq q, y_j^{(v)} = y_j^{(v-1)} + h \sum_{i=1}^{j-1} a_{ji} f(t + c_i h, y_i).$$

du rang 1, nous avons :

$$|y'_2 - y_2| = |y'_1 - y_1|.$$

Supposons alors que:  $\forall 1 \leq l \leq j-1, |y'_l - y_l| \leq |y'_1 - y_1| \sum_{k=0}^{l-1} (K|h|_x)^k,$

où  $\alpha = \max_{1 \leq j \leq q} \sum_{i=1}^{j-1} |a_{ji}|$ . Au rang  $j$ , il vient alors:

$$\begin{aligned} |y'_j - y_j| &\leq |y'_1 - y_1| + K|h|_x \max_{1 \leq i \leq j-1} |y'_i - y_i| \\ &\leq |y'_1 - y_1| \sum_{i=0}^{j-1} (K|h|_x)^i, \end{aligned}$$

Ce qui achève la récurrence. Il s'ensuit que:

$$\forall (y, y') \in \mathbb{R}^2, |\mathbb{E}(t, y', h) - \mathbb{E}(t, y, h)| \leq K \sum_{j=1}^q |b_j| \sum_{i=0}^{j-1} (K|h|_x)^i |y - y_1| \leq \Lambda |y' - y|,$$

Où  $\Lambda = K \sum_{j=1}^q |b_j| \sum_{i=0}^{j-1} (K|h|_x)^i$ . La méthode de Runge-Kutta est donc stable sous la condition  $h_{\max} \leq 1$ .

Remarque: La constante de stabilité de la méthode de Runge-Kutta dépend du pas maximal  $h_{\max}$ . D'après le lemme précédent, elle est de l'ordre de  $e^{\Lambda T}$ , où  $\Lambda = K \sum_{j=1}^q |b_j| \sum_{i=0}^{j-1} (K h_{\max})^i$ . Lorsque le pas maximal  $h_{\max}$  tend vers 0, nous vérifions que:

$$\Lambda \xrightarrow{h_{\max} \rightarrow 0} K \sum_{j=1}^q |b_j|$$

Lorsque les coefficients  $(b_j)_{1 \leq j \leq q}$  sont positifs, ce qui est naturel en pratique, et sous la condition de conservation  $\sum_{j=1}^q b_j = 1$ , nous obtenons:

$$\Lambda \xrightarrow{h_{\max} \rightarrow 0} K.$$

Cette limite est la meilleure possible puisqu'elle correspond au cas explicite d'une fonction  $f$  linéaire. Les méthodes de Runge-Kutta présentent donc le précieux avantage d'être très stables (au moins lorsque le pas maximal  $h_{\max}$  est suffisamment petit).

Les deux lemmes précédents permettent d'établir la convergence d'une méthode de Runge-Kutta.

Lemme: Supposons que la fonction  $f$  est globalement lipschitzienne par rapport à  $y \in \mathbb{R}$ , soit qu'il existe un nombre positif  $K$  tel que:

$\forall t \in [t_0, t_0 + T], \forall (y, y') \in \mathbb{R}^2, |f(t, y) - f(t, y')| \leq K |y - y'|$ ,  
 et considérons des points  $(c_j)_{1 \leq j \leq q} \in [0, 1]^q$ , et des coefficients réels  $(b_j)_{1 \leq j \leq q}$  et  $(a_{ji})_{1 \leq i < j \leq q}$  tels que:  
 $\sum_{j=1}^q b_j = 1$ .  
 La méthode de Runge - Kutta associée est convergente.

Preuve:

Immédiat.

Exemples: Les méthodes de Heun et de Runge - Kutta classiques sont convergentes.

Il est enfin permis d'étendre les résultats généraux sur l'ordre des méthodes numériques à un pas aux méthodes de Runge - Kutta.

Lémmes: La méthode de Runge - Kutta associée aux points  $(c_j)_{1 \leq j \leq q}$  et aux coefficients réels  $(b_j)_{1 \leq j \leq q}$  et  $(a_{ji})_{1 \leq i < j \leq q}$  sont:

(i) d'ordre (au moins égal à)  $\alpha$  si et seulement si:

$$\sum_{j=1}^q b_j = 1$$

(ii) d'ordre (au moins égal à)  $\alpha$  si et seulement si:

$$\sum_{j=1}^q b_j = 1, \sum_{j=1}^q b_j c_j = \sum_{j=1}^q \sum_{i=1}^{j-1} b_j a_{ji} = 1$$

Preuve:

(i) est une conséquence directe du résultat de consistance précédent. En ce qui

concerne (ii), rappelons que:

$$\forall t \in [t_0, t_0 + T], \forall (y, h) \in \mathbb{R}^2, \Phi(t, y, h) = \sum_{j=1}^q b_j f(t + c_j h, y_j),$$

où  $\forall 1 \leq j \leq q, y_j = y + h \sum_{i=1}^{j-1} a_{ji} f(t + c_i h, y_i)$ . Il vient donc:

$$d_t \Phi(t, y, h) = \sum_{j=1}^q b_j \left( c_j \frac{d_t}{dt} f(t + c_j h, y_j) + \frac{dy}{dt} f(t + c_j h, y_j) \right) \times \sum_{i=1}^{j-1} (a_{ji} f(t + c_i h, y_i) + h a_{ji} \frac{d_t}{dt} f(t + c_i h, y_i))$$

De sorte que:

$$d_t \Phi(t, y, 0) = \sum_{j=1}^q b_j \left( c_j \frac{d_t}{dt} f(t, y) + \frac{dy}{dt} f(t, y) \sum_{i=1}^{j-1} a_{ji} f(t, y) \right),$$

ce qui conduit à l'équivalence de (ii) par le résultat général sur l'ordre d'une méthode numérique à un pas.

En pratique, les méthodes de quadrature élémentaires sous-jacentes aux méthodes de Runge-Kutta sont au moins d'ordre 2, ce qui revient à dire que:

$$\sum_{j=1}^q b_j = 1 \text{ et } \forall 1 \leq j \leq q, c_j = \sum_{i=1}^{j-1} a_{ji}.$$

Dans ce cas, il est possible d'étendre le théorème précédent de la façon suivante.

Théorème: Considérons des points  $(c_j)_{1 \leq j \leq q} \in (0, 1]^q$  et des coefficients tels

$$(b_j)_{1 \leq j \leq q} \text{ et } (a_{ji})_{1 \leq i < j \leq q} \text{ tels que:}$$

$$\sum_{j=1}^q b_j = 1 \text{ et } \forall 1 \leq j \leq q, c_j = \sum_{i=1}^{j-1} a_{ji}.$$

La méthode de Runge-Kutta associée est d'ordre (au moins) égal à  $\alpha$ . Elle est d'ordre (au moins) égal à:

$$(i) 2 \text{ si et seulement si: } \sum_{j=1}^q b_j c_j = \frac{1}{2}.$$

$$(ii) 3 \text{ si et seulement si: } \sum_{j=1}^q b_j c_j^2 = 3 \sum_{j=1}^q b_j c_j = 6 \sum_{j=1}^q \sum_{i=1}^{j-1} a_{ji} b_j c_j = 1.$$

Preuve:

Elle est similaire à celle du théorème précédent.

Exemples: (i) Méthode de Heun

La méthode de Heun est d'ordre exactement égal à 2.

(ii) Méthode de Runge-Kutta classique

La méthode de Runge-Kutta classique est d'ordre exactement égal à 4. Elle présente l'avantage d'un ordre élevé, d'une complexité algorithmique faible et d'une grande stabilité, raisons pour lesquelles elle est souvent employée en pratique.

### III Difficultés numériques de la résolution approchée des équations différentielles ordinaires

Considérons le problème de Cauchy:

$$\begin{cases} \forall t \in [t_0, t_0 + T], y'(t) = f(t, y(t)) \\ \text{et } y(t_0) = y_0 \end{cases}$$

où  $y_0 \in \mathbb{R}$  et  $f: [t_0, t_0 + T] \times \mathbb{R} \rightarrow \mathbb{R}$ . Afin de résoudre ce problème, nous pouvons appliquer les méthodes numériques à un pas que nous avons introduits précédemment, et les résultats de convergence que nous avons démontrés garantiront alors que nous obtiendrons des solutions approchées raisonnables de la solution de ce problème de Cauchy.

Ces résultats reposent néanmoins sur des hypothèses précises sur la fonction  $f$ , lesquelles ne sont pas toujours vérifiées. De plus, même dans le cas où ces hypothèses sont vérifiées, les estimations d'erreur fournies par les résultats de convergence précédents peuvent être très mauvaises, parce que les constantes (par exemple, de stabilité) qui apparaissent dans ces estimations sont très grandes.

La mise en pratique des méthodes numériques à un pas se heurte ainsi à un certain nombre de difficultés numériques que nous allons maintenant décrire.

#### 1. Problèmes mathématiquement et numériquement bien posés

Definition: Soit  $f \in C^0([t_0, t_0 + T] \times \mathbb{R}, \mathbb{R})$ . Le problème de Cauchy

$$\begin{cases} \forall t \in [t_0, t_0 + T], y'(t) = f(t, y(t)), \\ \text{et } y(t_0) = y_0 \end{cases}$$

est mathématiquement bien posé si et seulement si:

- (i) il existe une unique solution  $y \in C^2([t_0, t_0 + T], \mathbb{R})$  pour

chaque donnée initiale  $y_0 \in \mathbb{R}$ .

(ii) la solution  $y$  dépend de manière continue de la donnée initiale, autrement dit, l'application  $y_0 \mapsto y(t)$  est continue pour chaque  $t \in [t_0, t_0 + T]$ .

Le théorème de Cauchy - Lipschitz (global) garantit qu'un problème de Cauchy est mathématiquement bien posé.

Théorème de Cauchy - Lipschitz (global): Si la fonction  $f \in C^0([t_0, t_0 + T] \times \mathbb{R}, \mathbb{R})$  est globalement lipschitzienne par rapport à la variable  $y \in \mathbb{R}$ , c'est-à-dire s'il existe un nombre positif  $K$  tel que:

$$\forall t \in [t_0, t_0 + T], \forall (y, y') \in \mathbb{R}^2, |f(t, y') - f(t, y)| \leq K|y' - y|,$$

alors, le problème de Cauchy:

$$\begin{cases} \forall t \in [t_0, t_0 + T], y'(t) = f(t, y(t)), \\ y(t_0) = y_0, \end{cases}$$

est mathématiquement bien posé.

Preuve:

Voir ci-dessus.

Certains problèmes de Cauchy ne sont pas mathématiquement bien posés.

Exemple: Le problème de Cauchy

$$\begin{cases} \forall t \in \mathbb{R}_+, y'(t) = 2\sqrt{|y(t)|} \\ y(0) = 0 \end{cases}$$

admet les solutions:

(i)  $\forall t \in \mathbb{R}_+, y(t) = 0$

(ii)  $\forall a \in \mathbb{R}_+, \forall t \in \mathbb{R}_+, y(t) = 0$  si  $0 \leq t \leq a$ ,  
et  $(t-a)^2$  sinon,

Il n'y a pas d'unicité de la solution.

Définition: Soit  $f \in C^0([t_0, t_0 + \tau] \times \mathbb{R}, \mathbb{R})$ . Le problème de Cauchy

$$\begin{cases} \forall t \in [t_0, t_0 + \tau], y'(t) = f(t, y(t)), \\ y(t_0) = y_0 \end{cases}$$

est numériquement bien posé si et seulement si la continuité d'une solution  $y$  par rapport à la donnée initiale  $y_0$  est suffisamment bonne pour que la solution ne soit pas perturbée par une erreur initiale ou une erreur d'arrondi sur la fonction  $f$ .

Cette définition n'est pas très précise. La continuité de la solution  $y$  par rapport à la donnée initiale  $y_0$  se mesure par la taille des constantes qui apparaissent dans le lemme de Gronwall. Ces constantes doivent être suffisamment petites par rapport à la précision des calculs. En particulier, cette définition dépend de la précision possible.

Il existe des problèmes de Cauchy qui sont mathématiquement bien posés, mais numériquement mal posés.

Exemple: Le problème de Cauchy:

$$\begin{cases} \forall t \in (0, 10), y'(t) = 3y(t) - 1, \\ y(0) = \frac{1}{2}, \end{cases}$$

à pour solution:

$$\forall t \in (0, 10), y(t) = t + \frac{1}{3}.$$

Pour le théorème de Cauchy-Lipschitz (global), il est mathématiquement bien posé. Cependant, la solution associée à la donnée initiale  $y_\varepsilon(0) = \frac{1}{2} + \varepsilon$  pour  $\varepsilon \geq 0$  est égale à:

$$\forall t \in (0, 10), y_\varepsilon(t) = t + \frac{1}{3} + \varepsilon e^{3t}.$$

On observe en particulier que:

$$y_\varepsilon(10) - y(10) = \varepsilon e^{30} \approx 10^{23} \varepsilon$$

Le problème est donc numériquement mal posé pour une précision des calculs de l'ordre de  $10^{-20}$ , mais redonne une



numériquement bien posé pour une précision de l'ordre de  $10^{-20}$ .

Notons que il n'est pas possible de résoudre numériquement un problème de Cauchy numériquement mal posé puisque les erreurs d'arrondi sont alors incontrôlables.

## 2. Problèmes bien conditionnés et problèmes raides

Même lorsqu'un problème de Cauchy est mathématiquement et numériquement bien posé, il peut s'avérer difficile d'en déterminer une solution approchée en un temps raisonnable.

Definition: Soit  $f \in C^0([t_0, t_0+T] \times \mathbb{R}, \mathbb{R})$ . Le problème de Cauchy

$$\begin{cases} \forall t \in [t_0, t_0+T], y'(t) = f(t, y(t)), \\ \text{et } y(t_0) = y_0, \end{cases}$$

est bien conditionné si les méthodes numériques classiques peuvent en donner la solution en un temps raisonnable.

Cette définition n'est pas très précise. Les méthodes numériques classiques dont il est question se réduisent en général aux méthodes numériques à un pas explicite (et parfois implicite). Les temps de calcul dépendent du matériel informatique utilisé.

Exemple: Le problème de Cauchy:

$$\begin{cases} \forall t \in (0, 1], y'(t) = -250 y(t) + 30, \\ \text{et } y(0) = \frac{1}{5}, \end{cases}$$

a pour unique solution:

$$\forall t \in (0, 1], y(t) = \frac{1}{5}.$$

Pour le théorème de Cauchy-Lipschitz (global), il est mathématiquement bien posé. De plus, l'unique solution  $y_\varepsilon$  associée à la donnée initiale  $y_\varepsilon(0) = \frac{1}{5} + \varepsilon$  pour  $\varepsilon > 0$  est égale à:

$$\forall t \in [0, 2], y_2(t) = \frac{1}{5} + 2e^{-250t},$$

De sorte que ce problème de Cauchy est numériquement bien posé. Cependant, la méthode d'Euler explicite avec un pas constant  $h = \frac{2}{N}$ , où  $N \in \mathbb{N}^*$ , fournit des valeurs approchées  $(y_n)_{0 \leq n \leq N}$  de la solution  $y$  aux temps  $(t_n = \frac{2n}{N})_{0 \leq n \leq N}$  données par la formule de récurrence:

$$\forall 0 \leq n \leq N-1, y_{n+1} = y_n + \frac{2}{N} (-250 y_n + 30),$$

soit par les formules:

$$\forall 0 \leq n \leq N, y_n = \frac{1}{5} + \left(2 - \frac{250}{N}\right)^n \left(y_0 - \frac{1}{5}\right),$$

L'erreur ainsi commise est de l'ordre de:

$$\max_{0 \leq n \leq N} |y_n - y(t_n)| = \left(\frac{250}{N} - 2\right)^N \left|y_0 - \frac{1}{5}\right|$$

Lorsque  $N \leq 75$ , quantité qui peut s'avérer très grande. Lorsque  $N$  est par exemple de l'ordre de 50. Les problèmes de Cauchy

sont ainsi mal conditionnés pour la méthode d'Euler explicite.

Au contraire, la méthode d'Euler implicite avec le même pas constant  $h = \frac{2}{N}$ , fournit des valeurs approchées  $(y_n)_{0 \leq n \leq N}$  de la solution  $y$  aux temps  $(t_n)_{0 \leq n \leq N}$  données par la formule de récurrence:

$$\forall 0 \leq n \leq N-1, y_{n+1} = y_n + \frac{2}{N} (-250 y_{n+1} + 30),$$

soit par les formules:

$$\forall 0 \leq n \leq N, y_n = \frac{1}{5} + \frac{y_0 - \frac{1}{5}}{\left(1 + \frac{250}{N}\right)^n}.$$

L'erreur ainsi commise est de l'ordre de:

$$\max_{0 \leq n \leq N} |y_n - y(t_n)| = \left|y_0 - \frac{1}{5}\right|.$$

Le problème de Cauchy est donc bien conditionné pour la méthode d'Euler implicite.

Il faut donc qu'il existe des problèmes de Cauchy qui sont mathématiquement et numériquement bien posés, mais mal conditionnés pour toutes les méthodes numériques classiques. Ce sont les problèmes raides. Nous renvoyons à l'ouvrage de Michel Crouzeix et Alain Hignat "Analyse numérique des équations différentielles" pour de plus amples détails sur ce sujet.