

Erreurs d'arrondi de méthodes d'intégration numérique explicites

Journées FASTRELAX, LIP6 - UPMC

Sylvie Boldo¹ Florian Faissole¹ Alexandre Chapoutot²

¹Inria - LRI, Univ. Paris-Sud et CNRS - Univ. Paris-Saclay

²U2IS, ÉNSTA ParisTech

Table des matières

- 1 Introduction
- 2 Méthodes et systèmes étudiés
- 3 Erreurs d'arrondi
- 4 Conclusion et Perspectives

Équations différentielles ordinaires

Équations différentielles ordinaires (ÉDOs) :

$$F(y, y', y'', \dots, y^{(n)}, t) = 0.$$

Utilisation : aérospatial, modèles de populations, économie, astrophysique, microbiologie ...



Résolution **exacte** souvent difficile \Rightarrow **méthodes numériques**.

Résolution numérique d'ÉDOs

Méthode numérique \implies solution **approchée discrétisée** :

$$y_0, y_1, \dots, y_N \quad (\forall 0 \leq i \leq N, y_i \simeq y(t_0 + ih)).$$

Méthodes explicites :

$$y_n \text{ calculé à partir de } y_{n-1}, y_{n-2}, \dots, y_0.$$

Implémentation (**arithmétique FP**) \implies deux sources d'**erreurs** :

- erreur de **méthode** ;
- erreur d'**arrondi**.

Analyse d'erreur : un état-de-l'art

Analyse d'erreurs sur méthodes numériques :

- résultat **proba.** : erreur proportionnelle à \sqrt{N} [Henrici, 1963] ;
- en pratique (RK implicite) : proportionnelle à N [Hairer, 2008] ;
- analyse par **intervalles** [Bouissou, Martel, 2006] ;
- **intégration numérique** (Newton-Cotes ...) [Fousse, 2006].

Analyse d'erreur : un état-de-l'art

Analyse d'erreurs sur méthodes numériques :

- résultat **proba.** : erreur proportionnelle à \sqrt{N} [Henrici, 1963] ;
- en pratique (RK implicite) : proportionnelle à N [Hairer, 2008] ;
- analyse par **intervalles** [Bouissou, Martel, 2006] ;
- **intégration numérique** (Newton-Cotes ...) [Fousse, 2006].

Notre approche :

- analyse en **grain fin** ;
- propriétés **mathématiques** des schémas (**stabilité**).

Table des matières

- 1 Introduction
- 2 Méthodes et systèmes étudiés**
- 3 Erreurs d'arrondi
- 4 Conclusion et Perspectives

Méthodes de Runge-Kutta étudiées

Équations différentielles linéaires :

$$\dot{y} = \lambda y$$

Méthodes de Runge-Kutta étudiées

Équations différentielles linéaires :

$$\dot{y} = \lambda y$$

Méthodes de Runge-Kutta :

$$y_{n+1} = y_n + h \sum_{i=1}^s b_i k_i, \quad k_i = \lambda \left(t_n + c_i h, y_n + h \sum_{j=1}^s a_{ij} k_j \right)$$

(a_{ij} , b_i , c_i consignés dans tableau de Butcher).

Méthodes de Runge-Kutta étudiées

Équations différentielles linéaires :

$$\dot{y} = \lambda y$$

Méthodes de Runge-Kutta :

$$y_{n+1} = y_n + h \sum_{i=1}^s b_i k_i, \quad k_i = \lambda \left(t_n + c_i h, y_n + h \sum_{j=1}^s a_{ij} k_j \right)$$

(a_{ij} , b_i , c_i consignés dans tableau de Butcher).

Méthode exacte :

Implémentation :

$$\begin{cases} y_0 \in \mathbb{C} \\ y_{n+1} = R(h, \lambda) y_n \end{cases}$$

Méthodes de Runge-Kutta étudiées

Équations différentielles linéaires :

$$\dot{y} = \lambda y$$

Méthodes de Runge-Kutta :

$$y_{n+1} = y_n + h \sum_{i=1}^s b_i k_i, \quad k_i = \lambda \left(t_n + c_i h, y_n + h \sum_{j=1}^s a_{ij} k_j \right)$$

(a_{ij} , b_i , c_i consignés dans tableau de Butcher).

Méthode exacte :

Implémentation :

$$\begin{cases} y_0 \in \mathbb{C} \\ y_{n+1} = R(h, \lambda) y_n \end{cases}$$

$$\begin{cases} \tilde{y}_0 \simeq y_0 \\ \widetilde{y_{n+1}} = \widetilde{R}(\tilde{h}, \tilde{\lambda}, \tilde{y}_n) \end{cases}$$

Méthodes de Runge-Kutta explicites classiques

Euler :

- $R(h, \lambda) = 1 + h\lambda$;
- $\tilde{R}(\tilde{h}, \tilde{\lambda}, \tilde{y}_n) = \tilde{y}_n \oplus \tilde{h} \otimes \tilde{\lambda} \otimes \tilde{y}_n$.

Méthodes de Runge-Kutta explicites classiques

Euler :

- $R(h, \lambda) = 1 + h\lambda$;
- $\tilde{R}(\tilde{h}, \tilde{\lambda}, \tilde{y}_n) = \tilde{y}_n \oplus \tilde{h} \otimes \tilde{\lambda} \otimes \tilde{y}_n$.

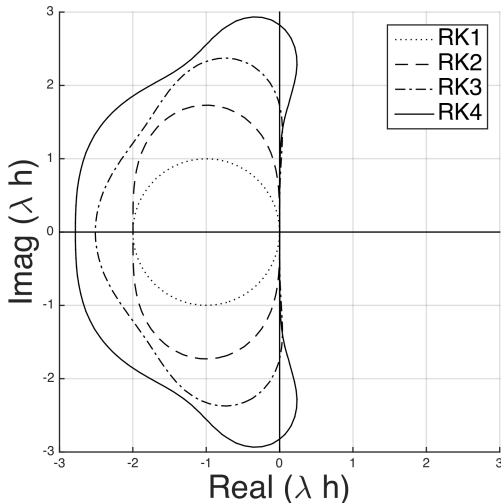
RK4 :

- $R(h, \lambda) = 1 + h\lambda + \frac{1}{2}(h\lambda)^2 + \frac{1}{6}(h\lambda)^3 + \frac{1}{24}(h\lambda)^4$;
- $\tilde{R}(\tilde{h}, \tilde{\lambda}, \tilde{y}_n) =$
 - $\left[\tilde{y}_n + \frac{\tilde{h}}{6}\tilde{\lambda}\tilde{y}_n + \frac{\tilde{h}}{3}\tilde{\lambda}\tilde{y}_n + \frac{\tilde{h}^2}{6}\tilde{\lambda}^2\tilde{y}_n + \frac{\tilde{h}}{3}\tilde{\lambda}\tilde{y}_n + \frac{\tilde{h}^2}{6}\tilde{\lambda}^2\tilde{y}_n + \frac{\tilde{h}^3}{12}\tilde{\lambda}^3\tilde{y}_n + \frac{\tilde{h}}{6}\tilde{\lambda}\tilde{y}_n + \frac{\tilde{h}^2}{6}\tilde{\lambda}^2\tilde{y}_n + \frac{\tilde{h}^3}{12}\tilde{\lambda}^3\tilde{y}_n + \frac{\tilde{h}^4}{24}\tilde{\lambda}^4\tilde{y}_n \right]$.

(+ de 60 opérations flottantes !)

Stabilité linéaire

Méthode de Runge-Kutta linéaire **stable** ssi $|R(h, \lambda)| < 1$:



Hypothèses de travail

- ÉDOs ;
 - du premier ordre,
 - linéaires,
 - inconnue à valeur dans \mathbb{R} (unidimensionnel),

Hypothèses de travail

- ÉDOs ;
 - du premier ordre,
 - linéaires,
 - inconnue à valeur dans \mathbb{R} (unidimensionnel),
- Méthodes ;
 - explicites,
 - à un pas constant,

Hypothèses de travail

- **ÉDOs** ;
 - du premier ordre,
 - linéaires,
 - inconnue à valeur dans \mathbb{R} (unidimensionnel),
- **Méthodes** ;
 - explicites,
 - à un pas constant,
- **Arithmétique à virgule flottante**.
 - ni underflow ni overflow,
 - radix 2 double précision ($u = 2^{-53}$),

Table des matières

- 1 Introduction
- 2 Méthodes et systèmes étudiés
- 3 Erreurs d'arrondi**
- 4 Conclusion et Perspectives

Erreurs locales et globales

Erreur locale :

$$\varepsilon_0 = |\tilde{y}_0 - y_0|$$

$$\forall n \in \mathbb{N}^*, \varepsilon_n = |\tilde{R}(\tilde{h}, \tilde{\lambda}, \widetilde{y_{n-1}}) - R(h, \lambda)\widetilde{y_{n-1}}|.$$

Erreur globale :

$$\forall n \in \mathbb{N}, E_n = \tilde{y}_n - y_n.$$

Erreurs locales et globales

Erreur locale :

$$\varepsilon_0 = |\widetilde{y}_0 - y_0|$$

$$\forall n \in \mathbb{N}^*, \varepsilon_n = |\widetilde{R}(\widetilde{h}, \widetilde{\lambda}, \widetilde{y_{n-1}}) - R(h, \lambda)\widetilde{y_{n-1}}|.$$

Erreur globale :

$$\forall n \in \mathbb{N}, E_n = \widetilde{y}_n - y_n.$$

Théorème 1 : Erreur globale absolue des schémas RK

Soit $C \in \mathbb{R}_+^*$. Supposons que $\forall n \in \mathbb{N}^*, \varepsilon_n \leq C|\widetilde{y_{n-1}}|$. Alors pour tout $n \in \mathbb{N}$,

$$|E_n| \leq (C + |R(h, \lambda)|)^n \left(\varepsilon_0 + n \frac{C|y_0|}{C + |R(h, \lambda)|} \right).$$

Erreurs relatives

Erreur relative :

$$\left| \frac{\tilde{y}_n - y_n}{y_n} \right| \leq \left(\frac{C + |R(h, \lambda)|}{|R(h, \lambda)|} \right)^n \left(\varepsilon_0 + n \frac{C|y_0|}{C + |R(h, \lambda)|} \right).$$

Erreurs relatives

Erreur relative :

$$\left| \frac{\tilde{y}_n - y_n}{y_n} \right| \leq \left(\frac{C + |R(h, \lambda)|}{|R(h, \lambda)|} \right)^n \left(\varepsilon_0 + n \frac{C|y_0|}{C + |R(h, \lambda)|} \right).$$

Si $C \ll |R(h, \lambda)|$, alors :

$$\left| \frac{\tilde{y}_n - y_n}{y_n} \right| \lesssim \varepsilon_0 + n \frac{C|y_0|}{|R(h, \lambda)|}.$$

En pratique (Euler, RK2, RK4) : $C \leq 200u$ et $200u \ll |R(h, \lambda)|$.

Erreurs locales

Lemme 2

Soit $y \in \mathbb{R}$. **Soit** $C_1, C_2 \in \mathbb{R}_+$. **Soit** $\alpha_1, \alpha_2 \in \mathbb{R}$. **Soit** $X_1, X_2 \in \mathbb{F}$ **t.q.** :
 $|X_1 - \alpha_1 y| \leq C_1 |y|$ **et** $|X_2 - \alpha_2| \leq C_2$.

Alors :

$$|X_1 \oplus (X_2 \otimes y) - (\alpha_1 + \alpha_2)y| \\ \leq |y| (C_1 + C_2 + u(|\alpha_1| + 2|\alpha_2| + C_1 + 2C_2) + u^2(C_2 + |\alpha_2|)).$$

Erreurs locales

Lemme 2

Soit $y \in \mathbb{R}$. *Soit* $C_1, C_2 \in \mathbb{R}_+$. *Soit* $\alpha_1, \alpha_2 \in \mathbb{R}$. *Soit* $X_1, X_2 \in \mathbb{F}$ *t.q.* :
 $|X_1 - \alpha_1 y| \leq C_1 |y|$ *et* $|X_2 - \alpha_2| \leq C_2$.

Alors :

$$|X_1 \oplus (X_2 \otimes y) - (\alpha_1 + \alpha_2)y| \\ \leq |y| (C_1 + C_2 + u(|\alpha_1| + 2|\alpha_2| + C_1 + 2C_2) + u^2(C_2 + |\alpha_2|)).$$

Exemple (erreur relative de la méthode d'Euler) :

$$\varepsilon_{n+1} = \left| \underbrace{\tilde{y}_n}_{X_1} \oplus \left(\underbrace{\tilde{h} \otimes \tilde{\lambda}}_{X_2} \otimes \underbrace{\tilde{y}_n}_y \right) - \left(\underbrace{1}_{\alpha_1} + \underbrace{h\lambda}_{\alpha_2} \right) \underbrace{\tilde{y}_n}_y \right|.$$

Erreurs locales

Lemme 2

Soit $y \in \mathbb{R}$. *Soit* $C_1, C_2 \in \mathbb{R}_+$. *Soit* $\alpha_1, \alpha_2 \in \mathbb{R}$. *Soit* $X_1, X_2 \in \mathbb{F}$ *t.q.* :
 $|X_1 - \alpha_1 y| \leq C_1 |y|$ *et* $|X_2 - \alpha_2| \leq C_2$.

Alors :

$$|X_1 \oplus (X_2 \otimes y) - (\alpha_1 + \alpha_2)y| \\ \leq |y| (C_1 + C_2 + u(|\alpha_1| + 2|\alpha_2| + C_1 + 2C_2) + u^2(C_2 + |\alpha_2|)).$$

Exemple (erreur relative de la méthode d'Euler) :

$$\varepsilon_{n+1} = \left| \underbrace{\tilde{y}_n}_{X_1} \oplus \left(\underbrace{\tilde{h} \otimes \tilde{\lambda}}_{X_2} \otimes \underbrace{\tilde{y}_n}_y \right) - \left(\underbrace{1}_{\alpha_1} + \underbrace{h\lambda}_{\alpha_2} \right) \underbrace{\tilde{y}_n}_y \right|.$$

$$\left| \underbrace{\tilde{y}_n}_{X_1} - \underbrace{1}_{\alpha_1} \underbrace{\tilde{y}_n}_y \right| \leq \underbrace{0}_{C_1} \left| \underbrace{\tilde{y}_n}_y \right|, \quad \left| \underbrace{\tilde{h} \otimes \tilde{\lambda}}_{X_2} - \underbrace{h\lambda}_{\alpha_2} \right| \leq \underbrace{6u}_{C_2} \text{ (Gappa)}.$$

Erreurs locales de méthodes d'ordre supérieur

Méthode **RK4** stable :

FPterm	C
\tilde{y}_n	0
○ $[\tilde{h}\frac{1}{6}\tilde{\lambda}]$	$2u$ (<i>Gappa</i>)
○ $[\tilde{y}_n + \tilde{h}\frac{1}{6}\tilde{\lambda}]$	$4u$
○ $[\tilde{h}\frac{1}{3}\tilde{\lambda}]$	$4u$ (<i>Gappa</i>)
○ $[\tilde{h}\tilde{h}\frac{1}{6}\tilde{\lambda}\tilde{\lambda}]$	$12u$ (<i>Gappa</i>)
○ $[\tilde{h}\tilde{h}\tilde{h}\frac{1}{12}\tilde{\lambda}\tilde{\lambda}\tilde{\lambda}]$	$28u$ (<i>Gappa</i>)
○ $[\tilde{h}\tilde{h}\tilde{h}\tilde{h}\frac{1}{24}\tilde{\lambda}\tilde{\lambda}\tilde{\lambda}\tilde{\lambda}]$	$53u$ (<i>Gappa</i>)
...	...
○ $[\tilde{y}_n + \dots + \tilde{h}\tilde{h}\tilde{h}\tilde{h}\frac{1}{24}\tilde{\lambda}\tilde{\lambda}\tilde{\lambda}\tilde{\lambda}\tilde{y}_n]$	$194u$

Erreurs locales de méthodes classiques

Lemme 3 : Erreur locale du schéma d'Euler

Supposons $-2 \leq h\lambda \leq -2^{-100}$ et $2^{-60} \leq h \leq 1$. Alors :
 $\forall n \in \mathbb{N}, \quad \varepsilon_{n+1}^{Euler} \leq 11.01u |\tilde{y}_n|.$

Lemme 4 : Erreur locale du schéma RK2

Supposons $-2 \leq h\lambda \leq -2^{-100}$ et $2^{-60} \leq h \leq 1$. Alors :
 $\forall n \in \mathbb{N}, \quad \varepsilon_{n+1}^{RK2} \leq 28.01u |\tilde{y}_n|.$

Lemme 5 : Erreur locale du schéma RK4

Supposons $-3 \leq h\lambda \leq -2^{-100}$ et $2^{-60} \leq h \leq 1$. Alors :
 $\forall n \in \mathbb{N}, \quad \varepsilon_{n+1}^{RK4} \leq 194u |\tilde{y}_n|.$

Erreurs globales de méthodes classiques

Bornes sur les **erreurs globales** de schémas classiques :

- **Euler** : $C = 11,01u$

$$|E_n| \leq (11,01u + |R(h, \lambda)|)^n \left(\varepsilon_0 + n \frac{11,01u|y_0|}{11,01u + |R(h, \lambda)|} \right) ;$$

- **RK2** : $C = 28,01u$

$$|E_n| \leq (28,01u + |R(h, \lambda)|)^n \left(\varepsilon_0 + n \frac{28,01u|y_0|}{28,01u + |R(h, \lambda)|} \right) ;$$

- **RK4** : $C = 194u$

$$|E_n| \leq (194u + |R(h, \lambda)|)^n \left(\varepsilon_0 + n \frac{194u|y_0|}{194u + |R(h, \lambda)|} \right).$$

⇒ pas de compensation ☹️ mais bornes raisonnables 😊.

Table des matières

- 1 Introduction
- 2 Méthodes et systèmes étudiés
- 3 Erreurs d'arrondi
- 4 Conclusion et Perspectives

Méthodologie générale et bilan

$$y_{n+1} = \sum_{i=0}^M \alpha_i y_n$$

$$\widetilde{y}_{n+1} = \oplus_{i=0}^M \widetilde{\alpha}_i \otimes \widetilde{y}_n$$

$$(\alpha_i = h\lambda, \frac{h\lambda}{3}, \frac{h\lambda}{6}, \frac{h^4\lambda^4}{24} \dots)$$

Méthodologie générale et bilan

$$y_{n+1} = \sum_{i=0}^M \alpha_i y_n \quad \widetilde{y}_{n+1} = \oplus_{i=0}^M \widetilde{\alpha}_i \otimes \widetilde{y}_n \quad (\alpha_i = h\lambda, \frac{h\lambda}{3}, \frac{h\lambda}{6}, \frac{h^4\lambda^4}{24} \dots)$$

Méthodologie pour borner les erreurs d'arrondi :

Méthodologie générale et bilan

$$y_{n+1} = \sum_{i=0}^M \alpha_i y_n \quad \widetilde{y}_{n+1} = \oplus_{i=0}^M \widetilde{\alpha}_i \otimes \widetilde{y}_n \quad \left(\alpha_i = h\lambda, \frac{h\lambda}{3}, \frac{h\lambda}{6}, \frac{h^4\lambda^4}{24} \dots \right)$$

Méthodologie pour borner les erreurs d'arrondi :

- 1) calculer l'erreur $\varepsilon_0 = |\widetilde{y}_0 - y_0|$;

Méthodologie générale et bilan

$$y_{n+1} = \sum_{i=0}^M \alpha_i y_n \quad \widetilde{y}_{n+1} = \oplus_{i=0}^M \widetilde{\alpha}_i \otimes \widetilde{y}_n \quad (\alpha_i = h\lambda, \frac{h\lambda}{3}, \frac{h\lambda}{6}, \frac{h^4\lambda^4}{24} \dots)$$

Méthodologie pour borner les erreurs d'arrondi :

- 1) calculer l'erreur $\varepsilon_0 = |\widetilde{y}_0 - y_0|$;
- 2) borner les erreurs sur chaque terme α_i (Gamma + stabilité) ;

Méthodologie générale et bilan

$$y_{n+1} = \sum_{i=0}^M \alpha_i y_n \quad \widetilde{y}_{n+1} = \oplus_{i=0}^M \widetilde{\alpha}_i \otimes \widetilde{y}_n \quad (\alpha_i = h\lambda, \frac{h\lambda}{3}, \frac{h\lambda}{6}, \frac{h^4\lambda^4}{24} \dots)$$

Méthodologie pour borner les erreurs d'arrondi :

- 1) calculer l'erreur $\varepsilon_0 = |\widetilde{y}_0 - y_0|$;
- 2) borner les erreurs sur chaque terme α_i (Gamma + stabilité) ;
- 3) borner l'erreur locale par M applications du Lemme 2 ;

Méthodologie générale et bilan

$$y_{n+1} = \sum_{i=0}^M \alpha_i y_n \quad \widetilde{y}_{n+1} = \oplus_{i=0}^M \widetilde{\alpha}_i \otimes \widetilde{y}_n \quad (\alpha_i = h\lambda, \frac{h\lambda}{3}, \frac{h\lambda}{6}, \frac{h^4\lambda^4}{24} \dots)$$

Méthodologie pour borner les erreurs d'arrondi :

- 1) calculer l'erreur $\varepsilon_0 = |\widetilde{y}_0 - y_0|$;
- 2) borner les erreurs sur chaque terme α_i (Gamma + stabilité) ;
- 3) borner l'erreur locale par M applications du Lemme 2 ;
- 4) borner l'erreur globale par instantiation du Théorème 1.

Méthodologie générale et bilan

$$y_{n+1} = \sum_{i=0}^M \alpha_i y_n \quad \widetilde{y}_{n+1} = \oplus_{i=0}^M \widetilde{\alpha}_i \otimes \widetilde{y}_n \quad (\alpha_i = h\lambda, \frac{h\lambda}{3}, \frac{h\lambda}{6}, \frac{h^4\lambda^4}{24} \dots)$$

Méthodologie pour borner les erreurs d'arrondi :

- 1) calculer l'erreur $\varepsilon_0 = |\widetilde{y}_0 - y_0|$;
- 2) borner les erreurs sur chaque terme α_i (Gamma + stabilité) ;
- 3) borner l'erreur locale par M applications du Lemme 2 ;
- 4) borner l'erreur globale par instantiation du Théorème 1.

Bilan :

- analyse **fine** et **mécanique** des erreurs d'arrondi ;
- **pas de compensation** (méthodes explicites à un pas) ;
- borne d'erreur **linéaire**.

Perspectives

- **overflows** et **underflows** à prendre en compte ;
- formalisation en Coq (assistant de preuves) :
 - utilisation bibliothèque **Flocq** [Boldo, Melquiond],
 - utilisation des tactiques **Gappa** et **interval** [Melquiond],
 - travail sur l'**automatisation** des preuves,
- **ÉDOs** + générales : complexes, matricielles, non linéaires ;
- **méthodes** + générales : multi-pas, pas variable, implicites ;
- somme **erreurs de méthodes** et **erreurs d'arrondi** ;
- équations aux dérivées partielles (éléments finis).