Cours d'analyse numérique

Laurent Seppecher & Grégory Vial

École Centrale de Lyon

Année 2022-2023 - Semestre 5



► Monitorat (romain.meaux@ecl21.ec-lyon.fr)

- ► Monitorat (romain.meaux@ecl21.ec-lyon.fr)
- ► Cours : pré-requis à regarder avant.

- ► Monitorat (romain.meaux@ecl21.ec-lyon.fr)
- Cours : pré-requis à regarder avant.
- ► TD : à préparer ! Utilisation des ordinateurs portables (Matlab).

- Monitorat (romain.meaux@ecl21.ec-lyon.fr)
- Cours : pré-requis à regarder avant.
- ► TD : à préparer ! Utilisation des ordinateurs portables (Matlab).
- Un tutoriel Matlab succinct est disponible (et à apprendre par cœur)

- Monitorat (romain.meaux@ecl21.ec-lyon.fr)
- Cours : pré-requis à regarder avant.
- ► TD : à préparer ! Utilisation des ordinateurs portables (Matlab).
- Un tutoriel Matlab succinct est disponible (et à apprendre par cœur)
- Contrôle des connaissances :



- Monitorat (romain.meaux@ecl21.ec-lyon.fr)
- Cours : pré-requis à regarder avant.
- ► TD : à préparer ! Utilisation des ordinateurs portables (Matlab).
- Un tutoriel Matlab succinct est disponible (et à apprendre par cœur)
- Contrôle des connaissances :
 - Savoir-faire 25% : mini-test (1H lors du TD4) à confirmer
 - Autorisés : documents cours/TD, programmes matlab.
 - Savoir 75%: examen terminal (1H30).
 - Autorisé : un A4 recto-verso manuscrit.





L'analyse numérique, c'est quoi? ça sert à quoi?

Contexte

- Contexte
 - Pas de solutions explicites en général.

- Contexte
 - Pas de solutions explicites en général.
 - Besoin de calcul approché.

- Contexte
 - Pas de solutions explicites en général.
 - Besoin de calcul approché.
- Analyse numérique

- Contexte
 - Pas de solutions explicites en général.
 - Besoin de calcul approché.
- Analyse numérique
 - Mise au point de méthodes numériques.

- Contexte
 - Pas de solutions explicites en général.
 - Besoin de calcul approché.
- Analyse numérique
 - Mise au point de méthodes numériques.
 - Preuve de convergence. Limitations des méthodes.

- Contexte
 - Pas de solutions explicites en général.
 - Besoin de calcul approché.
- Analyse numérique
 - Mise au point de méthodes numériques.
 - Preuve de convergence. Limitations des méthodes.
 - Quantification de la vitesse de convergence.



- Contexte
 - Pas de solutions explicites en général.
 - Besoin de calcul approché.
- Analyse numérique
 - Mise au point de méthodes numériques.
 - Preuve de convergence. Limitations des méthodes.
 - Quantification de la vitesse de convergence.
 - Tests par simulations numériques

- Contexte
 - Pas de solutions explicites en général.
 - Besoin de calcul approché.
- Analyse numérique
 - Mise au point de méthodes numériques.
 - Preuve de convergence. Limitations des méthodes.
 - Quantification de la vitesse de convergence.
 - Tests par simulations numériques
- Calcul scientifique



- Contexte
 - Pas de solutions explicites en général.
 - Besoin de calcul approché.
- Analyse numérique
 - Mise au point de méthodes numériques.
 - Preuve de convergence. Limitations des méthodes.
 - Quantification de la vitesse de convergence.
 - Tests par simulations numériques
- Calcul scientifique
 - Programmation sur problèmes réels.



- Contexte
 - Pas de solutions explicites en général.
 - Besoin de calcul approché.
- Analyse numérique
 - Mise au point de méthodes numériques.
 - Preuve de convergence. Limitations des méthodes.
 - Quantification de la vitesse de convergence.
 - Tests par simulations numériques
- ► Calcul scientifique
 - Programmation sur problèmes réels.
 - Optimisation, parallélisation, etc.



Exemples de vitesses de convergence

- **Méthode numérique** = suite (u_k) .
 - Convergence

 $u_k \longrightarrow u$.

Exemples de vitesses de convergence

- **Méthode numérique** = suite (u_k) .
 - Convergence

$$u_k \longrightarrow u$$
.

Convergence polynomiale

$$\exists C > 0, \ \exists p > 0, \ \forall k, \quad |u_k - u| \leqslant \frac{C}{k^p}$$

Exemples de vitesses de convergence

- **Méthode numérique** = suite (u_k) .
 - Convergence

$$u_k \longrightarrow u$$
.

Convergence polynomiale

$$\exists C > 0, \ \exists \rho > 0, \ \forall k, \quad |u_k - u| \leqslant \frac{C}{k\rho} \quad i.e. \quad \ln|u_k - u| \leqslant \ln C - \rho \ln k.$$

Exemples de vitesses de convergence

- **Méthode numérique** = suite (u_k) .
 - Convergence

$$u_k \longrightarrow u$$
.

Convergence polynomiale

$$\exists C > 0, \ \exists p > 0, \ \forall k, \quad |u_k - u| \leqslant \frac{C}{k^p} \quad i.e. \quad \ln|u_k - u| \leqslant \ln C - p \ln k.$$

Convergence géométrique

$$\exists C > 0, \ \exists \rho \in [0,1[, \ \forall k, \ |u_k - u| \leqslant C\rho^k]$$

Exemples de vitesses de convergence

- **Méthode numérique** = suite (u_k) .
 - Convergence

$$u_k \longrightarrow u$$
.

Convergence polynomiale

$$\exists C > 0, \ \exists p > 0, \ \forall k, \quad |u_k - u| \leqslant \frac{C}{k^p} \quad i.e. \quad \ln|u_k - u| \leqslant \ln C - p \ln k.$$

Convergence géométrique

$$\exists C > 0, \ \exists \rho \in [0,1[,\ \forall k,\ |u_k - u| \leqslant C\rho^k \quad i.e. \quad \ln|u_k - u| \leqslant \ln C + k \ln \rho.$$

Exemples de vitesses de convergence

- **Méthode numérique** = suite (u_k) .
 - Convergence

$$u_k \longrightarrow u$$
.

Convergence polynomiale

$$\exists C > 0, \ \exists p > 0, \ \forall k, \quad |u_k - u| \leqslant \frac{C}{kp} \quad i.e. \quad \ln|u_k - u| \leqslant \ln C - p \ln k.$$

Convergence géométrique

$$\exists C > 0, \ \exists \rho \in [0,1[\,,\ \forall k, \quad |u_k - u| \leqslant C\rho^k \quad i.e. \quad \ln|u_k - u| \leqslant \ln C + k \ln \rho.$$

Convergence quadratique

$$\exists C > 0, \ \exists \rho \in [0,1[\ , \ \forall k, \quad |u_k - u| \leqslant C \rho^{2^k}$$

Exemples de vitesses de convergence

- **Méthode numérique** = suite (u_k) .
 - Convergence

$$u_k \longrightarrow u$$
.

Convergence polynomiale

$$\exists C > 0, \ \exists p > 0, \ \forall k, \quad |u_k - u| \leqslant \frac{C}{kp} \quad i.e. \quad \ln|u_k - u| \leqslant \ln C - p \ln k.$$

Convergence géométrique

$$\exists C > 0, \ \exists \rho \in [0,1[, \ \forall k, \quad |u_k - u| \leqslant C\rho^k \quad i.e. \quad \ln|u_k - u| \leqslant \ln C + k \ln \rho.$$

Convergence quadratique

$$\exists C > 0, \ \exists \rho \in [0,1[,\ \forall k,\quad |u_k - u| \leqslant C\rho^{2^k} \quad i.e. \quad \ln|u_k - u| \leqslant \ln C + 2^k \ln \rho.$$

Plan du cours

- Cours 1. Systèmes linéaires et calcul de valeurs propres
- Cours 2. Interpolation et intégration numérique
- Cours 3. Optimisation numérique
- Cours 4. Approximation numérique des équations différentielles
- Cours 5. Discrétisation des EDP : l'équation de Laplace
- ▶ Cours 6. Discrétisation des EDP : l'équation de transport



Un conseil



Un conseil



« On ne peut pas faire de physique sans un bon niveau en mathématiques. Mais c'est vrai quelle que soit la science. Je vais même être caricatural : si vous voulez trouver du travail facilement, faites des maths. »



COURS 1

Systèmes linéaires Calcul de valeurs propres



Le modèle de Léontiev en comptabilité nationale (version production)



- ightharpoonup d branches : B_i produit P_i .
- ▶ 1 bien P_j nécessite a_{ij} biens P_i ,
- $ightharpoonup B_i$ produit une quantité q_i .



Le modèle de Léontiev en comptabilité nationale (version production)



Bilan pour la branche B_i :

- ightharpoonup d branches : B_i produit P_i .
- ▶ 1 bien P_j nécessite a_{ij} biens P_i ,
- $ightharpoonup B_i$ produit une quantité q_i .

$q_i - \sum_{i=1}^d a_{ij}q_j = (Mq)_i$ avec M = I - A.



Le modèle de Léontiev en comptabilité nationale (version production)



- ightharpoonup d branches : B_i produit P_i .
- ▶ 1 bien P_j nécessite a_{ij} biens P_i ,
- ► B_i produit une quantité q_i.

Bilan pour la branche B_i :

$$q_i - \sum_{i=1}^d a_{ij}q_j = (Mq)_i$$
 avec $M = I - A$.

Équilibre offre-demande :

$$Mq = \delta$$
.



Le modèle de Léontiev en comptabilité nationale (version production)



- ightharpoonup d branches : B_i produit P_i .
- ▶ 1 bien P_j nécessite a_{ij} biens P_i ,
- ► B_i produit une quantité q_i.

Bilan pour la branche B_i :

$$q_i - \sum_{i=1}^d a_{ij}q_j = (Mq)_i$$
 avec $M = I - A$.

Équilibre offre-demande :

$$Mq = \delta$$
.

Question : quelle influence de la demande sur les quantités produites ?

$$\Longrightarrow$$
 système linéaire $Mq=\delta$.



Le modèle de Léontiev en comptabilité nationale (version bénéfice)



- ightharpoonup d branches : B_i produit P_i .
- ▶ 1 bien P_i nécessite
 - $ightharpoonup a_{ji}$ biens P_j ,
 - $ightharpoonup \ell_i$ travailleurs au salaire w.
- $ightharpoonup B_i$ vend au prix π_i .
- salaire fonction linéaire des prix : $w = \sum c_j \pi_j$.



Le modèle de Léontiev en comptabilité nationale (version bénéfice)



- \triangleright d branches : B_i produit P_i .
- ▶ 1 bien P_i nécessite
 - $ightharpoonup a_{ji}$ biens P_j ,
 - $ightharpoonup \ell_i$ travailleurs au salaire w.
- \triangleright B_i vend au prix π_i .
- salaire fonction linéaire des prix : $w = \sum c_j \pi_j$.

Coût de production de la branche B_i :

$$\sum_{i=1}^d a_{ji}\pi_j + \ell_i w = \sum_{i=1}^d (a_{ji} + \ell_i c_j)\pi_j = (N\pi)_i \quad \text{avec} \quad N = A^T + \ell c^T.$$



Le modèle de Léontiev en comptabilité nationale (version bénéfice)



- \triangleright d branches : B_i produit P_i .
- ▶ 1 bien *P_i* nécessite
 - $ightharpoonup a_{ji}$ biens P_j ,
 - $ightharpoonup \ell_i$ travailleurs au salaire w.
- \triangleright B_i vend au prix π_i .
- salaire fonction linéaire des prix : $w = \sum c_j \pi_j$.

Coût de production de la branche B_i :

$$\sum_{j=1}^{d} a_{ji} \pi_{j} + \ell_{i} w = \sum_{j=1}^{d} (a_{ji} + \ell_{i} c_{j}) \pi_{j} = (N\pi)_{i} \quad \text{avec} \quad N = A^{T} + \ell c^{T}.$$

Profit relatif de la branche B_i : $\tau_i = \frac{\pi_i - (N\pi)_i}{(N\pi)_i}$.



Le modèle de Léontiev en comptabilité nationale (version bénéfice)



- \triangleright d branches : B_i produit P_i .
- ▶ 1 bien *P_i* nécessite
 - $ightharpoonup a_{ji}$ biens P_j ,
 - $ightharpoonup \ell_i$ travailleurs au salaire w.
- $ightharpoonup B_i$ vend au prix π_i .
- salaire fonction linéaire des prix : $w = \sum c_i \pi_i$.

Coût de production de la branche B_i :

$$\sum_{j=1}^d a_{ji}\pi_j + \ell_i w = \sum_{j=1}^d (a_{ji} + \ell_i c_j)\pi_j = (N\pi)_i \quad \text{avec} \quad N = A^T + \ell \, c^T.$$

Profit relatif de la branche B_i : $\tau_i = \frac{\pi_i - (N\pi)_i}{(N\pi)_i}$.

Question: profit équilibré parmi les branches?



Le modèle de Léontiev en comptabilité nationale (version bénéfice)



- \triangleright d branches : B_i produit P_i .
- ▶ 1 bien P_i nécessite
 - $ightharpoonup a_{ji}$ biens P_j ,
 - $ightharpoonup \ell_i$ travailleurs au salaire w.
- $ightharpoonup B_i$ vend au prix π_i .
- salaire fonction linéaire des prix : $w = \sum c_i \pi_i$.

Coût de production de la branche B_i :

$$\sum_{j=1}^d a_{ji}\pi_j + \ell_i w = \sum_{j=1}^d (a_{ji} + \ell_i c_j)\pi_j = (N\pi)_i \quad \text{avec} \quad N = A^T + \ell \, c^T.$$

Profit relatif de la branche
$$B_i$$
: $\tau_i = \frac{\pi_i - (N\pi)_i}{(N\pi)_i}$.

Question : profit équilibré parmi les branches?



CENTRALEL

Analyse numérique matricielle Résolution de systèmes linéaires



Analyse numérique matricielle Résolution de systèmes linéaires

▶ En théorie... Soit A est inversible. $\exists ! x$ t.q. Ax = b. On a $x = A^{-1}b$.



Résolution de systèmes linéaires

- ► En théorie... Soit A est inversible. $\exists ! x$ t.q. Ax = b. On a $x = A^{-1}b$.
- En pratique...
 - Soit $H \in \mathbb{R}^{d \times d}$ t.q. $H_{ij} = \frac{1}{i+i-1}$. [hilbert.m]



Rq. H est inversible et $H^{-1} \in \mathbb{Z}^{d \times d}$ (formule explicite).



Résolution de systèmes linéaires

- ► En théorie... Soit A est inversible. $\exists !x$ t.q. Ax = b. On a $x = A^{-1}b$.
- En pratique...
 - Soit $H \in \mathbb{R}^{d \times d}$ t.q. $H_{ij} = \frac{1}{i+j-1}$. [hilbert.m]



Rq. H est inversible et $H^{-1} \in \mathbb{Z}^{d \times d}$ (formule explicite).

⇒ comment expliquer le phénomène?



Résolution de systèmes linéaires

- ► En théorie... Soit A est inversible. $\exists ! x$ t.q. Ax = b. On a $x = A^{-1}b$.
- En pratique...
 - Soit $H \in \mathbb{R}^{d \times d}$ t.q. $H_{ij} = \frac{1}{i+i-1}$. [hilbert.m]



Rg. H est inversible et $H^{-1} \in \mathbb{Z}^{d \times d}$ (formule explicite).

- ⇒ comment expliquer le phénomène?
- Formules de Cramer :

$$x_i = \frac{\det(A_i)}{\det(A)}, \quad \text{avec } A_i = (C_1|\ldots|C_{i-1}|b|C_{i+1}|\ldots|C_d).$$



Résolution de systèmes linéaires

- ► En théorie... Soit A est inversible. $\exists ! x$ t.q. Ax = b. On a $x = A^{-1}b$.
- En pratique...
 - Soit $H \in \mathbb{R}^{d \times d}$ t.q. $H_{ij} = \frac{1}{i+i-1}$. [hilbert.m]



Rg. H est inversible et $H^{-1} \in \mathbb{Z}^{d \times d}$ (formule explicite).

- ⇒ comment expliquer le phénomène?
- Formules de Cramer :

$$x_i = \frac{\det(A_i)}{\det(A)}, \quad \text{avec } A_i = (C_1|\ldots|C_{i-1}|b|C_{i+1}|\ldots|C_d).$$

Coût si les det sont calculés par dvpt p/r lignes : $\simeq (d+1)!$.



Résolution de systèmes linéaires

- ► En théorie... Soit A est inversible. $\exists ! x$ t.q. Ax = b. On a $x = A^{-1}b$.
- En pratique...
 - Soit $H \in \mathbb{R}^{d \times d}$ t.q. $H_{ij} = \frac{1}{i+i-1}$. [hilbert.m]



Rg. H est inversible et $H^{-1} \in \mathbb{Z}^{d \times d}$ (formule explicite).

- ⇒ comment expliquer le phénomène?
- Formules de Cramer :

$$x_i = \frac{\det(A_i)}{\det(A)}, \quad \text{avec } A_i = (C_1|\ldots|C_{i-1}|b|C_{i+1}|\ldots|C_d).$$

Coût si les det sont calculés par dvpt p/r lignes : $\simeq (d+1)!$.

⇒ quels algorithmes? quelles performances?



Analyse numérique matricielle Conditionnement des systèmes linéaires

Question : comment quantifier la propagation des erreurs?



Analyse numérique matricielle Conditionnement des systèmes linéaires

Question: comment quantifier la propagation des erreurs?

Proposition. Soient $A \in \mathbb{R}^{d \times d}$ inversible, $b, \delta b \in \mathbb{R}^d$. On note $x \in \mathbb{R}^d$ et $x + \delta x \in \mathbb{R}^d$ les solutions de $Ax = b, \qquad A(x + \delta x) = b + \delta b.$ Alors $\frac{\|\delta x\|}{\|x\|} \leqslant \operatorname{cond}(A) \frac{\|\delta b\|}{\|b\|},$ avec $\operatorname{cond}(A) = \|\|A\| \times \||A^{-1}\||$.

$$Ax = b,$$
 $A(x + \delta x) = b + \delta b.$

$$\frac{\|\delta x\|}{\|x\|} \leqslant \operatorname{cond}(A) \frac{\|\delta b\|}{\|b\|},$$



Conditionnement des systèmes linéaires

Question: comment quantifier la propagation des erreurs?

Proposition. Soient $A \in \mathbb{R}^{d \times d}$ inversible, $b, \delta b \in \mathbb{R}^d$. On note $x \in \mathbb{R}^d$ et $x + \delta x \in \mathbb{R}^d$ les solutions de $Ax = b, \qquad A(x + \delta x) = b + \delta b.$ Alors $\frac{\|\delta x\|}{\|x\|} \leqslant \operatorname{cond}(A) \frac{\|\delta b\|}{\|b\|},$ avec $\operatorname{cond}(A) = \|\|A\| \times \||A^{-1}\||$.

$$Ax = b,$$
 $A(x + \delta x) = b + \delta b.$

$$\frac{\|\delta x\|}{\|x\|} \leqslant \operatorname{cond}(A) \frac{\|\delta b\|}{\|b\|},$$

avec cond(A) =
$$|||A||| \times |||A^{-1}|||$$
.

- $\|A\|$ est la norme subordonnée à $\|x\|$.
- cond(A) mesure la propagation des erreurs relatives.
- ightharpoonup Si $\|\cdot\| = \|\cdot\|_p$, on note cond_p.



Conditionnement des systèmes linéaires

Question: comment quantifier la propagation des erreurs?

Proposition. Soient $A \in \mathbb{R}^{d \times d}$ inversible, $b, \delta b \in \mathbb{R}^d$. On note $x \in \mathbb{R}^d$ et $x + \delta x \in \mathbb{R}^d$ les solutions de $Ax = b, \qquad A(x + \delta x) = b + \delta b.$ Alors $\frac{\|\delta x\|}{\|x\|} \leqslant \operatorname{cond}(A) \frac{\|\delta b\|}{\|b\|},$ avec $\operatorname{cond}(A) = \|\|A\| \times \|A^{-1}\|$.

$$Ax = b,$$
 $A(x + \delta x) = b + \delta b.$

$$\frac{\|\delta x\|}{\|x\|} \leqslant \operatorname{cond}(A) \frac{\|\delta b\|}{\|b\|},$$

Preuve. Par différence $A(\delta x) = \delta b$ d'où $\|\delta x\| \leq \|A^{-1}\| \times \|\delta b\|$.

Par ailleurs b = Ax d'où $||b|| \le |||A||| \times ||x||$, soit $\frac{1}{||x||} \le \frac{|||A|||}{||b||}$.

D'où le résultat par produit.



Analyse numérique matricielle Conditionnement des systèmes linéaires

$$\begin{aligned} \textbf{Proposition. Soit } A &\in \mathbb{R}^{d \times d} \text{ inversible.} \\ &\blacktriangleright \forall \|\cdot\|, \quad \operatorname{cond}(A) \geqslant 1. \\ &\blacktriangleright \text{ Si } A \text{ est symétrique,} \\ &\quad \operatorname{cond}_2(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}. \\ &\blacktriangleright \text{ Si } A \text{ est orthogonale, } \operatorname{cond}_2(A) = 1. \end{aligned}$$



Conditionnement des systèmes linéaires

Proposition. Soit
$$A \in \mathbb{R}^{d \times d}$$
 inversible.
$$\forall \|\cdot\|, \quad \mathsf{cond}(A) \geqslant 1.$$

$$\Rightarrow \mathsf{Si} \ A \ \mathsf{est} \ \mathsf{sym\acute{e}trique},$$

$$\mathsf{cond}_2(A) = \frac{\lambda_{\mathsf{max}}(A)}{\lambda_{\mathsf{min}}(A)}.$$

$$\Rightarrow \mathsf{Si} \ A \ \mathsf{est} \ \mathsf{orthogonale}, \ \mathsf{cond}_2(A) = 1.$$

Rmq. Si cond(A) est « grand », on dit que A est mal conditionnée.



Conditionnement des systèmes linéaires

$$\begin{array}{l} \textbf{Proposition. Soit } A \in \mathbb{R}^{d \times d} \text{ inversible.} \\ & \blacktriangleright \forall \| \cdot \|, \quad \operatorname{cond}(A) \geqslant 1. \\ & \blacktriangleright \text{ Si } A \text{ est symétrique,} \\ & \quad \operatorname{cond}_2(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}. \\ & \blacktriangleright \text{ Si } A \text{ est orthogonale, } \operatorname{cond}_2(A) = 1. \end{array}$$

Rmq. Si cond(A) est « grand », on dit que A est mal conditionnée.

d	3	5	10	20
$cond_2(H)$	524	4.8×10^5	1.6×10^{13}	2.5×10^{18}

Table - Conditionnement de la matrice de Hilbert.



L'algorithme du pivot de Gauss

- Principe.
 - éliminer les inconnues successivement.
 - on aboutit à un système triangulaire.



L'algorithme du pivot de Gauss

- Principe.
 - éliminer les inconnues successivement.
 - on aboutit à un système triangulaire.
- Description.



L'algorithme du pivot de Gauss

- Principe.
 - éliminer les inconnues successivement.
 - on aboutit à un système triangulaire.
- ▶ Description.

Initialisation



L'algorithme du pivot de Gauss

- Principe.
 - éliminer les inconnues successivement.
 - on aboutit à un système triangulaire.
- Description.

Initialisation

$$\begin{vmatrix} a_{1,1}^{[0]} & a_{1,2}^{[0]} & \dots & a_{1,d}^{[0]} \\ a_{2,1}^{[0]} & a_{2,2}^{[0]} & \dots & a_{2,d}^{[0]} & b_2^{[0]} \\ \vdots & & & & \vdots \\ a_{d,1}^{[0]} & a_{d,2}^{[0]} & \dots & a_{d,d}^{[0]} & b_d^{[0]} \end{vmatrix}$$



L'algorithme du pivot de Gauss

- Principe.
 - éliminer les inconnues successivement.
 - on aboutit à un système triangulaire.
- Description.

Étape 1.1 : si
$$a_{1,1}^{[0]} = 0$$
, on permute $L_1 \leftrightarrow L_k$ t.q. $a_{k,1}^{[0]} \neq 0$.

$$\begin{vmatrix} a_{1,1}^{[0]} & a_{1,2}^{[0]} & \dots & a_{1,d}^{[0]} \\ a_{2,1}^{[0]} & a_{2,2}^{[0]} & \dots & a_{2,d}^{[0]} & b_2^{[0]} \\ \vdots & & & & \vdots \\ a_{d,1}^{[0]} & a_{d,2}^{[0]} & \dots & a_{d,d}^{[0]} & b_d^{[0]} \end{vmatrix}$$



L'algorithme du pivot de Gauss

- Principe.
 - éliminer les inconnues successivement.
 - on aboutit à un système triangulaire.
- ▶ Description.

Étape 1.2 :
$$\tilde{a}_{1,1}^{[0]} \neq 0$$
; on effectue $L_k \leftarrow L_k - r_k L_1$ avec $r_k = \frac{\tilde{s}_{k,1}^{[0]}}{\tilde{s}_{1,1}^{[0]}}$.

$$\begin{vmatrix} \tilde{a}_{1,1}^{[0]} & \tilde{a}_{1,2}^{[0]} & \dots & \tilde{a}_{1,d}^{[0]} & & \tilde{b}_{1}^{[0]} \\ \tilde{a}_{2,1}^{[0]} & \tilde{a}_{2,2}^{[0]} & \dots & \tilde{a}_{2,d}^{[0]} & & \tilde{b}_{2}^{[0]} \\ \vdots & & & & \vdots \\ \tilde{a}_{d,1}^{[0]} & \tilde{a}_{d,2}^{[0]} & \dots & \tilde{a}_{d,d}^{[0]} & & \tilde{b}_{d}^{[0]} \end{vmatrix}$$



L'algorithme du pivot de Gauss

- Principe.
 - éliminer les inconnues successivement.
 - on aboutit à un système triangulaire.
- ▶ Description.

Étape 1.2 :
$$\tilde{a}_{1,1}^{[0]} \neq 0$$
; on effectue $L_k \leftarrow L_k - r_k L_1$ avec $r_k = \frac{\tilde{s}_{k,1}^{[0]}}{\tilde{s}_{1,1}^{[0]}}$.

$$\begin{vmatrix} \tilde{a}_{1,1}^{[0]} & \tilde{a}_{1,2}^{[0]} & \dots & \tilde{a}_{1,d}^{[0]} & & \tilde{b}_{1}^{[0]} \\ 0 & a_{2,2}^{[1]} & \dots & a_{2,d}^{[1]} & & b_{2}^{[1]} \\ \vdots & & & & \vdots \\ 0 & a_{d,2}^{[1]} & \dots & a_{d,d}^{[1]} & & b_{d}^{[1]} \end{vmatrix}$$



L'algorithme du pivot de Gauss

- Principe.
 - éliminer les inconnues successivement.
 - on aboutit à un système triangulaire.
- Description.

Étape 2 : On recommence avec la sous-matrice!

$\tilde{a}_{1,1}^{[0]}$	$ ilde{a}_{1,2}^{[0]} \dots$	$\tilde{a}_{1,d}^{[0]}$	$ ilde{b}_1^{[0]}$
0	$a_{2,2}^{[1]}$	$a_{2,d}^{[1]}$	$b_2^{[1]}$
:	i i		
0	$a_{d,2}^{[1]}$	$a_{d,d}^{[1]}$	$\mid \mid b_d^{[1]} \mid$



L'algorithme du pivot de Gauss

- Principe.
 - éliminer les inconnues successivement.
 - on aboutit à un système triangulaire.
- Synthèse.
 - Après d-1 étapes, on est ramené à $Uy = \bar{c}$ avec U tri. sup.
 - $Vy = \bar{c}$ est résolu par « remontée ».



L'algorithme du pivot de Gauss

- Principe.
 - éliminer les inconnues successivement.
 - on aboutit à un système triangulaire.
- Synthèse.
 - Après d-1 étapes, on est ramené à $Uy=\bar{c}$ avec U tri. sup.
 - $Vy = \bar{c}$ est résolu par « remontée ».
- Coût de calcul.
 - L'étape k affecte $(d k)^2$ coefficients de la matrice.
 - Coût total $\simeq \sum_{k=1}^{d-1} (d-k)^2 = \mathcal{O}(d^3)$.
 - Les opérations sur le second membre et la remontée sont négligeables.



La décomposition LU

Interprétation matricielle des opérations élémentaires

▶ Permutation
$$L_i \leftrightarrow L_k \iff A \mapsto PA$$

► Transvection
$$L_k \leftarrow L_k - rL_i \iff A \mapsto TA$$
 (pour $k > i$)



La décomposition LU

- Interprétation matricielle des opérations élémentaires
 - $\blacktriangleright \ \, \mathsf{Permutation} \ \, \boxed{L_i \leftrightarrow L_k} \Longleftrightarrow \boxed{A \mapsto PA}$
 - ▶ Transvection $L_k \leftarrow L_k rL_i \iff A \mapsto TA$ (pour k > i)

► Écriture matricielle de l'algorithme de Gauss

 $A \longmapsto M_1 \times M_2 \times \cdots \times M_q A = U$ avec les M_p de type P ou T.

La décomposition LU

- Interprétation matricielle des opérations élémentaires
 - ▶ Permutation $L_i \leftrightarrow L_k$ \iff $A \mapsto PA$
 - ► Transvection $L_k \leftarrow L_k rL_i \iff A \mapsto TA$ (pour k > i)

► Écriture matricielle de l'algorithme de Gauss

$$A \longmapsto M_1 \times M_2 \times \cdots \times M_q A = U$$
 avec les M_p de type P ou T .



La décomposition LU

- Interprétation matricielle des opérations élémentaires
 - ▶ Permutation $L_i \leftrightarrow L_k \iff A \mapsto PA$
 - ► Transvection $L_k \leftarrow L_k rL_i \iff A \mapsto TA$ (pour k > i)

► Écriture matricielle de l'algorithme de Gauss

$$A = N_1 \times N_2 \times \cdots \times N_q \times U$$
 avec les N_p de type P ou $T!$



La décomposition LU

- Interprétation matricielle des opérations élémentaires
 - ▶ Permutation $L_i \leftrightarrow L_k \iff A \mapsto PA$
 - ▶ Transvection $L_k \leftarrow L_k rL_i \iff A \mapsto TA$ (pour k > i)
- ► Écriture matricielle de l'algorithme de Gauss

$$A = N_1 \times N_2 \times \cdots \times N_q \times U$$
 avec les N_p de type P ou T !

▶ Si aucune permutation, les N_p sont tri. inf.



La décomposition LU

- Interprétation matricielle des opérations élémentaires
 - ▶ Permutation $L_i \leftrightarrow L_k \iff A \mapsto PA$
 - ► Transvection $L_k \leftarrow L_k rL_i \iff A \mapsto TA$ (pour k > i)
- Écriture matricielle de l'algorithme de Gauss

$$\boxed{A = N_1 \times N_2 \times \cdots \times N_q \times U} \quad \text{avec les } N_p \text{ de type } P \text{ ou } T \text{!}$$

- ▶ Si aucune permutation, les N_p sont tri. inf.
- Dans ce cas,

A = LU avec L tri. inf. et U tri. sup.



La décomposition LU

- Interprétation matricielle des opérations élémentaires
 - ▶ Permutation $L_i \leftrightarrow L_k \iff A \mapsto PA$
 - ▶ Transvection $L_k \leftarrow L_k rL_i \iff A \mapsto TA$ (pour k > i)
- Écriture matricielle de l'algorithme de Gauss

$$\boxed{A = N_1 \times N_2 \times \cdots \times N_q \times U} \quad \text{avec les } N_p \text{ de type } P \text{ ou } T \text{!}$$

- \triangleright Si aucune permutation, les N_p sont tri. inf.
- Dans ce cas,

$$A = LU$$
 avec L tri. inf. et U tri. sup.

► Comment caractériser ce cas?



La décomposition LU

Définition. Soit $A \in \mathbb{R}^{d \times d}$. On appelle mineur fondamental de taille $k \leqslant d$ le déterminant $\Delta_k = \det \left[(A_{ij})_{1 \leqslant i,j \leqslant k} \right].$

$$\Delta_k = \det \left[(A_{ij})_{1 \leqslant i,j \leqslant k} \right]$$

Théorème. Soit $A \in \mathbb{R}^{d \times d}$ telle que $\forall k \leqslant d, \quad \Delta_k \neq 0.$ Alors $\exists ! (L, U)$ t.q. $\blacktriangleright \ \, L \ \, \text{tri. inf. avec diagonale unité.}$ $\blacktriangleright \ \, U \ \, \text{tri. sup.}$ $\blacktriangleright \ \, A = LU.$

$$\forall k \leqslant d, \quad \Delta_k \neq 0.$$



Analyse numérique matricielle La décomposition LU

Intérêt algorithmique de la décomposition LU.



La décomposition LU

Intérêt algorithmique de la décomposition LU.

- ▶ On doit résoudre $Ax_p = b_p$ pour de nombreux p = 1, ..., P.
 - ▶ 1 fois : A = LU.
 - P fois : $Ly_p = b_p$ et $Ux_p = y_p$.
 - Coût total: $\mathcal{O}(d^3 + Pd^2)$.

La décomposition LU

Intérêt algorithmique de la décomposition LU.

- ▶ On doit résoudre $Ax_p = b_p$ pour de nombreux p = 1, ..., P.
 - ▶ 1 fois : A = LU.
 - ightharpoonup P fois : $Ly_p = b_p$ et $Ux_p = y_p$.
 - ightharpoonup Coût total: $\mathcal{O}(d^3 + Pd^2)$.
- ► Mieux que *P* fois Gauss.



La décomposition LU

Intérêt algorithmique de la décomposition LU.

- ▶ On doit résoudre $Ax_p = b_p$ pour de nombreux p = 1, ..., P.
 - ▶ 1 fois : A = LU.
 - ightharpoonup P fois : $Ly_p = b_p$ et $Ux_p = y_p$.
 - ightharpoonup Coût total: $\mathcal{O}(d^3 + Pd^2)$.
- ► Mieux que *P* fois Gauss.
- Vraiment utile seulement lorsque les b_p sont connus seulement successivement...



La décomposition LU

Intérêt algorithmique de la décomposition LU.

- ▶ On doit résoudre $Ax_p = b_p$ pour de nombreux p = 1, ..., P.
 - ightharpoonup 1 fois : A = LU.
 - ightharpoonup P fois : $Ly_p = b_p$ et $Ux_p = y_p$.
 - ightharpoonup Coût total: $\mathcal{O}(d^3 + Pd^2)$.
- ► Mieux que *P* fois Gauss.
- Vraiment utile seulement lorsque les b_p sont connus seulement successivement...

Extensions, optimisations

- ► Si A est symétrique, définie positive, $A = BB^{\mathsf{T}}$ avec B tri.inf.
- ▶ Si A a une structure particulière, les calculs peuvent être simplifiés.



La décomposition LU

Intérêt algorithmique de la décomposition LU.

- On doit résoudre $Ax_p = b_p$ pour de nombreux p = 1, ..., P.
 - ▶ 1 fois : A = LU.
 - ightharpoonup P fois : $Ly_p = b_p$ et $Ux_p = y_p$.
 - ightharpoonup Coût total: $\mathcal{O}(d^3 + Pd^2)$.
- ► Mieux que *P* fois Gauss.
- Vraiment utile seulement lorsque les b_p sont connus seulement successivement...

Extensions, optimisations

- ► Si A est symétrique, définie positive, $A = BB^{\mathsf{T}}$ avec B tri.inf.
- ▶ Si A a une structure particulière, les calculs peuvent être simplifiés.



Méthode utilisée par matlab? Voir help \



La décomposition LU

► Un système linéaire...

```
\begin{cases} \alpha x_1 + x_2 + \dots + x_d = 1 \\ x_1 + \alpha x_2 = 1 \\ x_1 + \alpha x_3 = 1 \\ \vdots & \dots & \vdots \\ x_1 + \alpha x_d = 1 \end{cases}
```



La décomposition LU

Un système linéaire...

$$\begin{cases} \alpha x_1 + x_2 + \dots + x_d = 1 \\ x_1 + \alpha x_2 = 1 \\ x_1 + \alpha x_3 = 1 \\ \vdots & \ddots & \vdots \\ x_1 + \alpha x_d = 1 \end{cases}$$

$$\begin{cases} \alpha x_1 + x_2 + \dots + x_d = 1 \\ x_1 + \alpha x_2 & = 1 \\ x_1 + \alpha x_3 & = 1 \\ \vdots & \dots & \vdots \\ x_1 & + \alpha x_d = 1 \end{cases} \begin{cases} \alpha y_d + y_{d-1} + \dots + y_1 = 1 \\ y_d + \alpha y_{d-1} & = 1 \\ y_d & + \alpha y_{d-2} & = 1 \\ \vdots & \dots & \vdots \\ y_d & + \alpha y_1 = 1 \end{cases}$$

La décomposition LU

Un système linéaire...

$$\begin{cases} \alpha x_1 + x_2 + \dots + x_d = 1 \\ x_1 + \alpha x_2 = 1 \\ x_1 + \alpha x_3 = 1 \\ \vdots & \dots & \vdots \\ x_1 + \alpha x_d = 1 \end{cases}$$

$$\begin{cases} \alpha x_1 + x_2 + \dots + x_d = 1 \\ x_1 + \alpha x_2 = 1 \\ \vdots + \alpha x_3 = 1 \\ \vdots \\ x_1 + \alpha x_d = 1 \end{cases} \begin{cases} y_1 + y_2 + \dots + \alpha y_d = 1 \\ \alpha y_{d-1} + y_d = 1 \\ \alpha y_{d-2} + y_d = 1 \\ \vdots \\ \alpha y_1 + y_d = 1 \end{cases}$$

La décomposition LU

Un système linéaire...

$$\begin{cases} \alpha x_1 + x_2 + \dots + x_d = 1 \\ x_1 + \alpha x_2 = 1 \\ x_1 + \alpha x_3 = 1 \\ \vdots & \dots & \vdots \\ x_1 + \alpha x_d = 1 \end{cases}$$

$$\begin{cases} \alpha x_1 + x_2 + \dots + x_d = 1 \\ x_1 + \alpha x_2 & = 1 \\ x_1 + \alpha x_3 & = 1 \\ \vdots & \dots & \vdots \\ x_1 & + \alpha x_d = 1 \end{cases} \begin{cases} \alpha y_1 + y_d = 1 \\ \vdots & \dots & \vdots \\ \alpha y_{d-2} + y_d = 1 \\ \alpha y_{d-1} + y_d = 1 \\ y_1 + y_2 + \dots + \alpha y_d = 1 \end{cases}$$

La décomposition LU

Un système linéaire...

$$\begin{cases} \alpha x_1 + x_2 + \dots + x_d = 1 \\ x_1 + \alpha x_2 = 1 \\ x_1 + \alpha x_3 = 1 \\ \vdots & \dots & \vdots \\ x_1 + \alpha x_d = 1 \end{cases}$$

$$\begin{cases} \alpha x_1 + x_2 + \dots + x_d = 1 \\ x_1 + \alpha x_2 & = 1 \\ x_1 + \alpha x_3 & = 1 \\ \vdots & \dots & \vdots \\ x_1 & + \alpha x_d = 1 \end{cases} \begin{cases} \alpha y_1 & + y_d = 1 \\ \vdots & \dots & \vdots \\ \alpha y_{d-2} + y_d = 1 \\ \vdots & \alpha y_{d-1} + y_d = 1 \\ y_1 + y_2 + \dots + \alpha y_d = 1 \end{cases}$$

... deux matrices...

$$A_1 = \left(\begin{array}{cc} \alpha & 1 - 1 \\ 1 \\ 1 \\ 1 \end{array}\right)$$

$$A_{1} = \begin{pmatrix} \alpha & 1 - 1 \\ 1 & & \\ 1 & \alpha \end{pmatrix} \qquad A_{2} = \begin{pmatrix} \alpha & 1 \\ & | \\ 1 - 1 & \alpha \end{pmatrix}$$

La décomposition LU

▶ Un système linéaire... $(\alpha > d)$

$$\begin{cases} \alpha x_1 + x_2 + \dots + x_d = 1 \\ x_1 + \alpha x_2 = 1 \\ x_1 + \alpha x_3 = 1 \\ \vdots & \dots & \vdots \\ x_1 + \alpha x_d = 1 \end{cases}$$

$$\begin{cases} \alpha x_1 + x_2 + \dots + x_d = 1 \\ x_1 + \alpha x_2 & = 1 \\ x_1 + \alpha x_3 & = 1 \\ \vdots & \dots & \vdots \\ x_1 + \alpha x_d = 1 \end{cases} \begin{cases} \alpha y_1 + y_d = 1 \\ \vdots & \dots & \vdots \\ \alpha y_{d-2} + y_d = 1 \\ \alpha y_{d-1} + y_d = 1 \\ y_1 + y_2 + \dots + \alpha y_d = 1 \end{cases}$$

... deux matrices...

$$A_{1} = \left(\begin{array}{cc} \alpha & 1 - 1 \\ 1 \\ 1 \\ 1 \end{array}\right)$$

$$A_{1} = \begin{pmatrix} \alpha & 1 - 1 \\ 1 & & \\ 1 & \alpha \end{pmatrix} \qquad A_{2} = \begin{pmatrix} \alpha & 1 \\ & | \\ 1 & 1 \\ 1 - 1 & \alpha \end{pmatrix}$$

LUfleches.m



D'autres méthodes de résolution de systèmes linéaires



D'autres méthodes de résolution de systèmes linéaires

But: construire une suite $(x^{(n)})$ qui converge vers x t.q. Ax = b.



D'autres méthodes de résolution de systèmes linéaires

- ▶ But : construire une suite $(x^{(n)})$ qui converge vers x t.q. Ax = b.
- ldée : si A = M N avec M inversible, alors

$$Ax = b \iff x = \underbrace{M^{-1}(Nx + b)}_{\varphi(x)}.$$



D'autres méthodes de résolution de systèmes linéaires

- ▶ But : construire une suite $(x^{(n)})$ qui converge vers x t.q. Ax = b.
- ldée : si A = M N avec M inversible, alors

$$Ax = b \iff x = \underbrace{M^{-1}(Nx + b)}_{\varphi(x)}.$$

Problème de point fixe : soit $x^{(0)} \in \mathbb{R}^d$, on définit

$$x^{(n+1)} = \varphi(x^{(n)}).$$



D'autres méthodes de résolution de systèmes linéaires

- ▶ But : construire une suite $(x^{(n)})$ qui converge vers x t.q. Ax = b.
- ▶ Idée : si A = M N avec M inversible, alors

$$Ax = b \iff x = \underbrace{M^{-1}(Nx + b)}_{\varphi(x)}.$$

Problème de point fixe : soit $x^{(0)} \in \mathbb{R}^d$, on définit

$$x^{(n+1)} = \varphi(x^{(n)}).$$

Théorème. La suite $(x^{(n)})$ converge vers la solution de Ax = b pour b quelconque ssi $\rho(M^{-1}N) < 1$.



D'autres méthodes de résolution de systèmes linéaires

- ▶ But : construire une suite $(x^{(n)})$ qui converge vers x t.q. Ax = b.
- ▶ Idée : si A = M N avec M inversible, alors

$$Ax = b \iff x = \underbrace{M^{-1}(Nx + b)}_{\varphi(x)}.$$

Problème de point fixe : soit $x^{(0)} \in \mathbb{R}^d$, on définit

$$x^{(n+1)} = \varphi(x^{(n)}).$$

Théorème. La suite $(x^{(n)})$ converge vers la solution de Ax = b pour b quelconque ssi $\rho(M^{-1}N) < 1$.

Preuve. On pose
$$e^{(n)} = x - x^{(n)}$$
: $e^{(n+1)} = M^{-1}Ne^{(n)}$ donc $e^{(n)} = (M^{-1}N)^n e^{(0)}$.

Le cas où $M^{-1}N$ est diagonalisable est clair.



D'autres méthodes de résolution de systèmes linéaires

- ▶ But : construire une suite $(x^{(n)})$ qui converge vers x t.q. Ax = b.
- ldée : si A = M N avec M inversible, alors

$$Ax = b \iff x = \underbrace{M^{-1}(Nx + b)}_{\varphi(x)}.$$

Problème de point fixe : soit $x^{(0)} \in \mathbb{R}^d$, on définit

$$x^{(n+1)} = \varphi(x^{(n)}).$$

Théorème. La suite $(x^{(n)})$ converge vers la solution de Ax = b pour b quelconque ssi $\rho(M^{-1}N) < 1$.

Remarque. Si $\rho(M^{-1}N) < 1$ alors $\exists ||| \cdot ||| \text{ t.q. } |||M^{-1}N||| < 1$.

Rappel. $\rho(M^{-1}N) \leq |||M^{-1}N|||$ pour toute norme subordonnée.



D'autres méthodes de résolution de systèmes linéaires

$$x^{(n+1)} = M^{-1}(Nx^{(n)} + b).$$

En pratique.



D'autres méthodes de résolution de systèmes linéaires

$$x^{(n+1)} = M^{-1}(Nx^{(n)} + b).$$

En pratique.

► Système linéaire $Mx^{(n+1)} = Nx^{(n)} + b$ à chaque itération.



D'autres méthodes de résolution de systèmes linéaires

$$x^{(n+1)} = M^{-1}(Nx^{(n)} + b).$$

En pratique.

- ► Système linéaire $Mx^{(n+1)} = Nx^{(n)} + b$ à chaque itération.
- ► *M* doit être « facile à inverser ».



D'autres méthodes de résolution de systèmes linéaires

$$x^{(n+1)} = M^{-1}(Nx^{(n)} + b).$$

En pratique.

- ► Système linéaire $Mx^{(n+1)} = Nx^{(n)} + b$ à chaque itération.
- M doit être « facile à inverser ».
 - ► Méthode de Jacobi

M = diagonale de A si non nulle!



D'autres méthodes de résolution de systèmes linéaires

$$x^{(n+1)} = M^{-1}(Nx^{(n)} + b).$$

En pratique.

- Système linéaire $Mx^{(n+1)} = Nx^{(n)} + b$ à chaque itération.
- M doit être « facile à inverser ».
 - Méthode de JacobiM = diagonale de A si non nulle!
 - Méthode de Gauss-Seidel
 M = partie triangulaire inférieure de A si diagonale non nulle!



D'autres méthodes de résolution de systèmes linéaires

$$x^{(n+1)} = M^{-1}(Nx^{(n)} + b).$$

En pratique.

- Système linéaire $Mx^{(n+1)} = Nx^{(n)} + b$ à chaque itération.
- M doit être « facile à inverser ».
 - Méthode de Jacobi
 M = diagonale de A si non nulle!
 - Méthode de Gauss-Seidel
 M = partie triangulaire inférieure de A si diagonale non nulle!
- ▶ Si $\rho(M^{-1}N)$ < 1, convergence en $\mathcal{O}\left(\rho(M^{-1}N)^n\right)$.



Analyse numérique matricielle Recherche d'éléments propres

Problématique : Rechercher quelques vp de A, et les \overrightarrow{vp} associés.



Recherche d'éléments propres

Problématique: Rechercher quelques vp de A, et les \overrightarrow{vp} associés.

► Constat. si
$$A = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$$
 avec $\lambda_1 > \lambda_2 > 0$, et $x^0 = \begin{pmatrix} a \\ b \end{pmatrix}$,



Recherche d'éléments propres

Problématique: Rechercher quelques vp de A, et les \overrightarrow{vp} associés.

► Constat. si
$$A = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$$
 avec $\lambda_1 > \lambda_2 > 0$, et $x^0 = \begin{pmatrix} a \\ b \end{pmatrix}$,
$$A^n x^0 \sim a \lambda_1^n \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad \text{si } a \neq 0.$$



Recherche d'éléments propres

Problématique: Rechercher quelques vp de A, et les \overrightarrow{vp} associés.

► Constat. si
$$A = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$$
 avec $\lambda_1 > \lambda_2 > 0$, et $x^0 = \begin{pmatrix} a \\ b \end{pmatrix}$,

$$A^n x^0 \sim a \lambda_1^n \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$
, si $a \neq 0$.

Ainsi,

$$\frac{A^n x^0}{\|A^n x^0\|} \sim \mathbf{e_1}$$
 avec $A\mathbf{e_1} = \lambda_1 \mathbf{e_1}$!



Recherche d'éléments propres

Problématique: Rechercher quelques vp de A, et les \overrightarrow{vp} associés.

► Constat. si
$$A = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$$
 avec $\lambda_1 > \lambda_2 > 0$, et $x^0 = \begin{pmatrix} a \\ b \end{pmatrix}$,

$$A^n x^0 \sim a \lambda_1^n \left(egin{array}{c} 1 \\ 0 \end{array}
ight), \qquad {
m si} \ a
eq 0.$$

Ainsi,

$$\frac{A^n x^0}{\|A^n x^0\|} \sim \mathbf{e_1}$$
 avec $A\mathbf{e_1} = \lambda_1 \mathbf{e_1}$!

▶ Algorithme. Soit $x^0 \in \mathbb{R}^d$, on construit

$$y^{n+1} = Ax^n$$
 et $x^{n+1} = \frac{y^{n+1}}{\|y^{n+1}\|}$.



La méthode de la puissance

Algorithme. Soit $x^0 \in \mathbb{R}^d$, on construit

$$y^{n+1} = Ax^n$$
 et $x^{n+1} = \frac{y^{n+1}}{\|y^{n+1}\|}$.



La méthode de la puissance

Algorithme. Soit $x^0 \in \mathbb{R}^d$, on construit

$$y^{n+1} = Ax^n$$
 et $x^{n+1} = \frac{y^{n+1}}{\|y^{n+1}\|}$.

Théorème. Soit A de vp (complexes) $|\lambda_1| \leq |\lambda_2| \leq \cdots \leq |\lambda_{d-1}| < |\lambda_d|$ Pour presque[†] tout $x^0 \in \mathbb{R}^n$, la suite $(Ax^n|x^n)$ converge vers λ_d . De plus, la suite

$$(\tilde{x}^n) := \left(\frac{|\lambda_d|}{\lambda_d}\right)^n x^n$$

 $(\tilde{x}^n):=\left(\frac{|\lambda_d|}{\lambda_d}\right)^nx^n$ converge vers un vecteur propre associé à λ_d .



La méthode de la puissance

Algorithme. Soit $x^0 \in \mathbb{R}^d$, on construit

$$y^{n+1} = Ax^n$$
 et $x^{n+1} = \frac{y^{n+1}}{\|y^{n+1}\|}$.

Théorème. Soit A de vp (complexes) $|\lambda_1| \leq |\lambda_2| \leq \cdots \leq |\lambda_{d-1}| < |\lambda_d|$

Pour presque[†] tout $x^0 \in \mathbb{R}^n$, la suite $(Ax^n|x^n)$ converge vers λ_d . De plus, la suite

$$\left(\tilde{x}^{n}\right):=\left(\frac{\left|\lambda_{d}\right|}{\lambda_{d}}\right)^{n}x^{n}$$

converge vers un vecteur propre associé à λ_d .

[†] La condition est : la projection de x^0 sur le sous-espace propre associé à λ_1 n'est pas 0.

Attention à l'hypothèse $|\lambda_{d-1}| < |\lambda_d|$ qui est nécessaire.





COURS 2

Intégration numérique



Intégration numérique

Modélisation : calcul d'un temps de parcours



Deux toboggans... deux vitesses...



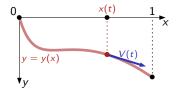


Quel temps de descente pour chaque toboggan?



Intégration numérique

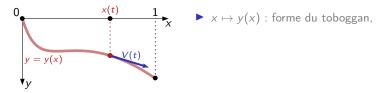
Modélisation : calcul d'un temps de parcours





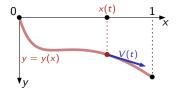
Intégration numérique

Modélisation : calcul d'un temps de parcours





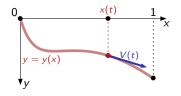
Modélisation : calcul d'un temps de parcours



- $x \mapsto y(x)$: forme du toboggan,
- $ightharpoonup t\mapsto ig(x(t),y(x(t)ig) : {\sf position},$



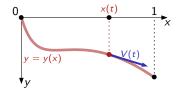
Modélisation : calcul d'un temps de parcours



- $ightharpoonup x\mapsto y(x)$: forme du toboggan,
- ▶ $t \mapsto (x(t), y(x(t)))$: position,
- T : temps de descente,



Modélisation : calcul d'un temps de parcours



$$ightharpoonup x \mapsto y(x)$$
: forme du toboggan,

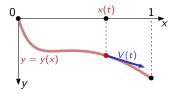
▶
$$t \mapsto (x(t), y(x(t)))$$
 : position,

T: temps de descente.

$$|V(t)| = |x'(t)|\sqrt{1 + y'^2(x(t))}$$



Modélisation : calcul d'un temps de parcours



- $ightharpoonup x \mapsto y(x)$: forme du toboggan,
- $ightharpoonup t\mapsto (x(t),y(x(t)))$: position,
- T: temps de descente,

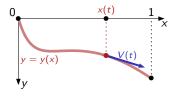
Vitesse:
$$V(t) = \begin{bmatrix} x'(t) \\ y'(x(t))x'(t) \end{bmatrix}$$
, $|V(t)| = |x'(t)|\sqrt{1 + y'^2(x(t))}$

$$|V(t)| = |x'(t)|\sqrt{1 + y'^2(x(t))}$$

Énergies :
$$E_c(t) = \frac{1}{2}mx'(t)^2 \left[1 + y'(x(t))^2\right], \quad E_p(t) = -mgy(x(t)).$$



Modélisation : calcul d'un temps de parcours



- $\triangleright x \mapsto y(x)$: forme du toboggan,
- $ightharpoonup t\mapsto (x(t),y(x(t)))$: position,
- T: temps de descente,

Vitesse:
$$V(t) = \begin{bmatrix} x'(t) \\ y'(x(t))x'(t) \end{bmatrix}$$
, $|V(t)| = |x'(t)|\sqrt{1 + y'^2(x(t))}$

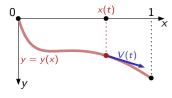
$$|V(t)| = |x'(t)|\sqrt{1 + y'^2(x(t))}$$

$$\mathsf{E}_{\mathsf{nergies}} : E_{\mathsf{c}}(t) = \frac{1}{2} m x'(t)^2 \left[1 + y'(x(t))^2 \right], \quad E_{\mathsf{p}}(t) = - m g y(x(t)).$$

Bilan entre
$$t = 0$$
 et $t = T$: $\frac{1}{2}m(x'(t))^2[1 + y'(x(t))^2] - mgy(x(t)) = 0$.
 $\implies x'(t) = \sqrt{2gy(x(t))}/\sqrt{1 + y'(x(t))^2}$.



Modélisation : calcul d'un temps de parcours



- $\triangleright x \mapsto y(x)$: forme du toboggan,
- \blacktriangleright $t \mapsto (x(t), y(x(t)) : position,$
- T: temps de descente,

$$|V(t)| = |x'(t)|\sqrt{1 + y'^2(x(t))}$$

Énergies :
$$E_c(t) = \frac{1}{2}mx'(t)^2 \left[1 + y'(x(t))^2\right], \quad E_p(t) = -mgy(x(t)).$$

Bilan entre
$$t = 0$$
 et $t = T$: $\frac{1}{2}m(x'(t))^2[1 + y'(x(t))^2] - mgy(x(t)) = 0$.
 $\implies x'(t) = \sqrt{2gy(x(t))}/\sqrt{1 + y'(x(t))^2}$.

Temps de descente :
$$T = \int_0^T dt = \int_0^1 \frac{\sqrt{1 + y'(x)^2}}{\sqrt{2gy(x)}} dx$$
.



Modélisation : calcul d'un temps de parcours

Problématique. La forme $x\mapsto y(x)$ du toboggan étant donnée, comment calculer le temps de descente T? $T=\int_0^1 \frac{\sqrt{1+y'(x)^2}}{\sqrt{2gy(x)}}\,\mathrm{d}x.$

$$T = \int_0^1 \frac{\sqrt{1 + y'(x)^2}}{\sqrt{2gy(x)}} dx.$$



Modélisation : calcul d'un temps de parcours

Problématique. La forme $x\mapsto y(x)$ du toboggan étant donnée, comment calculer le temps de descente T? $T=\int_0^1 \frac{\sqrt{1+y'(x)^2}}{\sqrt{2gy(x)}}\,\mathrm{d}x.$

$$T = \int_0^1 \frac{\sqrt{1 + y'(x)^2}}{\sqrt{2gy(x)}} \, \mathrm{d}x$$

Problème général : calcul de $\int_a^b f(x) dx$.



Modélisation : calcul d'un temps de parcours

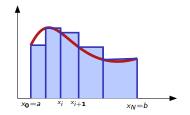
Problématique. La forme $x\mapsto y(x)$ du toboggan étant donnée, comment calculer le temps de descente T?

$$T = \int_0^1 \frac{\sqrt{1 + y'(x)^2}}{\sqrt{2gy(x)}} dx.$$

- Problème général : calcul de $\int_a^b f(x) dx$.
- ► Enjeux :
 - ► Mise en place d'une méthode numérique,
 - Étude de la précision,
 - Analyse du temps de calcul,
 - Quelle méthode pour quelle fonction f?

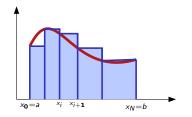


Une première idée : la méthode des rectangles





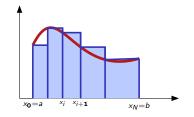
Une première idée : la méthode des rectangles



$$\int_a^b f(x) dx \simeq \sum_{i=0}^{N-1} (x_{i+1} - x_i) f(x_i).$$



Une première idée : la méthode des rectangles



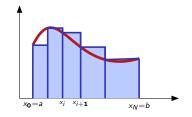
$$\int_{a}^{b} f(x) dx \simeq \sum_{i=0}^{N-1} (x_{i+1} - x_{i}) f(x_{i}).$$

► Cas équidistant : $x_{i+1} - x_i = h = \frac{b-a}{N}$,

$$\int_a^b f(x) dx \simeq h \sum_{i=0}^{N-1} f(x_i).$$



Une première idée : la méthode des rectangles



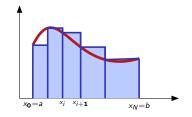
$$\int_{a}^{b} f(x) dx \simeq \sum_{i=0}^{N-1} (x_{i+1} - x_{i}) f(x_{i}).$$

► Cas équidistant : $x_{i+1} - x_i = h = \frac{b-a}{N}$,

$$\int_a^b f(x) dx \simeq h \sum_{i=0}^{N-1} f(x_i).$$

▶ Est-ce que l'approximation converge lorsque $N \to +\infty$?

Une première idée : la méthode des rectangles



$$\int_{a}^{b} f(x) dx \simeq \sum_{i=0}^{N-1} (x_{i+1} - x_i) f(x_i).$$

► Cas équidistant : $x_{i+1} - x_i = h = \frac{b-a}{N}$,

$$\int_a^b f(x) dx \simeq h \sum_{i=0}^{N-1} f(x_i).$$

- **E**st-ce que l'approximation converge lorsque $N \to +\infty$?
- Si oui, avec quelle précision?



Convergence de la méthode des rectangles

► Cas équidistant (rappel : Nh = b - a)

$$I_N(f) = h \sum_{i=0}^{N-1} f(x_i).$$



Convergence de la méthode des rectangles

► Cas équidistant (rappel : Nh = b - a)

$$I_N(f) = h \sum_{i=0}^{N-1} f(x_i).$$

▶ Si f est continue, $I_N(f)$ est une somme de Riemann.

$$\Longrightarrow I_N(f) \xrightarrow[N \to +\infty]{} \int_a^b f(x) dx.$$



Convergence de la méthode des rectangles

ightharpoonup Cas équidistant (rappel : Nh = b - a)

$$I_N(f) = h \sum_{i=0}^{N-1} f(x_i).$$

▶ Si f est continue, $I_N(f)$ est une somme de Riemann.

$$\Longrightarrow I_N(f) \xrightarrow[N \to +\infty]{} \int_a^b f(x) dx.$$

C'est lié à l'approximation uniforme des fonctions continues par des fonctions en escalier / à la continuité uniforme des fonctions continues sur un segment.



Convergence de la méthode des rectangles

ightharpoonup Cas équidistant (rappel : Nh = b - a)

$$I_N(f) = h \sum_{i=0}^{N-1} f(x_i).$$

▶ Si f est continue, $I_N(f)$ est une somme de Riemann.

$$\Longrightarrow I_N(f) \xrightarrow[N \to +\infty]{} \int_a^b f(x) dx.$$

C'est lié à l'approximation uniforme des fonctions continues par des fonctions en escalier / à la continuité uniforme des fonctions continues sur un segment.

▶ Peut-on quantifier la précision de l'approximation ?



Précision de la méthode des rectangles

$$I = \int_a^b f(x) dx, \qquad I_N = h \sum_{i=0}^{N-1} f(a+ih).$$



Précision de la méthode des rectangles

$$I = \int_a^b f(x) dx, \qquad I_N = h \sum_{i=0}^{N-1} f(a+ih).$$

Erreur:

$$|I - I_N| = \left| \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} \left[f(x) - f(x_i) \right] dx \right|.$$



Précision de la méthode des rectangles

$$I = \int_{a}^{b} f(x) dx, \qquad I_{N} = h \sum_{i=0}^{N-1} f(a+ih).$$

► Erreur :

$$|I - I_N| = \left| \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} \left[f(x) - f(x_i) \right] dx \right|.$$

▶ Inégalité des accroissements finis : si $f \in \mathscr{C}^1([a,b])$,

$$|f(x)-f(x_i)|\leqslant M_1|x-x_i|$$
 avec $M_1=\sup_{t\in[a,b]}|f'(t)|.$



Précision de la méthode des rectangles

$$I = \int_{a}^{b} f(x) dx, \qquad I_{N} = h \sum_{i=0}^{N-1} f(a+ih).$$

► Erreur :

$$|I - I_N| = \left| \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} \left[f(x) - f(x_i) \right] dx \right|.$$

▶ Inégalité des accroissements finis : si $f \in \mathscr{C}^1([a,b])$,

$$|f(x) - f(x_i)| \le M_1 |x - x_i|$$
 avec $M_1 = \sup_{t \in [a,b]} |f'(t)|$.

Estimation d'erreur

$$|I - I_N| \leqslant \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} |f(x) - f(x_i)| dx$$



Précision de la méthode des rectangles

$$I = \int_{a}^{b} f(x) dx, \qquad I_{N} = h \sum_{i=0}^{N-1} f(a+ih).$$

► Erreur :

$$|I - I_N| = \left| \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} \left[f(x) - f(x_i) \right] dx \right|.$$

▶ Inégalité des accroissements finis : si $f \in \mathscr{C}^1([a,b])$,

$$|f(x) - f(x_i)| \le M_1 |x - x_i|$$
 avec $M_1 = \sup_{t \in [a,b]} |f'(t)|$.

Estimation d'erreur

$$|I - I_N| \leqslant M_1 \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} |x - x_i| dx$$



Précision de la méthode des rectangles

$$I = \int_{a}^{b} f(x) dx, \qquad I_{N} = h \sum_{i=0}^{N-1} f(a+ih).$$

► Erreur :

$$|I - I_N| = \left| \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} \left[f(x) - f(x_i) \right] dx \right|.$$

▶ Inégalité des accroissements finis : si $f \in \mathscr{C}^1([a,b])$,

$$|f(x) - f(x_i)| \le M_1 |x - x_i|$$
 avec $M_1 = \sup_{t \in [a,b]} |f'(t)|$.

Estimation d'erreur

$$|I - I_N| \leqslant M_1 \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} (x - x_i) dx$$



Précision de la méthode des rectangles

$$I = \int_{a}^{b} f(x) dx, \qquad I_{N} = h \sum_{i=0}^{N-1} f(a+ih).$$

► Erreur :

$$|I - I_N| = \left| \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} \left[f(x) - f(x_i) \right] dx \right|.$$

▶ Inégalité des accroissements finis : si $f \in \mathscr{C}^1([a,b])$,

$$|f(x) - f(x_i)| \le M_1 |x - x_i|$$
 avec $M_1 = \sup_{t \in [a,b]} |f'(t)|$.

Estimation d'erreur

$$|I - I_N| \le M_1 \sum_{i=0}^{N-1} \left[\frac{(x - x_i)^2}{2} \right]_{x_i}^{x_{i+1}} dx$$



Précision de la méthode des rectangles

$$I = \int_{a}^{b} f(x) dx, \qquad I_{N} = h \sum_{i=0}^{N-1} f(a+ih).$$

► Erreur :

$$|I - I_N| = \left| \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} \left[f(x) - f(x_i) \right] dx \right|.$$

▶ Inégalité des accroissements finis : si $f \in \mathcal{C}^1([a,b])$,

$$|f(x) - f(x_i)| \le M_1 |x - x_i|$$
 avec $M_1 = \sup_{t \in [a,b]} |f'(t)|$.

► Estimation d'erreur

$$|I-I_N|\leqslant M_1\sum_{i=0}^{N-1}\frac{h^2}{2}$$



Précision de la méthode des rectangles

$$I = \int_{a}^{b} f(x) dx, \qquad I_{N} = h \sum_{i=0}^{N-1} f(a+ih).$$

► Erreur :

$$|I - I_N| = \left| \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} \left[f(x) - f(x_i) \right] dx \right|.$$

▶ Inégalité des accroissements finis : si $f \in \mathscr{C}^1([a,b])$,

$$|f(x) - f(x_i)| \leqslant M_1 |x - x_i|$$
 avec $M_1 = \sup_{t \in [a,b]} |f'(t)|$.

Estimation d'erreur

$$|I-I_N|\leqslant M_1N\frac{h^2}{2}$$



Précision de la méthode des rectangles

$$I = \int_{a}^{b} f(x) dx, \qquad I_{N} = h \sum_{i=0}^{N-1} f(a+ih).$$

► Erreur :

$$|I-I_N| = \left|\sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} \left[f(x) - f(x_i) \right] dx \right|.$$

▶ Inégalité des accroissements finis : si $f \in \mathscr{C}^1([a,b])$,

$$|f(x) - f(x_i)| \le M_1 |x - x_i|$$
 avec $M_1 = \sup_{t \in [a,b]} |f'(t)|$.

Estimation d'erreur

$$|I-I_N|\leqslant \frac{M_1(b-a)}{2}h$$



Précision de la méthode des rectangles

$$I = \int_{a}^{b} f(x) dx, \qquad I_{N} = h \sum_{i=0}^{N-1} f(a+ih).$$

► Erreur :

$$|I - I_N| = \left| \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} \left[f(x) - f(x_i) \right] dx \right|.$$

▶ Inégalité des accroissements finis : si $f \in \mathscr{C}^1([a,b])$,

$$|f(x) - f(x_i)| \leqslant M_1 |x - x_i|$$
 avec $M_1 = \sup_{t \in [a,b]} |f'(t)|$.

Estimation d'erreur

$$|I-I_N|\leqslant \frac{M_1(b-a)^2}{2N}$$



Précision de la méthode des rectangles

Si
$$f \in \mathscr{C}^1([a,b])$$
, alors, pour la méthode des rectangles,
$$|I-I_N| \leqslant \frac{(b-a)^2}{2N} \sup_{t \in [a,b]} |f'(t)|$$



Précision de la méthode des rectangles

Si
$$f \in \mathscr{C}^1([a,b])$$
, alors, pour la méthode des rectangles,
$$|I-I_N| \leqslant \frac{(b-a)^2}{2N} \sup_{t \in [a,b]} |f'(t)|$$

À retenir :

- ▶ Hypothèse de régularité : $f \in \mathcal{C}^1([a, b])$.
- ► Convergence en $\mathcal{O}\left(\frac{1}{N}\right)$.



Précision de la méthode des rectangles

Si
$$f \in \mathscr{C}^1([a,b])$$
, alors, pour la méthode des rectangles,
$$|I-I_N| \leqslant \frac{(b-a)^2}{2N} \sup_{t \in [a,b]} |f'(t)|$$

À retenir :

- ▶ Hypothèse de régularité : $f \in \mathscr{C}^1([a, b])$.
- ► Convergence en $\mathcal{O}\left(\frac{1}{N}\right)$.

Question : que se passe-t-il si $f \notin \mathcal{C}^1([a,b])$?

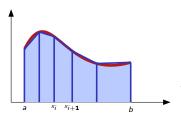
- Simulation pour diverses fonctions.

► Analyse théorique : cf. TD.



Une autre méthode

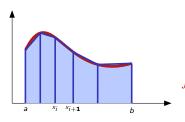
► Méthode des trapèzes



$$\int_a^b f(x) dx \simeq \sum_{i=0}^{N-1} (x_{i+1} - x_i) \frac{f(x_i) + f(x_{i+1})}{2}.$$

Une autre méthode

Méthode des trapèzes



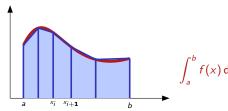
$$\int_a^b f(x) dx \simeq \sum_{i=0}^{N-1} (x_{i+1} - x_i) \frac{f(x_i) + f(x_{i+1})}{2}.$$

▶ Erreur observée dans le cas général : $\mathcal{O}(N^{-2})$.



Une autre méthode

Méthode des trapèzes



$$\int_a^b f(x) dx \simeq \sum_{i=0}^{N-1} (x_{i+1} - x_i) \frac{f(x_i) + f(x_{i+1})}{2}.$$

- ▶ Erreur observée dans le cas général : $\mathcal{O}(N^{-2})$.
- Estimation d'erreur théorique?

 \rightarrow estimation sur $[x_i, x_{i+1}]$ de la différence courbe/corde.

Intégration numérique Méthode des trapèzes

Estimation d'erreur (subdivision uniforme de pas h)



Intégration numérique Méthode des trapèzes

Estimation d'erreur (subdivision uniforme de pas h)

$$|I - I_N| \le \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} |f(x) - p(x)| dx,$$

où
$$p \in \mathbb{P}_1$$
 satisfait $p(x_i) = f(x_i)$ et $p(x_{i+1}) = f(x_{i+1})$.



Méthode des trapèzes

Estimation d'erreur (subdivision uniforme de pas h)

$$|I - I_N| \le \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} |f(x) - p(x)| dx,$$

où
$$p \in \mathbb{P}_1$$
 satisfait $p(x_i) = f(x_i)$ et $p(x_{i+1}) = f(x_{i+1})$.

▶ On pose $\psi = f - p$. Th. Rolle : $\exists \xi_i \in [x_i, x_{i+1}], \ \psi'(\xi_i) = 0$.

$$\forall x \in [x_i, x_{i+1}], \ \big|\psi'(x)\big| \leqslant |x - \xi_i| \sup_{[x_i, x_{i+1}]} \big|\psi''\big| \leqslant h \sup_{[a, b]} \big|f''\big|.$$



Méthode des trapèzes

Estimation d'erreur (subdivision uniforme de pas h)

$$|I - I_N| \le \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} |f(x) - p(x)| dx,$$

où $p \in \mathbb{P}_1$ satisfait $p(x_i) = f(x_i)$ et $p(x_{i+1}) = f(x_{i+1})$.

- ▶ On pose $\psi = f p$. Th. Rolle : $\exists \xi_i \in [x_i, x_{i+1}], \ \psi'(\xi_i) = 0$. $\forall x \in [x_i, x_{i+1}], \ \big|\psi'(x)\big| \leqslant |x \xi_i| \sup_{[x_i, x_{i+1}]} \big|\psi''\big| \leqslant h \sup_{[a,b]} \big|f''\big|.$
- ▶ Donc $\forall x \in [x_i, x_{i+1}]$,

$$|\psi(x)| \leqslant \int_{x_i}^x |\psi'(x)| \leqslant h^2 \sup_{[a,b]} |f''|.$$



Méthode des trapèzes

Estimation d'erreur (subdivision uniforme de pas h)

$$|I - I_N| \le \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} |f(x) - p(x)| dx,$$

où $p \in \mathbb{P}_1$ satisfait $p(x_i) = f(x_i)$ et $p(x_{i+1}) = f(x_{i+1})$.

- ▶ On pose $\psi = f p$. Th. Rolle : $\exists \xi_i \in [x_i, x_{i+1}], \ \psi'(\xi_i) = 0$. $\forall x \in [x_i, x_{i+1}], \ |\psi'(x)| \leq |x \xi_i| \sup_{[x_i, x_{i+1}]} |\psi''| \leq h \sup_{[a,b]} |f''|$.
- ▶ Donc $\forall x \in [x_i, x_{i+1}]$,

$$|\psi(x)| \leqslant \int_{x_i}^x |\psi'(x)| \leqslant h^2 \sup_{[a,b]} |f''|.$$

► Conclusion : $|I - I_N| \le h^2 (b - a) \sup_{[a,b]} |f''| = \frac{(b - a)^3}{N^2} M_2$.



D'autres méthodes

lacktriangle Idée : remplacer arphi par un polynôme π_{arphi} et

$$\int_0^1 arphi(t) \, \mathrm{d}t \simeq \int_0^1 \pi_arphi(t) \, \mathrm{d}t.$$



D'autres méthodes

ightharpoonup Idée : remplacer arphi par un polynôme π_{arphi} et

$$\int_0^1 arphi(t) \, \mathrm{d}t \simeq \int_0^1 \pi_arphi(t) \, \mathrm{d}t.$$

▶ Rappel: pour $0 \leqslant t_0 < t_1 < \dots < t_k \leqslant 1$ et $b_0, b_1, \dots, b_k \in \mathbb{R}$ fixés, $\exists ! p \in \mathbb{P}_k, \quad \forall q = 0, 1, \dots, k, \quad p(t_q) = b_q.$

Dans la base de Lagrange,

$$ho(t) = \sum_{q=0}^k b_q L_q(t), \quad ext{avec} \quad L_q(t) = \prod_{\ell
eq q} rac{t-t_\ell}{t_q-t_\ell}.$$



D'autres méthodes

ldée : remplacer φ par un polynôme π_{φ} et

$$\int_0^1 arphi(t) \, \mathrm{d}t \simeq \int_0^1 \pi_arphi(t) \, \mathrm{d}t.$$

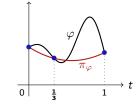
▶ Rappel: pour $0 \le t_0 < t_1 < \dots < t_k \le 1$ et $b_0, b_1, \dots, b_k \in \mathbb{R}$ fixés, $\exists ! p \in \mathbb{P}_k, \quad \forall q = 0, 1, \dots, k, \quad p(t_q) = b_q.$

Dans la base de Lagrange,

$$p(t) = \sum_{q=0}^k b_q L_q(t), \quad ext{avec} \quad L_q(t) = \prod_{\ell
eq q} rac{t-t_\ell}{t_q-t_\ell}.$$

Exemple:

 $\pi_{arphi}=$ polynôme d'interpolation $\mbox{de Lagrange de } arphi$ en 0, $\frac{1}{3}$, 1.



Intermède : interpolation de Lagrange

$$\int_0^1 arphi(t) \, \mathrm{d}t \simeq \int_0^1 \pi_arphi(t) \, \mathrm{d}t.$$



Intermède : interpolation de Lagrange

$$\int_0^1 arphi(t) \, \mathrm{d}t \simeq \int_0^1 \pi_arphi(t) \, \mathrm{d}t.$$

l n'est pas nécessaire d'expliciter π_{φ} pour approcher l'intégrale :

$$\pi_{arphi}(t) = \sum_{q=0}^k arphi(t_q) L_q(t) \quad \Longrightarrow \quad \int_0^1 \pi_{arphi}(t) \, \mathrm{d}t = \sum_{q=0}^k \underbrace{\left(\int_0^1 L_q(t) \, \mathrm{d}t
ight)}_{=w_q} arphi(t_q)$$



Intermède : interpolation de Lagrange

$$\int_0^1 arphi(t) \, \mathrm{d}t \simeq \int_0^1 \pi_arphi(t) \, \mathrm{d}t.$$

▶ Il n'est pas nécessaire d'expliciter π_{φ} pour approcher l'intégrale :

$$\pi_{\varphi}(t) = \sum_{q=0}^{k} \varphi(t_q) L_q(t) \quad \Longrightarrow \quad \int_0^1 \pi_{\varphi}(t) dt = \sum_{q=0}^{k} \underbrace{\left(\int_0^1 L_q(t) dt\right)}_{=w_q} \varphi(t_q)$$

d'où l'approximation

$$\int_0^1 \varphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^k w_q \varphi(t_q).$$



Intermède : interpolation de Lagrange



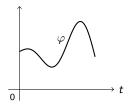
Intermède : interpolation de Lagrange

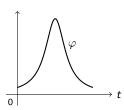
▶ Que se passe-t-il lorsque $k \to \infty$? (points (t_q) équidistants)

$$\varphi(t) = 1 + \frac{t}{2}\cos(8t)$$
 $\qquad \qquad \varphi(t) = \frac{2}{1 + 10(2t - 1)^2}$



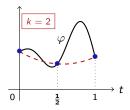
Intermède : interpolation de Lagrange

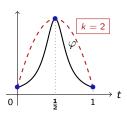






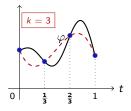
Intermède : interpolation de Lagrange

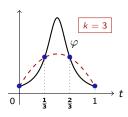






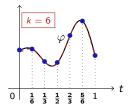
Intermède : interpolation de Lagrange

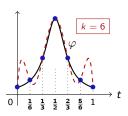






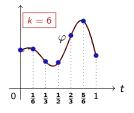
Intermède : interpolation de Lagrange

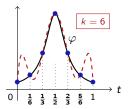






Intermède : interpolation de Lagrange



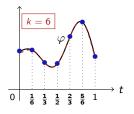


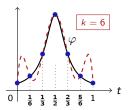
Non convergence lorsque $k \to +\infty$ (phénomène de Runge)



Intermède : interpolation de Lagrange

▶ Que se passe-t-il lorsque $k \to \infty$? (points (t_q) équidistants)



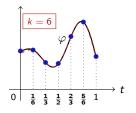


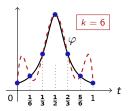
Non convergence lorsque $k \to +\infty$ (phénomène de Runge)

▶ Pas d'espoir d'obtenir une approximation qui converge lorsque $k \to +\infty$.



Intermède : interpolation de Lagrange





Non convergence lorsque $k \to +\infty$ (phénomène de Runge)

- ▶ Pas d'espoir d'obtenir une approximation qui converge lorsque $k \to +\infty$.
- ▶ On va conserver *k* fixé et travailler sur une subdivision de l'intervalle.



$$\int_a^b f(x) \, \mathrm{d}x$$

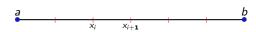




Principe des méthodes composées

$$\int_{a}^{b} f(x) dx = \sum_{i=0}^{N-1} \int_{x_{i}}^{x_{i+1}} f(x) dx$$

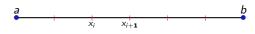
▶ On subdivise l'intervalle [a, b] (pas uniforme h).





$$\int_{a}^{b} f(x) dx = \sum_{i=0}^{N-1} \int_{x_{i}}^{x_{i+1}} f(x) dx$$

- On subdivise l'intervalle [a, b] (pas uniforme h).
- On choisit un modèle élémentaire sur [0,1]

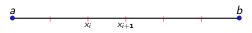






$$\int_{a}^{b} f(x) dx = \sum_{i=0}^{N-1} \int_{x_{i}}^{x_{i+1}} f(x) dx$$

- ▶ On subdivise l'intervalle [a, b] (pas uniforme h).
- $lackbox{ On choisit un modèle élémentaire sur } [0,1]: \int_0^1 arphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^k w_q arphi(t_q)$

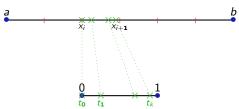






$$\int_a^b f(x) dx = \sum_{i=0}^{N-1} \int_{x_i}^{x_{i+1}} f(x) dx = \sum_{i=0}^{N-1} h \int_0^1 f(x_i + th) dt \simeq h \sum_{i=0}^{N-1} \sum_{q=0}^k w_q f(x_i + t_q h).$$

- On subdivise l'intervalle [a, b] (pas uniforme h).
- $lackbox{ On choisit un modèle élémentaire sur } [0,1]: \int_0^1 arphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^\kappa w_q arphi(t_q)$
- On transporte sur chaque sous-intervalle.





Méthodes composées

Exemple:

$$\blacktriangleright$$
 $(k = 1)$. $t_0 = 0$, $t_1 = 1$,



Méthodes composées

Exemple:

- \blacktriangleright (k=1). $t_0=0$, $t_1=1$,
- $ightharpoonup w_0 = \frac{1}{2}, \ w_1 = \frac{1}{2}.$



Méthodes composées

Exemple:

- \blacktriangleright (k=1). $t_0=0$, $t_1=1$,
- $ightharpoonup w_0 = \frac{1}{2}, \ w_1 = \frac{1}{2}.$

$$\int_{a}^{b} f(x) dx \simeq h \sum_{i=0}^{N-1} \sum_{q=0}^{k} w_{q} f(x_{i} + t_{q} h).$$

$$= h \sum_{i=0}^{N-1} \frac{f(x_{i}) + f(x_{i+1})}{2}$$



Méthodes composées

Exemple:

- $(k=1). t_0=0, t_1=1,$
- $ightharpoonup w_0 = \frac{1}{2}, \ w_1 = \frac{1}{2}.$

$$\int_{a}^{b} f(x) dx \simeq h \sum_{i=0}^{N-1} \sum_{q=0}^{k} w_{q} f(x_{i} + t_{q} h).$$

$$= h \sum_{i=0}^{N-1} \frac{f(x_{i}) + f(x_{i+1})}{2}$$

⇒ Méthode des trapèzes!



Méthodes composées

- Bilan
 - ► Modèle élémentaire sur [0,1] :

$$(M_e)$$

$$\int_0^1 \varphi(t) dt \simeq \sum_{q=0}^k w_q \varphi(t_q).$$



Méthodes composées

- Bilan
 - ► Modèle élémentaire sur [0,1] :

$$(M_e) \qquad \qquad \int_0^1 \varphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

► Méthode composée sur [a, b] :

(FQ_c)
$$I = \int_a^b f(x) dx \simeq I_N = h \sum_{i=0}^{N-1} \sum_{q=0}^k w_q f(x_i + t_q h).$$



Méthodes composées

- Bilan
 - ► Modèle élémentaire sur [0,1] :

$$\int_0^1 arphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^k w_q arphi(t_q).$$

Méthode composée sur [a, b] :

(FQ_c)
$$I = \int_a^b f(x) dx \simeq I_N = h \sum_{i=0}^{N-1} \sum_{q=0}^k w_q f(x_i + t_q h).$$

A priori : N(k+1) évaluations de f.



Méthodes composées

- Bilan
 - ► Modèle élémentaire sur [0,1] :

$$(M_e) \qquad \qquad \int_0^1 \varphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

Méthode composée sur [a, b] :

(FQ_c)
$$I = \int_a^b f(x) dx \simeq I_N = h \sum_{i=0}^{N-1} \sum_{q=0}^k w_q f(x_i + t_q h).$$

- A priori : N(k+1) évaluations de f.
- ▶ Quelle précision pour (FQ_c) ?



Méthodes composées

Une remarque

$$I - I_N = h \sum_{i=0}^{N-1} \left(\int_0^1 \underbrace{f(x_i + th)}_{\varphi_i(t)} dt - \int_0^1 \pi_{\varphi_i}(t) dt \right),$$

- $\triangleright \varphi_i : t \mapsto f(x_i + th),$
- $ightharpoonup orall arphi, \ \pi_{arphi} \in \mathbb{P}_k \ ext{satisfait} \ \pi_{arphi}(t_q) = arphi(t_q) \ (0 \leqslant q \leqslant k).$



Méthodes composées

▶ Une remarque

$$I - I_N = h \sum_{i=0}^{N-1} \left(\int_0^1 \underbrace{f(x_i + th)}_{\varphi_i(t)} dt - \int_0^1 \pi_{\varphi_i}(t) dt \right),$$

- $\triangleright \varphi_i : t \mapsto f(x_i + th),$
- $ightharpoonup orall arphi, \ \pi_arphi \in \mathbb{P}_k \ ext{satisfait} \ \pi_arphi(t_q) = arphi(t_q) \ (0 \leqslant q \leqslant k).$

Théorème. pour $\varphi \in \mathscr{C}^{k+1}([0,1])$,

$$|arphi(t)-\pi_{arphi}(t)|\leqslant rac{\displaystyle\sup_{t\in[0,1]}\left|arphi^{(k+1)}(t)
ight|}{(k+1)!}.$$



Méthodes composées

▶ Une remarque

$$I - I_N = h \sum_{i=0}^{N-1} \left(\int_0^1 \underbrace{f(x_i + th)}_{\varphi_i(t)} dt - \int_0^1 \pi_{\varphi_i}(t) dt \right),$$

- $\qquad \qquad \varphi_i: t \mapsto f(x_i + th), \ \varphi_i^{(k+1)}(t) = h^{k+1} f^{(k+1)}(x_i + th).$
- $ightharpoonup orall arphi, \, \pi_{arphi} \in \mathbb{P}_k \, ext{ satisfait } \pi_{arphi}(t_q) = arphi(t_q) \, \, (0 \leqslant q \leqslant k).$

Théorème. pour $\varphi \in \mathscr{C}^{k+1}([0,1])$,

$$|arphi(t) - \pi_{arphi}(t)| \leqslant rac{\displaystyle \sup_{t \in [0,1]} \left| arphi^{(k+1)}(t)
ight|}{(k+1)!}.$$

► Conséquence : si $f \in \mathscr{C}^{k+1}[a,b]$,

$$|I - I_N| \leqslant \frac{h^{k+1}(b-a)}{(k+1)!} \sup_{x \in [a,b]} |f^{(k+1)}(x)|.$$



Intégration numérique Méthodes composées

$$(M_e)$$

$$\int_0^1 \varphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

(FQ_c)
$$I = \int_{a}^{b} f(x) dx \simeq I_{N} = h \sum_{i=0}^{N-1} \sum_{q=0}^{k} w_{q} f(x_{i} + t_{q} h).$$

Théorème. Si (M_e) est exacte pour tout $\varphi \in \mathbb{P}_\ell$ et si $f \in \mathscr{C}^{\ell+1}([a,b])$, alors

$$|I - I_N| \leqslant \frac{h^{\ell+1}(b-a)}{(\ell+1)!} \sup_{x \in [a,b]} \left| f^{(\ell+1)}(x) \right|.$$



Méthodes composées

$$(M_e)$$

$$\int_0^1 arphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^k w_q arphi(t_q).$$

(FQ_c)
$$I = \int_{a}^{b} f(x) dx \simeq I_{N} = h \sum_{i=0}^{N-1} \sum_{q=0}^{k} w_{q} f(x_{i} + t_{q} h).$$

Théorème. Si (M_e) est exacte pour tout $\varphi \in \mathbb{P}_\ell$ et si $f \in \mathscr{C}^{\ell+1}([a,b])$, alors

$$|I - I_N| \leqslant rac{h^{\ell+1}(b-a)}{(\ell+1)!} \sup_{x \in [a,b]} \left| f^{(\ell+1)}(x) \right|.$$

Remarque. La constante n'est pas optimale. . .

Méthodes composées

$$(M_e)$$

$$\int_0^1 \varphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

$$(FQ_c) I = \int_a^b f(x) dx \simeq I_N = h \sum_{i=0}^{N-1} \sum_{q=0}^k w_q f(x_i + t_q h).$$

Théorème. Si (M_e) est exacte pour tout $\varphi \in \mathbb{P}_\ell$ et si $f \in \mathscr{C}^{\ell+1}([a,b])$, alors

$$|I - I_N| \leqslant rac{h^{\ell+1}(b-a)}{(\ell+1)!} \sup_{x \in [a,b]} \left| f^{(\ell+1)}(x) \right|.$$

Remarque. La constante n'est pas optimale. . .

Question. A priori, $\ell = k$ N(k+1) évaluations de f, précision en $\mathcal{O}\left(\frac{1}{N^{k+1}}\right)$.



Méthodes composées

$$(M_e)$$

$$\int_0^1 \varphi(t) dt \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

(FQ_c)
$$I = \int_{a}^{b} f(x) dx \simeq I_{N} = h \sum_{i=0}^{N-1} \sum_{q=0}^{k} w_{q} f(x_{i} + t_{q} h).$$

Théorème. Si (M_e) est exacte pour tout $\varphi \in \mathbb{P}_\ell$ et si $f \in \mathscr{C}^{\ell+1}([a,b])$, alors

$$|I - I_N| \leqslant \frac{h^{\ell+1}(b-a)}{(\ell+1)!} \sup_{x \in [a,b]} \left| f^{(\ell+1)}(x) \right|.$$

Remarque. La constante n'est pas optimale. . .

Question. A priori, $\ell = k$ N(k+1) évaluations de f, précision en $\mathcal{O}\left(\frac{1}{N^{k+1}}\right)$.

Peut-on faire mieux?



Optimalité des choix des points/poids de quadrature?

Points de quadrature fixés : comment calculer les poids?

$$(M_e)$$

$$\int_0^1 \varphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^k w_q \varphi(t_q).$$



Optimalité des choix des points/poids de quadrature?

Points de quadrature fixés : comment calculer les poids?

$$(M_e)$$

$$\int_0^1 \varphi(t) dt \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

 \blacktriangleright k et $(t_q)_{q=0,...,k}$ donnés.



Optimalité des choix des points/poids de quadrature?

Points de quadrature fixés : comment calculer les poids?

$$(M_e)$$

$$\int_0^1 arphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^k w_q arphi(t_q).$$

- \triangleright k et $(t_q)_{q=0,...,k}$ donnés.
- (M_e) exact sur $\mathbb{P}_k \iff (M_e)$ exact pour la base $(1, t, \dots, t^k)$,

$$\iff \forall i=0,1,\ldots,k, \quad \sum_{q=0}^k w_q t_q^i = \frac{1}{i+1}.$$

(système linéaire de Vandermonde)

Optimalité des choix des points/poids de quadrature?

Points de quadrature fixés : comment calculer les poids?

$$\left(\mathit{M}_{e}
ight) \qquad \qquad \int_{0}^{1} arphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^{k} \mathit{w}_{q} arphi(t_{q}).$$

- \triangleright k et $(t_q)_{q=0,\ldots,k}$ donnés.
- (M_e) exact sur $\mathbb{P}_k \iff (M_e)$ exact pour la base $(1, t, \dots, t^k)$,

$$\iff \forall i=0,1,\ldots,k, \quad \sum_{q=0}^k w_q t_q^i = \frac{1}{i+1}.$$

(système linéaire de Vandermonde)

 \implies les poids (w_q) sont uniquement déterminés!



Optimalité des choix des points/poids de quadrature?

Points de quadrature fixés : comment calculer les poids?

$$(M_e)$$

$$\int_0^1 arphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^k w_q arphi(t_q).$$

- \triangleright k et $(t_q)_{q=0,...,k}$ donnés.
- (M_e) exact sur $\mathbb{P}_k \iff (M_e)$ exact pour la base $(1, t, \dots, t^k)$,

$$\iff \forall i=0,1,\ldots,k, \quad \sum_{q=0}^k w_q t_q^i = \frac{1}{i+1}.$$

(système linéaire de Vandermonde)

 \implies les poids (w_q) sont uniquement déterminés!

Conclusion. Sauf « miracle », (M_e) est exact sur \mathbb{P}_k seulement...



Optimalité des choix des points/poids de quadrature?

Avec k+1 points de quadrature‡ ${\Bbb R}$: peut-on être exact au delà de ${\Bbb P}_k$?

$$(M_e) \qquad \qquad \int_0^1 \varphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^k w_q \varphi(t_q).$$



Optimalité des choix des points/poids de quadrature?

Avec k+1 points de quadrature $\ddagger \mathbb{R}:$ peut-on être exact au delà de \mathbb{P}_k ?

$$(M_e) \qquad \int_0^1 \varphi(t) dt \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

▶ Remarque. pour $p(t) = \prod_{q=0}^k (t-t_q)^2 \in \mathbb{P}_{2k+2}$,

$$\int_0^1 p(t) dt > 0, \quad \sum_{q=0}^k w_q p(t_q) = 0.$$



Optimalité des choix des points/poids de quadrature?

 $\boxed{\mathsf{A}\mathsf{vec}\,\,k+1\;\mathsf{points}\;\mathsf{de}\;\mathsf{quadrature} \overset{.}{\downarrow} \mathbb{R}\;:\;\mathsf{peut}\text{-}\mathsf{on}\;\mathsf{\^{e}tre}\;\mathsf{exact}\;\mathsf{au}\;\mathsf{del} \grave{\mathsf{a}}\;\mathsf{de}\;\mathbb{P}_k\,?}$

$$(M_e) \qquad \int_0^1 \varphi(t) dt \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

▶ Remarque. pour $p(t) = \prod_{q=0}^k (t-t_q)^2 \in \mathbb{P}_{2k+2}$,

$$\int_0^1 p(t) dt > 0, \quad \sum_{q=0}^k w_q p(t_q) = 0.$$

 $\Longrightarrow (M_e)$ ne peut pas être exact \mathbb{P}_{2k+2} !



Optimalité des choix des points/poids de quadrature?

Avec k+1 points de quadrature‡ ${\Bbb R}$: peut-on être exact au delà de ${\Bbb P}_k$?

$$(M_e) \qquad \int_0^1 \varphi(t) dt \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

▶ Remarque. pour $p(t) = \prod_{q=0}^k (t-t_q)^2 \in \mathbb{P}_{2k+2}$,

$$\int_0^1 p(t) \, \mathrm{d}t > 0, \quad \sum_{q=0}^k w_q p(t_q) = 0.$$

- $\Longrightarrow (M_e)$ ne peut pas être exact \mathbb{P}_{2k+2} !
- ▶ Peut-on espérer être exact \mathbb{P}_{2k+1} ?



Optimalité des choix des points/poids de quadrature?

$$(M_e)$$

$$\int_0^1 \varphi(t) dt \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

 $ightharpoonup (M_e)$ exact sur \mathbb{P}_{2k+1} ?



Optimalité des choix des points/poids de quadrature?

$$(M_e)$$

$$\int_0^1 \varphi(t) dt \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

 $\qquad \qquad (M_e) \text{ exact sur } \mathbb{P}_{2k+1} \text{? Soit } p \in \mathbb{P}_{2k+1} \text{ et } \pi(t) = \prod_{q=0}^k (t-t_q) \in \mathbb{P}_{k+1}.$



$$(M_e)$$

$$\int_0^1 \varphi(t) dt \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

- $\qquad \qquad (\textit{M}_e) \text{ exact sur } \mathbb{P}_{2k+1} \text{? Soit } p \in \mathbb{P}_{2k+1} \text{ et } \pi(t) = \prod_{q=0}^{\kappa} (t-t_q) \in \mathbb{P}_{k+1}.$
- ▶ Division euclidienne : $p = \pi u + r$, avec $d^{\circ}u \leq k$ et $d^{\circ}r \leq k$.



$$(M_e)$$

$$\int_0^1 \varphi(t) dt \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

- (M_e) exact sur \mathbb{P}_{2k+1} ? Soit $p \in \mathbb{P}_{2k+1}$ et $\pi(t) = \prod_{i=0}^{n} (t t_q) \in \mathbb{P}_{k+1}$.
- ▶ Division euclidienne : $p = \pi u + r$, avec $d^{\circ}u \leq k$ et $d^{\circ}r \leq k$.

$$\sum_{q=0}^k w_q p(t_q) = \sum_{q=0}^k w_q r(t_q). \qquad \text{car} \pi(t_q) = 0.$$



$$(M_e)$$

$$\int_0^1 \varphi(t) dt \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

- $\qquad \qquad (\textit{M}_e) \text{ exact sur } \mathbb{P}_{2k+1} \text{? Soit } p \in \mathbb{P}_{2k+1} \text{ et } \pi(t) = \prod_{q=0}^{\kappa} (t-t_q) \in \mathbb{P}_{k+1}.$
- ▶ Division euclidienne : $p = \pi u + r$, avec $d^{\circ}u \leq k$ et $d^{\circ}r \leq k$.

$$\sum_{q=0}^k w_q p(t_q) = \sum_{q=0}^k w_q r(t_q). \qquad \text{car} \pi(t_q) = 0.$$

$$\int_0^1 p(t) dt = \int_0^1 \pi(t) u(t) dt + \int_0^1 r(t) dt.$$



Optimalité des choix des points/poids de quadrature?

$$(M_e)$$

$$\int_0^1 \varphi(t) dt \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

- $\qquad \qquad (\textit{M}_e) \text{ exact sur } \mathbb{P}_{2k+1} \text{? Soit } p \in \mathbb{P}_{2k+1} \text{ et } \pi(t) = \prod_{q=0}^{\kappa} (t-t_q) \in \mathbb{P}_{k+1}.$
- ▶ Division euclidienne : $p = \pi u + r$, avec $d^{\circ}u \leq k$ et $d^{\circ}r \leq k$.

$$\sum_{q=0}^k w_q p(t_q) = \sum_{q=0}^k w_q r(t_q). \qquad \text{car} \pi(t_q) = 0.$$

$$\int_0^1 p(t) dt = \int_0^1 \pi(t) u(t) dt + \int_0^1 r(t) dt.$$

▶ Si (M_e) est exact sur \mathbb{P}_k et $\pi \perp \mathbb{P}_k$, alors (M_e) est exact sur \mathbb{P}_{2k+1} !



$$(M_e)$$

$$\int_0^1 \varphi(t) dt \simeq \sum_{q=0}^k w_q \varphi(t_q).$$

- $\qquad \qquad (\textit{M}_e) \text{ exact sur } \mathbb{P}_{2k+1} \text{? Soit } p \in \mathbb{P}_{2k+1} \text{ et } \pi(t) = \prod_{q=0}^{\kappa} (t-t_q) \in \mathbb{P}_{k+1}.$
- ▶ Division euclidienne : $p = \pi u + r$, avec $d^{\circ}u \leq k$ et $d^{\circ}r \leq k$.

$$\sum_{q=0}^k w_q p(t_q) = \sum_{q=0}^k w_q r(t_q). \qquad \text{car} \pi(t_q) = 0.$$

$$\int_0^1 p(t) dt = \int_0^1 \pi(t) u(t) dt + \int_0^1 r(t) dt.$$

- ▶ Si (M_e) est exact sur \mathbb{P}_k et $\pi \perp \mathbb{P}_k$, alors (M_e) est exact sur \mathbb{P}_{2k+1} !
- \implies (t_q) = racines du $(k+1)^e$ polynôme orthogonal pour $L^2(0,1)$.



Méthodes de Gauss

$$(\varphi,\psi)=\int_0^1 \varphi(t)\psi(t)\,\mathrm{d}t.$$

Théorème. Soit (P_0,P_1,\ldots) famille de polynômes orthogonaux (avec $d^{\circ}P_i=i)$ pour le produit scalaire $(\varphi,\psi)=\int_0^1\varphi(t)\psi(t)\,\mathrm{d}t.$ On note $(t_q)_{q=0,\ldots,k}$ les racines de P_{k+1} . On détermine les poids $(w_q)_{q=0,\ldots,k}$ t.q. (M_e) soit exact \mathbb{P}_k , i.e. solution de $\forall i=0,1,\ldots,k, \quad \sum_{q=0}^k w_q t_q^i=\frac{1}{i+1}.$ Alors la formule $(M_e)\qquad \qquad \int_0^1\varphi(t)\,\mathrm{d}t \simeq \sum_{q=0}^k w_q\varphi(t_q).$ est exacte \mathbb{P}_{2k+1} . $(M\acute{e}thode\ de\ Gauss-Legendre).$

$$orall i=0,1,\ldots,k, \quad \sum_{q=0}^k w_q t_q^i = rac{1}{i+1}.$$

$$\int_0^1 arphi(t) \, \mathrm{d}t \simeq \sum_{q=0}^k w_q arphi(t_q).$$

Exemple. Méthode de Gauss-Legendre exacte \mathbb{P}_5 .

$$\int_0^1 \varphi(t)\,\mathrm{d}t \simeq \frac{5}{18}\varphi\left(\frac{1}{2}-\frac{\sqrt{15}}{10}\right) + \frac{4}{9}\varphi\left(\frac{1}{2}\right) + \frac{5}{18}\varphi\left(\frac{1}{2}+\frac{\sqrt{15}}{10}\right).$$

La méthode composée correspondante est en $\mathcal{O}(N^{-6})$.

Pour k grand, le calcul des points n'est pas explicite.





Exemple. Méthode de Gauss-Legendre exacte \mathbb{P}_5 .

$$\int_0^1 \varphi(t)\,\mathrm{d}t \simeq \frac{5}{18}\varphi\left(\frac{1}{2}-\frac{\sqrt{15}}{10}\right) + \frac{4}{9}\varphi\left(\frac{1}{2}\right) + \frac{5}{18}\varphi\left(\frac{1}{2}+\frac{\sqrt{15}}{10}\right).$$

La méthode composée correspondante est en $\mathcal{O}(N^{-6})$.

Pour k grand, le calcul des points n'est pas explicite.



Remarque. Si le *poids* $\omega \ge 0$ est tel que

$$orall p \in \mathbb{P}, \quad \int_0^1 |p(t)| \omega(t) \, \mathrm{d}t < +\infty,$$

alors on peut faire de même pour le produit scalaire $L^2(\omega(t) dt)$.



Culture : la méthode de Monte-Carlo

$$\mathbb{E}\left[f(X)\right] = \frac{1}{b-a} \int_a^b f(x) \, \mathrm{d}x.$$

Rappel. Si $X \hookrightarrow \mathbb{U}([a,b])$ et $f:[a,b] \to \mathbb{R}$, alors $\mathbb{E}\left[f(X)\right] = \frac{1}{b-a} \int_a^b f(x) \, \mathrm{d}x.$ Si (x_0,x_1,\ldots,x_{N-1}) est une réalisation d'un échantillon issu de la loi de X, alors $\frac{b-a}{N} \sum_{i=0}^{N-1} f(x_i) \quad \text{approche} \quad \int_a^b f(x) \, \mathrm{d}x.$

$$\frac{b-a}{N} \sum_{i=0}^{N-1} f(x_i)$$
 approche $\int_a^b f(x) dx$.



Culture : la méthode de Monte-Carlo

$$\mathbb{E}\left[f(X)\right] = \frac{1}{b-a} \int_a^b f(x) \, \mathrm{d}x$$

Rappel. Si $X \hookrightarrow \mathbb{U}([a,b])$ et $f:[a,b] \to \mathbb{R}$, alors $\mathbb{E}\left[f(X)\right] = \frac{1}{b-a} \int_a^b f(x) \, \mathrm{d}x.$ Si (x_0,x_1,\ldots,x_{N-1}) est une réalisation d'un échantillon issu de la loi de X, alors $\frac{b-a}{N} \sum_{i=0}^{N-1} f(x_i) \quad \text{approche} \quad \int_a^b f(x) \, \mathrm{d}x.$

$$\frac{b-a}{N} \sum_{i=0}^{N-1} f(x_i)$$
 approche $\int_a^b f(x) dx$.

Sous matlab:

Méthode en $\mathcal{O}(N^{-\frac{1}{2}})$, surtout utile en grande dimension.





COURS 3

Optimisation numérique



Optimiser : une démarche universelle

► Mécanique :



- ► Mécanique :
 - Équilibre d'un système minimisant une énergie



- ► Mécanique :
 - ► Équilibre d'un système minimisant une énergie
 - Structure qui maximise la résistance



- Mécanique :
 - Équilibre d'un système minimisant une énergie
 - Structure qui maximise la résistance
- ► Transport :



- ► Mécanique :
 - Équilibre d'un système minimisant une énergie
 - ► Structure qui maximise la résistance
- ► Transport :
 - Minimiser la distance, le temps de trajet, le coût,...



- Mécanique :
 - Équilibre d'un système minimisant une énergie
 - Structure qui maximise la résistance
- ► Transport :
 - Minimiser la distance, le temps de trajet, le coût,...
- Economie, gestion :



- Mécanique :
 - Équilibre d'un système minimisant une énergie
 - Structure qui maximise la résistance
- ► Transport :
 - Minimiser la distance, le temps de trajet, le coût,...
- Economie, gestion :
 - Maximiser le profit d'une entreprise



- Mécanique :
 - Équilibre d'un système minimisant une énergie
 - Structure qui maximise la résistance
- ► Transport :
 - Minimiser la distance, le temps de trajet, le coût,...
- Economie, gestion :
 - ► Maximiser le profit d'une entreprise
 - Minimiser les risques d'un placement boursier



- Mécanique :
 - Équilibre d'un système minimisant une énergie
 - Structure qui maximise la résistance
- ► Transport :
 - Minimiser la distance, le temps de trajet, le coût,...
- Economie, gestion :
 - ► Maximiser le profit d'une entreprise
 - Minimiser les risques d'un placement boursier
- Médecine :



- Mécanique :
 - Équilibre d'un système minimisant une énergie
 - Structure qui maximise la résistance
- ► Transport :
 - Minimiser la distance, le temps de trajet, le coût,...
- Economie, gestion :
 - ► Maximiser le profit d'une entreprise
 - Minimiser les risques d'un placement boursier
- Médecine :
 - Optimiser une thérapie sous contraintes de dosage



- Mécanique :
 - Équilibre d'un système minimisant une énergie
 - Structure qui maximise la résistance
- ► Transport :
 - Minimiser la distance, le temps de trajet, le coût,...
- Economie, gestion :
 - ► Maximiser le profit d'une entreprise
 - Minimiser les risques d'un placement boursier
- Médecine :
 - Optimiser une thérapie sous contraintes de dosage
- Imagerie, problèmes inverses :



Optimiser : une démarche universelle

- Mécanique :
 - Équilibre d'un système minimisant une énergie
 - Structure qui maximise la résistance
- ► Transport :
 - Minimiser la distance, le temps de trajet, le coût,...
- Economie, gestion :
 - Maximiser le profit d'une entreprise
 - Minimiser les risques d'un placement boursier
- Médecine :
 - Optimiser une thérapie sous contraintes de dosage
- Imagerie, problèmes inverses :
 - Restaurer une image abimée par minimisation de variation



Optimiser : une démarche universelle

- Mécanique :
 - Équilibre d'un système minimisant une énergie
 - Structure qui maximise la résistance
- ► Transport :
 - Minimiser la distance, le temps de trajet, le coût,...
- Economie, gestion :
 - Maximiser le profit d'une entreprise
 - Minimiser les risques d'un placement boursier
- ► Médecine :
 - Optimiser une thérapie sous contraintes de dosage
- Imagerie, problèmes inverses :
 - Restaurer une image abimée par minimisation de variation
 - Résoudre des équations de forme $\varphi(x) = y$ en minimisant

$$\mathbf{f}(x) = \|\varphi(x) - y\|^2$$



Position d'un câble soumis à la gravité





Quelle est la position exacte, la tension (mécanique), la distance au sol... des lignes HT?

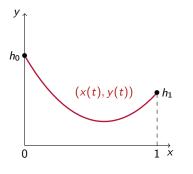


Idem pour les cables d'un téléphérique?

Les cables minimisent leur énergie potentielle



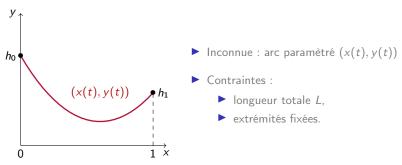
Modélisation : câbles HT



- ▶ Inconnue : arc paramètré (x(t), y(t))
- ► Contraintes :
 - ► longueur totale *L*,
 - extrémités fixées.



Modélisation : câbles HT

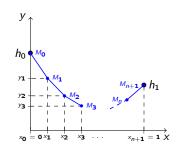


Théorème de l'énergie potentielle.

Le cable est à l'équilibre s'il minimise son énergie potentielle



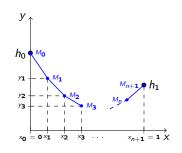
Modélisation du câble HT : discrétisation



Discrétisation de l'espace



Modélisation du câble HT : discrétisation

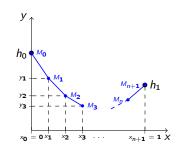


Discrétisation de l'espace

► Câble affine par morceaux



Modélisation du câble HT : discrétisation

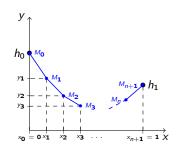


Discrétisation de l'espace

- ► Câble affine par morceaux
- $\blacktriangleright M_i = (x_i, y_i).$



Modélisation du câble HT : discrétisation



Discrétisation de l'espace

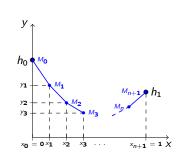
- ► Câble affine par morceaux
- $\qquad M_i = (x_i, y_i).$
- ▶ On note

$$\mathbf{x} = (x_1, \dots, x_n), \ \mathbf{y} = (y_1, \dots, y_n).$$

et $y_0 = h_0, \ y_{n+1} = h_1.$



Modélisation du câble HT : discrétisation



Discrétisation de l'espace

- ► Câble affine par morceaux
- $\qquad M_i = (x_i, y_i).$
- ► On note

$$\mathbf{x} = (x_1, \dots, x_n), \ \mathbf{y} = (y_1, \dots, y_n).$$

et
$$y_0 = h_0, \ y_{n+1} = h_1.$$

Problème discret : $\min \{f(x,y) ; (x,y) \in K\}$

avec

$$f(x,y) = \sum_{i=1}^n y_i$$

$$\qquad \qquad \mathbf{K} = \Big\{ (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{2n} \; ; \; (x_{i+1} - x_i)^2 + (y_{i+1} - y_i)^2 = \frac{L^2}{(n+1)^2}, \; \forall i = 0 \dots n \Big\}.$$



Formalisme général

Problème d'optimisation.

 $\min_{u \in \mathcal{K}} f(u)$

avec $K\subset\mathbb{R}^N$ et $\mathbf{f}:K o\mathbb{R}$.



Formalisme général

Problème d'optimisation.

$$\min_{\mathbf{u} \in K} \mathbf{f}(\mathbf{u})$$

avec $K \subset \mathbb{R}^N$ et $\mathbf{f}: K \to \mathbb{R}$.

Vocabulaire:

▶ f : fonction objectif ou énergie.



Formalisme général

Problème d'optimisation.

$$\min_{u \in \mathcal{K}} f(u)$$

avec $K \subset \mathbb{R}^N$ et $\mathbf{f}: K \to \mathbb{R}$.

Vocabulaire:

- ▶ f : fonction objectif ou énergie.
- $ightharpoonup \mathbb{R}^N$: espace des paramètres.



Formalisme général

Problème d'optimisation.



avec $K \subset \mathbb{R}^N$ et $\mathbf{f}: K \to \mathbb{R}$.

Vocabulaire:

- ▶ f : fonction objectif ou énergie.
- $ightharpoonup \mathbb{R}^N$: espace des paramètres.
- K : ensemble admissible.



Formalisme général

Problème d'optimisation.

$$\min_{\mathbf{u}\in K}\mathbf{f}(\mathbf{u})$$

avec $K \subset \mathbb{R}^N$ et $\mathbf{f}: K \to \mathbb{R}$.

Vocabulaire:

- f : fonction objectif ou énergie.
- $ightharpoonup \mathbb{R}^N$: espace des paramètres.
- K : ensemble admissible.
- ▶ Solution : point $\mathbf{u}^* \in K$ qui minimise \mathbf{f} sur K; on écrit

$$f(u^*) = \min_{u \in \mathcal{K}} f(u) \quad \text{ et aussi } \quad u^* = \underset{u \in \mathcal{K}}{\mathsf{argmin}} \ f(u)$$



Formalisme général

Problème d'optimisation.

$$\min_{u \in \mathcal{K}} f(u)$$

avec $K \subset \mathbb{R}^N$ et $\mathbf{f}: K \to \mathbb{R}$.

Vocabulaire:

▶ f : fonction objectif ou énergie.

 $ightharpoonup \mathbb{R}^N$: espace des paramètres.

K : ensemble admissible.

▶ Solution : point $\mathbf{u}^* \in K$ qui minimise \mathbf{f} sur K; on écrit

$$f(u^*) = \min_{u \in \mathcal{K}} f(u) \quad \text{ et aussi } \quad u^* = \underset{u \in \mathcal{K}}{\text{argmin }} f(u)$$

 $ightharpoonup f(u^*)$: minimum de f sur K.



Formalisme général

Problème d'optimisation.

$$\min_{u \in \mathcal{K}} f(u)$$

avec $K \subset \mathbb{R}^N$ et $\mathbf{f}: K \to \mathbb{R}$.

Vocabulaire:

f : fonction objectif ou énergie.

 $ightharpoonup \mathbb{R}^N$: espace des paramètres.

K : ensemble admissible.

▶ Solution : point $\mathbf{u}^* \in K$ qui minimise \mathbf{f} sur K; on écrit

$$f(u^*) = \min_{u \in \mathcal{K}} f(u) \quad \text{ et aussi } \quad u^* = \underset{u \in \mathcal{K}}{\text{argmin }} f(u)$$

 $ightharpoonup f(u^*)$: minimum de f sur K.

Attention: Il peut y avoir zéro, une ou plusieurs solutions au problème.



Optimisation libre dans \mathbb{R}



Quelques définitions

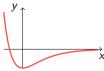
Définition. Une fonction $f: \mathbb{R} \mapsto \mathbb{R}$ est unimodale s'il existe $x^* \in \mathbb{R}$ tel que f soit strictement décroissante sur $]-\infty, x^*[$ et strictement croissante sur $]x^*, +\infty[$.



Quelques définitions

Définition. Une fonction $f: \mathbb{R} \mapsto \mathbb{R}$ est unimodale s'il existe $x^* \in \mathbb{R}$ tel que f soit strictement décroissante sur $]-\infty, x^*[$ et strictement croissante sur $]x^*, +\infty[$.

Exemple:
$$x \mapsto e^{-2x} - 2e^{-x}$$



Quelques définitions

Définition. Une fonction $f: \mathbb{R} \mapsto \mathbb{R}$ est unimodale s'il existe $x^* \in \mathbb{R}$ tel que f soit strictement décroissante sur $]-\infty, x^*[$ et strictement croissante sur $]x^*, +\infty[$.

Définition. Une fonction $f: \mathbb{R} \mapsto \mathbb{R}$ est coercive si $\lim_{|x| \to +\infty} f(x) = +\infty$.

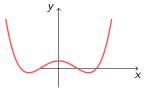


Quelques définitions

Définition. Une fonction $f: \mathbb{R} \mapsto \mathbb{R}$ est unimodale s'il existe $x^* \in \mathbb{R}$ tel que f soit strictement décroissante sur $]-\infty, x^*[$ et strictement croissante sur $]x^*, +\infty[$.

Définition. Une fonction $f : \mathbb{R} \to \mathbb{R}$ est coercive si $\lim_{|x| \to +\infty} f(x) = +\infty$.

Exemple: $x \mapsto x^4 - 5x^3 + 4$





Quelques définitions

Définition. Une fonction $f: \mathbb{R} \mapsto \mathbb{R}$ est unimodale s'il existe $x^* \in \mathbb{R}$ tel que f soit strictement décroissante sur $]-\infty, x^*[$ et strictement croissante sur $]x^*, +\infty[$.

Définition. Une fonction $f: \mathbb{R} \mapsto \mathbb{R}$ est coercive si $\lim_{|x| \to +\infty} f(x) = +\infty$.

Définition. Une fonction $f: \mathbb{R} \mapsto \mathbb{R}$ de classe \mathscr{C}^2 est fortement convexe ssi $\exists \alpha > 0$ tel que pour tout $x \in \mathbb{R}$, $f''(x) \geqslant \alpha$. (On dit aussi α -convexe).



Quelques définitions

Définition. Une fonction $f: \mathbb{R} \mapsto \mathbb{R}$ est unimodale s'il existe $x^* \in \mathbb{R}$ tel que f soit strictement décroissante sur $]-\infty, x^*[$ et strictement croissante sur $]x^*, +\infty[$.

Définition. Une fonction $f : \mathbb{R} \to \mathbb{R}$ est coercive si $\lim_{|x| \to +\infty} f(x) = +\infty$.

Définition. Une fonction $f: \mathbb{R} \mapsto \mathbb{R}$ de classe \mathscr{C}^2 est fortement convexe ssi $\exists \alpha > 0$ tel que pour tout $x \in \mathbb{R}$, $f''(x) \geqslant \alpha$. (On dit aussi α -convexe).

Exemple: $x \mapsto e^x$ strictement convexe, mais pas fortement convexe.





Quelques définitions

Définition. Une fonction $f: \mathbb{R} \to \mathbb{R}$ est unimodale s'il existe $x^* \in \mathbb{R}$ tel que f soit strictement décroissante sur $]-\infty, x^*[$ et strictement croissante sur $]x^*, +\infty[$.

I Définition. Une fonction $f : \mathbb{R} \mapsto \mathbb{R}$ est coercive si $\lim_{|x| \to +\infty} f(x) = +\infty$.

Définition. Une fonction $f: \mathbb{R} \mapsto \mathbb{R}$ de classe \mathscr{C}^2 est fortement convexe ssi $\exists \alpha > 0$ tel que pour tout $x \in \mathbb{R}$, $f''(x) \geqslant \alpha$. (On dit aussi α -convexe).

Proposition.

- ▶ fortement convexe ⇒ coercive
 ▶ fortement convexe et coercive ⇒ unimodale



Méthode de dichotomie 1

But : résoudre F(x) = 0

▶ $F : \mathbb{R} \to \mathbb{R}$ croissante et continue, s'annule en $x^* \in [a_0, b_0]$.



31

Méthode de dichotomie 1

But : résoudre F(x) = 0

- ▶ $F : \mathbb{R} \to \mathbb{R}$ croissante et continue, s'annule en $x^* \in [a_0, b_0]$.
- On construit (a_n) , (b_n) par l'algorithme

On pose
$$x_n = \frac{a_n + b_n}{2}$$
.
Si $F(x_n) \le 0$,

or
$$F(x_n) \leq 0$$
,

$$a_{n+1}=x_n, \quad \text{et} \quad b_{n+1}=b_n.$$

▶ Si
$$F(x_n) > 0$$
,

► Si
$$F(x_n) > 0$$
,
 $a_{n+1} = a_n$, et $b_{n+1} = x_n$.

Méthode de dichotomie 1

But : résoudre F(x) = 0

- ▶ $F : \mathbb{R} \to \mathbb{R}$ croissante et continue, s'annule en $x^* \in [a_0, b_0]$.
- ▶ On construit (a_n) , (b_n) par l'algorithme

On pose
$$x_n = \frac{a_n + b_n}{2}$$
.
Si $F(x_n) \le 0$,
 $a_{n+1} = x_n$, et $b_{n+1} = b_n$.
Si $F(x_n) > 0$,
 $a_{n+1} = a_n$, et $b_{n+1} = x_n$.

Proposition. La méthode converge à vitesse géométrique :

$$\forall n \in \mathbb{N}, \quad |x_n - x^*| \le \frac{b_0 - a_0}{2^{n+1}}.$$



Méthode de dichotomie 2

But : résoudre $\min f(x)$

▶ $f : \mathbb{R} \to \mathbb{R}$ unimodale et dérivable, minimale en $x^* \in [a_0, b_0]$.



Méthode de dichotomie 2

But : résoudre min f(x)

- ▶ $f : \mathbb{R} \to \mathbb{R}$ unimodale et dérivable, minimale en $x^* \in [a_0, b_0]$.
- ▶ On construit (a_n) , (b_n) par l'algorithme

On pose
$$x_n = \frac{a_n + b_n}{2}$$
.

Si
$$f'(x_n) \le 0$$
,
 $a_{n+1} = x_n$, et $b_{n+1} = b_n$.

► Si
$$f'(x_n) > 0$$
, $a_{n+1} = a_n$, et $b_{n+1} = x_n$.



Méthode de dichotomie 2

But : résoudre min f(x)

- ▶ $f : \mathbb{R} \to \mathbb{R}$ unimodale et dérivable, minimale en $x^* \in [a_0, b_0]$.
- ▶ On construit (a_n) , (b_n) par l'algorithme

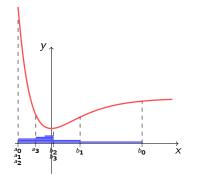
On pose
$$x_n = \frac{a_n + b_n}{2}$$
.
Si $f'(x_n) \le 0$,
 $a_{n+1} = x_n$, et $b_{n+1} = b_n$.
Si $f'(x_n) > 0$,
 $a_{n+1} = a_n$, et $b_{n+1} = x_n$.

Proposition. La méthode converge à vitesse géométrique :

$$\forall n \in \mathbb{N}, \quad |x_n - x^*| \le \frac{b_0 - a_0}{2^{n+1}}.$$



Méthode de dichotomie 2



Avantages:

- ► Très simple à mettre en œuvre.
- Convergence rapide.

Inconvénient:

► Évaluation de f'.

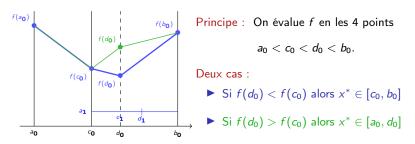
Coût : À chaque itération, on évalue seulement une fois la dérivée de f.





Méthode du nombre d'or

Soit $f : \mathbb{R} \to \mathbb{R}$ unimodale, minimale en $x^* \in [a_0, b_0]$.



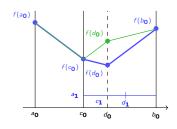
- La taille de l'intervalle diminue à chaque étape.
- ► Coût : À chaque itération, une seule nouvelle évaluation de f.





Méthode du nombre d'or

Quel rapport avec le nombre d'or?



► Par symétrie on impose

$$b_0 - c_0 = d_0 - a_0.$$

► Pour conserver le ratio des intervalles

$$\gamma = \frac{b_0 - a_0}{b_0 - c_0} = \frac{b_1 - a_1}{b_1 - c_1} = \frac{b_0 - c_0}{b_0 - d_0}$$

$$\frac{1}{\gamma} = \frac{b_0 - d_0}{b_0 - c_0} = \gamma - \frac{a_0 - d_0}{b_0 - c_0} = \gamma - 1$$

Donc $\gamma^2 - \gamma - 1 = 0$ et ainsi

$$\gamma = \frac{1 + \sqrt{5}}{2}$$



Méthode du nombre d'or

- ▶ $f : \mathbb{R} \to \mathbb{R}$ unimodale, minimale en $x^* \in [a_0, b_0]$.
- ▶ On construit (a_n) , (b_n) par l'algorithme

$$\begin{array}{|c|c|c|} \blacktriangleright & \text{Si } f(c_n) \leq f(d_n), \\ \hline & a_{n+1} &= a_n \\ & b_{n+1} &= d_n \\ & c_{n+1} &= b_{n+1} - \frac{b_{n+1} - a_{n+1}}{\gamma} \\ & d_{n+1} &= c_n \\ \hline \end{array} \quad \begin{array}{|c|c|c|} \hline & \blacktriangleright & \text{Si } f(c_n) > f(d_n), \\ \hline & a_{n+1} &= c_n \\ & b_{n+1} &= b_n \\ & c_{n+1} &= d_n \\ & d_{n+1} &= a_{n+1} + \frac{b_{n+1} - a_{n+1}}{\gamma} \\ \hline \end{array}$$

Proposition. La méthode converge à vitesse géométrique :

$$\forall n \in \mathbb{N}, \quad \left| \frac{a_n + b_n}{2} - x^* \right| \leq \frac{1}{2} (b_0 - a_0) \gamma^{-n}.$$



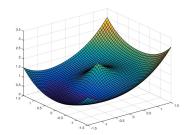
Optimisation libre dans \mathbb{R}^N



Méthodes de descentes : principe

Exemple:

$$f(x) = \left(\sqrt{x_1^2 + x_2^2} - 1\right)^2 + \frac{x_1}{2} + 1.$$



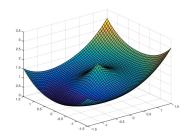


Méthodes de descentes : principe

Exemple:

$$\mathbf{f(x)} = \left(\sqrt{x_1^2 + x_2^2} - 1\right)^2 + \frac{x_1}{2} + 1.$$

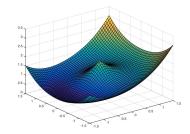
À partir d'un point initial $\mathbf{x}^0 \in \mathbb{R}^N$,



Méthodes de descentes : principe

Exemple:

$$f(\textbf{x}) = \left(\sqrt{x_1^2 + x_2^2} - 1\right)^2 + \frac{x_1}{2} + 1.$$



À partir d'un point initial $x^0 \in \mathbb{R}^N$,

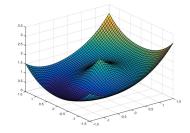
▶ on choisit une direction de descente $\mathbf{d}^0 \in \mathbb{R}^N$ et un pas $\rho_0 > 0$,



Méthodes de descentes : principe

Exemple:

$$f(x) = \left(\sqrt{x_1^2 + x_2^2} - 1\right)^2 + \frac{x_1}{2} + 1.$$



À partir d'un point initial $x^0 \in \mathbb{R}^N$,

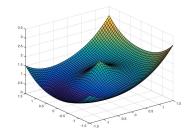
- on choisit une direction de descente $\mathbf{d}^0 \in \mathbb{R}^N$ et un pas $\rho_0 > 0$,
- on obtient un nouveau point $x^1 = x^0 + \rho_0 d^0$,



Méthodes de descentes : principe

Exemple:

$$f(x) = \left(\sqrt{x_1^2 + x_2^2} - 1\right)^2 + \frac{x_1}{2} + 1.$$



À partir d'un point initial $\mathbf{x}^0 \in \mathbb{R}^N$,

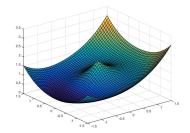
- on choisit une direction de descente $\mathbf{d}^0 \in \mathbb{R}^N$ et un pas $\rho_0 > 0$,
- on obtient un nouveau point $x^1 = x^0 + \rho_0 d^0$,
- on itère le processus : $\mathbf{x}^{n+1} = \mathbf{x}^n + \rho_n \mathbf{d}^n$



Méthodes de descentes : principe

Exemple:

$$f(x) = \left(\sqrt{x_1^2 + x_2^2} - 1\right)^2 + \frac{x_1}{2} + 1.$$



À partir d'un point initial $\mathbf{x}^0 \in \mathbb{R}^N$,

- ▶ on choisit une direction de descente $\mathbf{d}^0 \in \mathbb{R}^N$ et un pas $\rho_0 > 0$,
- on obtient un nouveau point $x^1 = x^0 + \rho_0 d^0$,
- on itère le processus : $x^{n+1} = x^n + \rho_n d^n$

But. faire en sorte que $f(x^{n+1}) \leq f(x^n)$.





Optimisation numérique Gradient à pas fixe

Comment choisir la direction \mathbf{d}^n ?

▶ On veut $(f(x^n))_n$ décroissante.



Gradient à pas fixe

Comment choisir la direction \mathbf{d}^n ?

- ▶ On veut $(f(x^n))_n$ décroissante.
- ▶ On choisit \mathbf{d}^n telle que $\varphi: t \mapsto \mathbf{f}(\mathbf{x}^n + t\mathbf{d}^n)$ soit décroissante au voisinage de 0.



Gradient à pas fixe

Comment choisir la direction \mathbf{d}^n ?

- ▶ On veut $(f(x^n))_n$ décroissante.
- ▶ On choisit \mathbf{d}^n telle que $\varphi: t \mapsto \mathbf{f}(\mathbf{x}^n + t\mathbf{d}^n)$ soit décroissante au voisinage de 0.



Gradient à pas fixe

Comment choisir la direction \mathbf{d}^n ?

- ▶ On veut $(f(x^n))_n$ décroissante.
- ▶ On choisit \mathbf{d}^n telle que $\varphi: t \mapsto \mathbf{f}(\mathbf{x}^n + t\mathbf{d}^n)$ soit décroissante au voisinage de 0.
- ► Un choix simple est

$$\mathbf{d}^n = -\nabla \mathbf{f}(\mathbf{x}^n).$$



Gradient à pas fixe

Comment choisir la direction \mathbf{d}^n ?

- ▶ On veut $(f(x^n))_n$ décroissante.
- ▶ On choisit \mathbf{d}^n telle que $\varphi: t \mapsto \mathbf{f}(\mathbf{x}^n + t\mathbf{d}^n)$ soit décroissante au voisinage de 0.
- ► Un choix simple est

$$\mathbf{d}^n = -\nabla \mathbf{f}(\mathbf{x}^n).$$

Méthode du gradient à pas fixe. Elle est définie par

$$\begin{cases} \mathbf{x}^0 \in \mathbb{R}^N, \\ \mathbf{x}^{n+1} = \mathbf{x}^n - \rho \nabla \mathbf{f}(\mathbf{x}^n). \end{cases}$$

avec a > 0



Gradient à pas fixe : convergence

Théorème. Soit
$$\mathbf{f} \in \mathscr{C}^1(\mathbb{R}^N, \mathbb{R})$$
 strictement convexe et coercive telle que $(*)$ $\exists M>0, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^N, \quad \|\nabla \mathbf{f}(\mathbf{x}) - \nabla \mathbf{f}(\mathbf{y})\|_2 \leq M \|\mathbf{x} - \mathbf{y}\|_2.$ Si $0<\rho<\frac{2}{M}$ alors la méthode du gradient à pas fixe converge. La limite est l'unique point de minimum global de \mathbf{f} .



Gradient à pas fixe : convergence

Théorème. Soit
$$\mathbf{f} \in \mathscr{C}^1(\mathbb{R}^N, \mathbb{R})$$
 strictement convexe et coercive telle que $(*)$ $\exists M>0, \quad \forall \mathbf{x},\mathbf{y} \in \mathbb{R}^N, \quad \|\nabla \mathbf{f}(\mathbf{x})-\nabla \mathbf{f}(\mathbf{y})\|_2 \leq M \|\mathbf{x}-\mathbf{y}\|_2.$ Si $0<\rho<\frac{2}{M}$ alors la méthode du gradient à pas fixe converge. La limite est l'unique point de minimum global de \mathbf{f} .

Remarques:

- ▶ Condition (*) : $\mathbf{x} \mapsto \nabla \mathbf{f}(\mathbf{x})$ est globalement *M*-lipschitzienne sur \mathbb{R}^N .
- ► Si on ajoute l'hypothèse (f est fortement convexe)

$$\exists \alpha > 0, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^N, \quad (H\mathbf{f}(\mathbf{x})\mathbf{y}|\mathbf{y}) \ge \alpha \|\mathbf{y}\|^2.$$

et $\rho < \frac{2\alpha}{\mathit{M}^2}$ alors la convergence est géométrique.



Rappel sur la convexité

Pour des fonctions de classe \mathscr{C}^2

$f:\mathbb{R} o\mathbb{R}$	$\mathbf{f}:\mathbb{R}^N o\mathbb{R}$	
$\mathbf{f}'' \geq 0$	$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^N, \ (H\mathbf{f}(\mathbf{x})\mathbf{y} \mathbf{y}) \geq 0$	convexe
f'' > 0	$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^N, \ y \neq 0, \ (Hf(\mathbf{x})\mathbf{y} \mathbf{y}) > 0$	⇒ strictement convexe
$\mathbf{f}'' \geq \alpha$	$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^N, \ (Hf(\mathbf{x})\mathbf{y} \mathbf{y}) \ge \alpha \ \mathbf{y}\ _2^2$	fortement convexe ($lpha > 0$)



Optimisation numérique Gradient à pas optimal

Peut-on faire mieux?

On peut chercher le meilleur pas possible à chaque étape :

$$\rho_n$$
 minimise $\rho \mapsto f(x^n - \rho \nabla f(x^n))$.

Vocabulaire : cette étape est dite recherche linéaire.

Méthode: type dichotomie.



42

Optimisation numérique Gradient à pas optimal

Peut-on faire mieux?

On peut chercher le meilleur pas possible à chaque étape :

$$\rho_n$$
 minimise $\rho \mapsto f(x^n - \rho \nabla f(x^n))$.

Vocabulaire : cette étape est dite recherche linéaire.

Méthode: type dichotomie.

Méthode du gradient à pas optimal. Elle est définie par

$$\begin{cases} \mathbf{x}^0 \in \mathbb{R}^N, \\ \rho_n = \arg\min_{\rho \in \mathbb{R}} \mathbf{f}(\mathbf{x}^n - \rho \nabla \mathbf{f}(\mathbf{x}^n)) \\ \mathbf{x}^{n+1} = \mathbf{x}^n - \rho_n \nabla \mathbf{f}(\mathbf{x}^n). \end{cases}$$



Optimisation numérique Gradient à pas optimal

Peut-on faire mieux?

On peut chercher le meilleur pas possible à chaque étape :

$$\rho_n$$
 minimise $\rho \mapsto f(x^n - \rho \nabla f(x^n))$.

Vocabulaire : cette étape est dite recherche linéaire.

Méthode: type dichotomie.

Méthode du gradient à pas optimal. Elle est définie par

$$\left\{ \begin{array}{l} \mathbf{x}^0 \in \mathbb{R}^N, \\ \rho_n = \mathop{\arg\min}_{\rho \in \mathbb{R}} \mathbf{f}(\mathbf{x}^n - \rho \nabla \mathbf{f}(\mathbf{x}^n)) \\ \mathbf{x}^{n+1} = \mathbf{x}^n - \rho_n \nabla \mathbf{f}(\mathbf{x}^n). \end{array} \right.$$

Remarque. même type de résultats de convergence que pour le pas fixe.



Gradient à pas optimal pour les moindres carrés

Exemple : Minimisation de $f(x) = ||Ax - b||_2^2$.

ightharpoonup On parle de résolution de Ax = b au sens des moindres carrés.



Gradient à pas optimal pour les moindres carrés

Exemple : Minimisation de $f(x) = ||Ax - b||_2^2$.

- ▶ On parle de résolution de Ax = b au sens des moindres carrés.
- Éventuellement, A est rectangulaire (injective).



43

Gradient à pas optimal pour les moindres carrés

Exemple: Minimisation de $f(x) = ||Ax - b||_2^2$.

- ▶ On parle de résolution de Ax = b au sens des moindres carrés.
- Éventuellement, A est rectangulaire (injective).
- Lien avec l'approximation polynomiale.



Gradient à pas optimal pour les moindres carrés

Exemple: Minimisation de $f(x) = ||Ax - b||_2^2$.

- ▶ On parle de résolution de Ax = b au sens des moindres carrés.
- Éventuellement, A est rectangulaire (injective).
- Lien avec l'approximation polynomiale.



Gradient à pas optimal pour les moindres carrés

Exemple: Minimisation de $f(x) = ||Ax - b||_2^2$.

- On parle de résolution de Ax = b au sens des moindres carrés.
- Éventuellement, A est rectangulaire (injective).
- Lien avec l'approximation polynomiale.
- ► Rappel : $\nabla \mathbf{f}(\mathbf{x}) = 2A^T(A\mathbf{x} \mathbf{b})$

$$\min_{\rho>0} \mathbf{f}(\mathbf{x} - \rho \nabla \mathbf{f}(\mathbf{x})).$$



Gradient à pas optimal pour les moindres carrés

Exemple: Minimisation de $f(x) = ||Ax - b||_2^2$.

- ▶ On parle de résolution de Ax = b au sens des moindres carrés.
- Éventuellement, A est rectangulaire (injective).
- Lien avec l'approximation polynomiale.
- $Rappel : \nabla f(x) = 2A^T (Ax b)$

$$\min_{\rho>0} \mathbf{f}(\mathbf{x} - \rho \nabla \mathbf{f}(\mathbf{x})).$$

$$f(x - \rho \nabla f(x)) = \|Ax - b - \rho A \nabla f(x)\|_{2}^{2}$$



Gradient à pas optimal pour les moindres carrés

Exemple: Minimisation de $f(x) = ||Ax - b||_2^2$.

- ▶ On parle de résolution de Ax = b au sens des moindres carrés.
- Éventuellement, A est rectangulaire (injective).
- Lien avec l'approximation polynomiale.
- ► Rappel : $\nabla \mathbf{f}(\mathbf{x}) = 2\mathbf{A}^T(\mathbf{A}\mathbf{x} \mathbf{b})$

$$\min_{\rho>0} \mathbf{f}(\mathbf{x} - \rho \nabla \mathbf{f}(\mathbf{x})).$$

$$f(x - \rho \nabla f(x)) = \left\| Ax - \mathbf{b} \right\|_{2}^{2} - 2\rho (Ax - \mathbf{b}|A\nabla f(x)) + \rho^{2} \left\| A\nabla f(x) \right\|_{2}^{2}$$



Gradient à pas optimal pour les moindres carrés

Exemple: Minimisation de $f(x) = ||Ax - b||_2^2$.

- ▶ On parle de résolution de Ax = b au sens des moindres carrés.
- Éventuellement, A est rectangulaire (injective).
- Lien avec l'approximation polynomiale.
- ► Rappel : $\nabla \mathbf{f}(\mathbf{x}) = 2\mathbf{A}^T(\mathbf{A}\mathbf{x} \mathbf{b})$

$$\min_{\rho>0} \mathbf{f}(\mathbf{x} - \rho \nabla \mathbf{f}(\mathbf{x})).$$

$$f(x - \rho \nabla f(x)) = \|Ax - b\|_2^2 - 2\rho(A^T Ax - A^T b|\nabla f(x)) + \rho^2 \|A\nabla f(x)\|_2^2$$



Gradient à pas optimal pour les moindres carrés

Exemple: Minimisation de $f(x) = ||Ax - b||_2^2$.

- ▶ On parle de résolution de Ax = b au sens des moindres carrés.
- Éventuellement, A est rectangulaire (injective).
- Lien avec l'approximation polynomiale.
- $Rappel : \nabla f(x) = 2A^T (Ax b)$

$$\min_{\rho>0} \mathbf{f}(\mathbf{x} - \rho \nabla \mathbf{f}(\mathbf{x})).$$

$$f(\mathbf{x} - \rho \nabla f(\mathbf{x})) = \rho^{2} \left\| A \nabla f(\mathbf{x}) \right\|_{2}^{2} - \rho \left\| \nabla f(\mathbf{x}) \right\|_{2}^{2} + \left\| A\mathbf{x} - \mathbf{b} \right\|_{2}^{2}$$



Gradient à pas optimal pour les moindres carrés

Exemple: Minimisation de $f(x) = ||Ax - b||_2^2$.

- On parle de résolution de Ax = b au sens des moindres carrés.
- Éventuellement, A est rectangulaire (injective).
- Lien avec l'approximation polynomiale.
- ► Rappel : $\nabla \mathbf{f}(\mathbf{x}) = 2\mathbf{A}^T(\mathbf{A}\mathbf{x} \mathbf{b})$

Pour trouver le pas optimal, il faut résoudre

$$\min_{\rho>0} \mathbf{f}(\mathbf{x} - \rho \nabla \mathbf{f}(\mathbf{x})).$$

$$f(x - \rho \nabla f(x)) = \rho^{2} \left\| A \nabla f(x) \right\|_{2}^{2} - \rho \left\| \nabla f(x) \right\|_{2}^{2} + \left\| Ax - \mathbf{b} \right\|_{2}^{2}$$

C'est un polynôme de degré 2 en ρ !



Gradient à pas optimal pour les moindres carrés

Exemple: Minimisation de $f(x) = ||Ax - b||_2^2$.

- On parle de résolution de Ax = b au sens des moindres carrés.
- Éventuellement, A est rectangulaire (injective).
- Lien avec l'approximation polynomiale.
- ► Rappel : $\nabla \mathbf{f}(\mathbf{x}) = 2A^T(A\mathbf{x} \mathbf{b})$

Pour trouver le pas optimal, il faut résoudre

$$\min_{\rho>0} \mathbf{f}(\mathbf{x} - \rho \nabla \mathbf{f}(\mathbf{x})).$$

$$f(x - \rho \nabla f(x)) = \rho^{2} \|A \nabla f(x)\|_{2}^{2} - \rho \|\nabla f(x)\|_{2}^{2} + \|Ax - \mathbf{b}\|_{2}^{2}$$

C'est un polynôme de degré 2 en ρ !

$$\rho^* = \frac{\|\nabla f(x)\|_2^2}{2\|A\nabla f(x)\|_2^2}.$$



Méthode de Newton

▶ Résolution d'une équation F(x) = 0 avec $F \in \mathscr{C}^1(\mathbb{R}^N, \mathbb{R}^N)$.

Méthode de Newton pour les zeros de F

$$\begin{cases} x^0 \in \mathbb{R}^N, \\ x^{n+1} = x^n - JF(x^n)^{-1}F(x^n). \end{cases}$$



Méthode de Newton

▶ Résolution d'une équation $\mathbf{F}(\mathbf{x}) = \mathbf{0}$ avec $\mathbf{F} \in \mathscr{C}^1(\mathbb{R}^N, \mathbb{R}^N)$.

Méthode de Newton pour les zeros de F

$$\left\{ \begin{array}{l} \mathbf{x}^0 \in \mathbb{R}^N, \\[1mm] \mathbf{x}^{n+1} = \mathbf{x}^n - J\mathbf{F}(\mathbf{x}^n)^{-1}\mathbf{F}(\mathbf{x}^n). \end{array} \right.$$

Faire un dessin 1D



Méthode de Newton

▶ Résolution d'une équation F(x) = 0 avec $F \in \mathscr{C}^1(\mathbb{R}^N, \mathbb{R}^N)$.

Méthode de Newton pour les zeros de F

$$\begin{cases} x^0 \in \mathbb{R}^N, \\ x^{n+1} = x^n - JF(x^n)^{-1}F(x^n). \end{cases}$$

▶ Minimisation de $\mathbf{f} \in \mathscr{C}^2(\mathbb{R}^N, \mathbb{R})$. On résout $\nabla \mathbf{f}(\mathbf{x}) = \mathbf{0}$.

Méthode de Newton pour minimiser f

$$\begin{cases} x^0 \in \mathbb{R}^N, \\ x^{n+1} = x^n - Hf(x^n)^{-1} \nabla f(x^n). \end{cases}$$



Méthode de Newton

▶ Résolution d'une équation F(x) = 0 avec $F \in \mathscr{C}^1(\mathbb{R}^N, \mathbb{R}^N)$.

Méthode de Newton pour les zeros de F

$$\begin{cases} \mathbf{x}^0 \in \mathbb{R}^N, \\ \mathbf{x}^{n+1} = \mathbf{x}^n - J\mathbf{F}(\mathbf{x}^n)^{-1}\mathbf{F}(\mathbf{x}^n). \end{cases}$$

▶ Minimisation de $\mathbf{f} \in \mathscr{C}^2(\mathbb{R}^N, \mathbb{R})$. On résout $\nabla \mathbf{f}(\mathbf{x}) = \mathbf{0}$.

Méthode de Newton pour minimiser f

$$\left\{ \begin{array}{l} x^0 \in \mathbb{R}^N, \\ \\ x^{n+1} = x^n - Hf(x^n)^{-1} \nabla f(x^n). \end{array} \right.$$

Remarques:

- La méthode est locale, mais rapide,
- ► Ne distingue pas min et max.



Optimisation sous contraintes



Exemple simple

Rappel : Un problème d'optimisation sous contraintes s'écrit

 $\min_{\mathbf{u} \in K} \mathbf{f}(\mathbf{u}), \quad \text{avec} \quad K \subsetneq \mathbb{R}^N.$



Exemple simple

Rappel : Un problème d'optimisation sous contraintes s'écrit

$$\min_{\mathbf{u} \in K} \mathbf{f}(\mathbf{u}), \quad \text{avec} \quad K \subsetneq \mathbb{R}^N.$$

En général, les algorithmes précédents vont converger en dehors de K!



Exemple simple

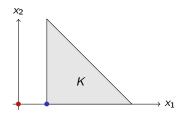
Rappel: Un problème d'optimisation sous contraintes s'écrit

$$\min_{\mathbf{u} \in \mathcal{K}} \mathbf{f}(\mathbf{u}), \quad \text{avec} \quad \mathcal{K} \subsetneq \mathbb{R}^{N}.$$

En général, les algorithmes précédents vont converger en dehors de K!

Exemple:
$$f(x) = x_1^2 + x_2^2$$
.

$$K = \left\{ x \in \mathbb{R}^2, \ x_1 \geq \frac{1}{2}, \ x_2 \geq 0, \ x_1 + x_2 \leq 1 \right\}$$







Projection sur un convexe

Soit $K \subset \mathbb{R}^N$ convexe fermé et non vide.

Définition. Pour tout $\mathbf{x} \in \mathbb{R}^N$, le projeté $\pi_K(\mathbf{x})$ est l'unique point de K qui minimise sa distance à \mathbf{x} . $\|\pi_K(\mathbf{x}) - \mathbf{x}\| = \min \Big\{ \|\mathbf{y} - \mathbf{x}\| \; ; \; \mathbf{y} \in K \Big\}.$

$$\|\pi_{\mathcal{K}}(\mathbf{x}) - \mathbf{x}\| = \min \left\{ \|\mathbf{y} - \mathbf{x}\| \; ; \; \mathbf{y} \in \mathcal{K} \right\}$$



Projection sur un convexe

Soit $K \subset \mathbb{R}^N$ convexe fermé et non vide.

Définition. Pour tout $\mathbf{x} \in \mathbb{R}^N$, le projeté $\pi_K(\mathbf{x})$ est l'unique point de K qui minimise sa distance à \mathbf{x} . $\|\pi_K(\mathbf{x}) - \mathbf{x}\| = \min \Big\{ \|\mathbf{y} - \mathbf{x}\| \; ; \; \mathbf{y} \in K \Big\}.$

$$\|\pi_{\mathcal{K}}(\mathbf{x}) - \mathbf{x}\| = \min \left\{ \|\mathbf{y} - \mathbf{x}\| \; ; \; \mathbf{y} \in \mathcal{K} \right\}$$

Remarques:

- ightharpoonup si $x \in K$ alors $\pi_K(x) = x$,
- ▶ si $\mathbf{x} \notin K$ alors $\pi_K(\mathbf{x}) \in \partial K$,
- \blacktriangleright π_{κ} est contractante donc lipschitzienne donc continue.

Méthode du gradient projeté. Elle est définie par

$$\left\{ \begin{array}{l} x^0 \in \mathbb{R}^N, \\ \\ \mathbf{x}^{n+1} = \pi_K \big(\mathbf{x}^n - \rho \nabla \mathbf{f}(\mathbf{x}^n) \big). \end{array} \right.$$
 avec $\rho > 0.$



Méthode du gradient projeté. Elle est définie par

$$\left\{ \begin{array}{l} \mathsf{x}^0 \in \mathbb{R}^{\it N}, \\ \mathsf{x}^{\it n+1} = \pi_{\it K} \big(\mathsf{x}^{\it n} - \rho \nabla \mathsf{f}(\mathsf{x}^{\it n})\big). \end{array} \right.$$
 avec $\rho > 0.$

$$\exists M > 0, \quad \forall x, y \in \mathbb{R}^N, \quad \left\| \nabla f(x) - \nabla f(y) \right\|_2 \le M \left\| x - y \right\|_2.$$

Théorème. Soient $\mathbf{f} \in \mathscr{C}^1(\mathbb{R}^N,\mathbb{R})$ strictement convexe et coercive, et K convexe fermé non vide, avec $\exists M>0, \quad \forall \mathbf{x},\mathbf{y} \in \mathbb{R}^N, \quad \|\nabla \mathbf{f}(\mathbf{x})-\nabla \mathbf{f}(\mathbf{y})\|_2 \leq M \left\|\mathbf{x}-\mathbf{y}\right\|_2.$ Si $0<\rho<\frac{2}{M}$ alors la méthode du gradient projeté converge vers l'unique point de minimum global de \mathbf{f} sur K.



Avantage:

convergence vers la solution du problème contraint.



Avantage:

convergence vers la solution du problème contraint.

Inconvénients:

- vitesse de convergence non garantie,
- le projecteur π_K peut être (très) difficile à calculer,



Avantage:

convergence vers la solution du problème contraint.

Inconvénients:

- vitesse de convergence non garantie,
- le projecteur π_K peut être (très) difficile à calculer,

Alternative possible:

Méthode de pénalisation.



Pénalisation

Définition. On appelle fonction de pénalisation de K toute fonction Definition. On appendix. $\beta: \mathbb{R}^N \to \mathbb{R}$ telle que $\beta \in \beta$ est continue $\beta \geq 0$ sur \mathbb{R}^N , $\beta(\mathbf{x}) = 0 \Leftrightarrow \mathbf{x} \in K$



Pénalisation

Définition. On appelle fonction de pénalisation de K toute fonction

- $\beta: \mathbb{R}^N \to \mathbb{R}$ telle que $\beta \text{ est continue}$ $\beta \geq 0 \text{ sur } \mathbb{R}^N,$ $\beta(\mathbf{x}) = 0 \Leftrightarrow \mathbf{x} \in K$

Remarque : Si c'est possible on choisira β convexe. (Il faut K convexe).



Pénalisation

Définition. On appelle fonction de pénalisation de K toute fonction $\beta: \mathbb{R}^N \to \mathbb{R}$ telle que $\beta \text{ est continue}$ $\beta \geq 0 \text{ sur } \mathbb{R}^N,$ $\beta(\mathbf{x}) = 0 \Leftrightarrow \mathbf{x} \in K$

Remarque : Si c'est possible on choisira β convexe. (Il faut K convexe).

Définition. Soit le problème d'optimisation sous contraintes

$$\min_{\mathbf{x} \in K} \mathbf{f}(\mathbf{x}),$$

$$\min_{\mathsf{x} \in \mathbb{R}^N} \mathsf{f}(\mathsf{x}) + \frac{1}{\varepsilon} \beta(\mathsf{x}).$$



Pénalisation

Définition. On appelle fonction de pénalisation de K toute fonction $\beta: \mathbb{R}^N \to \mathbb{R}$ telle que $\beta \text{ est continue}$ $\beta \geq 0 \text{ sur } \mathbb{R}^N,$ $\beta(\mathbf{x}) = 0 \Leftrightarrow \mathbf{x} \in K$

Remarque : Si c'est possible on choisira β convexe. (Il faut K convexe).

Définition. Soit le problème d'optimisation sous contraintes

$$\min_{\mathbf{x} \in K} \mathbf{f}(\mathbf{x})$$

le problème pénalisé associé est

$$\min_{\mathbf{x} \in \mathbb{R}^N} \mathbf{f}(\mathbf{x}) + \frac{1}{\varepsilon} \beta(\mathbf{x}).$$

Pourquoi ca marche? On fait payer (pénalise) le fait de ne pas être dans K.



Pénalisation

▶ Si $K_1 = \{ \mathbf{x} \in \mathbb{R}^N, \ \varphi_i(\mathbf{x}) = 0, \ i = 1 \dots p \},$ on peut choisir

$$\beta(\mathbf{x}) = \sum_{i=1}^{p} \left[\varphi_i(\mathbf{x}) \right]^2.$$



51

Pénalisation

▶ Si $K_1 = \{ \mathbf{x} \in \mathbb{R}^N, \ \varphi_i(\mathbf{x}) = 0, \ i = 1 \dots p \},$ on peut choisir

$$\beta(\mathbf{x}) = \sum_{i=1}^{p} \left[\varphi_i(\mathbf{x}) \right]^2.$$

▶ Si $K_2 = \{ \mathbf{x} \in \mathbb{R}^N, \ \psi_j(\mathbf{x}) \leq 0, \ i = 1 \dots q \}$, on peut choisir

$$\beta(\mathbf{x}) = \sum_{j=1}^{q} \left[\psi_j^+(\mathbf{x}) \right]^2.$$

 $\mathsf{avec}\ \psi_j^+ = \mathsf{max}(\mathsf{0}, \psi_j).$



Pénalisation

▶ Si $K_1 = \{ \mathbf{x} \in \mathbb{R}^N, \ \varphi_i(\mathbf{x}) = 0, \ i = 1 \dots p \}$, on peut choisir

$$\beta(\mathbf{x}) = \sum_{i=1}^{p} \left[\varphi_i(\mathbf{x}) \right]^2.$$

▶ Si $K_2 = \{ \mathbf{x} \in \mathbb{R}^N, \ \psi_j(\mathbf{x}) \leq 0, \ i = 1 \dots q \}$, on peut choisir

$$\beta(\mathbf{x}) = \sum_{j=1}^{q} \left[\psi_j^+(\mathbf{x}) \right]^2.$$

avec $\psi_i^+ = \max(0, \psi_i)$.

▶ Si $K = K_1 \cap K_2$, on somme :

$$\beta(\mathbf{x}) = \sum_{i=1}^{p} \left[\varphi_i(\mathbf{x})\right]^2 + \sum_{i=1}^{q} \left[\psi_i^+(\mathbf{x})\right]^2.$$



Méthode du gradient pénalisé

On résout le problème pénalisé (sans contrainte!)

$$\min_{\mathbf{x} \in \mathbb{R}^N} \mathbf{f}(\mathbf{x}) + \frac{1}{\varepsilon} \beta(\mathbf{x}).$$

par la méthode du gradient à pas fixe.



52

Méthode du gradient pénalisé

On résout le problème pénalisé (sans contrainte!)

$$\min_{\mathbf{x} \in \mathbb{R}^N} \mathbf{f}(\mathbf{x}) + \frac{1}{\varepsilon} \beta(\mathbf{x}).$$

par la méthode du gradient à pas fixe.

Méthode du gradient pénalisé. Elle est définie par

$$\left\{ \begin{array}{l} \mathsf{x}^0 \in \mathbb{R}^N, \\ \mathsf{x}^{n+1} = \mathsf{x}^n - \rho \nabla \mathsf{f}(\mathsf{x}^n) - \frac{\rho}{\varepsilon} \nabla \beta(\mathsf{x}). \end{array} \right.$$
 avec $\rho > 0$.



Méthode du gradient pénalisé

On résout le problème pénalisé (sans contrainte!)

$$\min_{\mathbf{x} \in \mathbb{R}^N} \mathbf{f}(\mathbf{x}) + \frac{1}{\varepsilon} \beta(\mathbf{x}).$$

par la méthode du gradient à pas fixe.

Méthode du gradient pénalisé. Elle est définie par

$$\left\{ \begin{array}{l} \mathsf{x}^0 \in \mathbb{R}^N, \\ \mathsf{x}^{n+1} = \mathsf{x}^n - \rho \nabla \mathsf{f}(\mathsf{x}^n) - \frac{\rho}{\varepsilon} \nabla \beta(\mathsf{x}). \end{array} \right.$$
 avec $\rho > 0$.

Remarque:

- ightharpoonup Choix de ε difficile!
- Conditionne le choix de ρ ...



Méthode du gradient pénalisé

Difficulté : La méthode du gradient pénalisé va converger vers \mathbf{x}_{ε} solution de

$$\min_{\mathbf{u}\in\mathbb{R}^N}\mathbf{f}(\mathbf{x})+\frac{1}{\varepsilon}\beta(\mathbf{x}).$$

Si ε et petit, on s'attend a ce que \mathbf{x}_{ε} soit proche de \mathbf{x}^* , solution de

$$\min_{x \in K} f(x)$$
.



53

Méthode du gradient pénalisé

Difficulté : La méthode du gradient pénalisé va converger vers x_{ε} solution de

$$\min_{\mathbf{u}\in\mathbb{R}^N}\mathbf{f}(\mathbf{x})+\frac{1}{\varepsilon}\beta(\mathbf{x}).$$

Si ε et petit, on s'attend a ce que x_{ε} soit proche de x^* , solution de

$$\min_{x \in K} f(x)$$
.

Théorème. Soient $\mathbf{f} \in \mathscr{C}^1(\mathbb{R}^N,\mathbb{R})$ strictement convexe et coercive, et $K \subset \mathbb{R}^N$ convexe fermé non vide. Soit β une pénalisation de K. Alors pour tout $\varepsilon > 0$, le problème $\min_{\mathbf{x} \in \mathbb{R}^N} \mathbf{f}(\mathbf{x}) + \frac{1}{\varepsilon} \beta(\mathbf{x})$ admet une unique solution \mathbf{x}_{ε} . Elle vérifie

$$\min_{\mathsf{x}\in\mathbb{R}^N}\mathsf{f}(\mathsf{x})+rac{1}{arepsilon}eta(\mathsf{x})$$

$$\lim_{\varepsilon \to 0} \mathsf{x}_{\varepsilon} = \mathsf{x}^*$$

 $\lim_{\varepsilon\to 0} x_\varepsilon = x^*$ où x^* est l'unique solution du problème initial $\min_{x\in K} f(x).$



COURS 4

Approximation numérique des équations différentielles



Quelques usages

► Mécanique :



- ► Mécanique :
 - Prédiction de trajectoire (planètes, comètes, astéroïdes, débris,...)



- Mécanique :
 - Prédiction de trajectoire (planètes, comètes, astéroïdes, débris,...)
 - Propriet des satellites, manœuvre spatiales, rentrée atmosphérique,



Quelques usages

► Mécanique :

- Prédiction de trajectoire (planètes, comètes, astéroïdes, débris,...)
- Propriet des satellites, manœuvre spatiales, rentrée atmosphérique,
- Militaire (balistique),



- ► Mécanique :
 - Prédiction de trajectoire (planètes, comètes, astéroïdes, débris,...)
 - Probite des satellites, manœuvre spatiales, rentrée atmosphérique,
 - Militaire (balistique),
- ▶ Biologie, environnement, :



Quelques usages

► Mécanique :

- Prédiction de trajectoire (planètes, comètes, astéroïdes, débris,...)
- Probite des satellites, manœuvre spatiales, rentrée atmosphérique,
- Militaire (balistique),

► Biologie, environnement, :

Evolution d'un écosystème, d'un population, croissance et extinction,



- Mécanique :
 - Prédiction de trajectoire (planètes, comètes, astéroïdes, débris,...)
 - Probite des satellites, manœuvre spatiales, rentrée atmosphérique,
 - Militaire (balistique),
- ▶ Biologie, environnement, :
 - Evolution d'un écosystème, d'un population, croissance et extinction,
- Chimie, industrie :



- Mécanique :
 - Prédiction de trajectoire (planètes, comètes, astéroïdes, débris,...)
 - Orbite des satellites, manœuvre spatiales, rentrée atmosphérique,
 - Militaire (balistique),
- ► Biologie, environnement, :
 - Evolution d'un écosystème, d'un population, croissance et extinction,
- Chimie, industrie :
 - Dynamique des processus,



- Mécanique :
 - Prédiction de trajectoire (planètes, comètes, astéroïdes, débris,...)
 - Orbite des satellites, manœuvre spatiales, rentrée atmosphérique,
 - Militaire (balistique),
- ▶ Biologie, environnement, :
 - Evolution d'un écosystème, d'un population, croissance et extinction,
- Chimie, industrie :
 - Dynamique des processus,
- Economie, gestion :



- Mécanique :
 - Prédiction de trajectoire (planètes, comètes, astéroïdes, débris,...)
 - Orbite des satellites, manœuvre spatiales, rentrée atmosphérique,
 - Militaire (balistique),
- ▶ Biologie, environnement, :
 - Evolution d'un écosystème, d'un population, croissance et extinction,
- ► Chimie, industrie :
 - Dynamique des processus,
- Economie, gestion :
 - Modélisation macroéconomique,



- ► Mécanique :
 - Prédiction de trajectoire (planètes, comètes, astéroïdes, débris,...)
 - Orbite des satellites, manœuvre spatiales, rentrée atmosphérique,
 - Militaire (balistique),
- ► Biologie, environnement, :
 - Evolution d'un écosystème, d'un population, croissance et extinction,
- Chimie, industrie :
 - Dynamique des processus,
- ► Economie, gestion :
 - Modélisation macroéconomique,
 - Modèles financiers,



- Mécanique :
 - Prédiction de trajectoire (planètes, comètes, astéroïdes, débris,...)
 - Orbite des satellites, manœuvre spatiales, rentrée atmosphérique,
 - Militaire (balistique),
- ▶ Biologie, environnement, :
 - Evolution d'un écosystème, d'un population, croissance et extinction,
- Chimie, industrie :
 - Dynamique des processus,
- ► Economie, gestion :
 - Modélisation macroéconomique,
 - Modèles financiers.
- ► Médecine :



Quelques usages

Mécanique :

- Prédiction de trajectoire (planètes, comètes, astéroïdes, débris,...)
- Orbite des satellites, manœuvre spatiales, rentrée atmosphérique,
- Militaire (balistique),

▶ Biologie, environnement, :

- Evolution d'un écosystème, d'un population, croissance et extinction,
- Chimie, industrie :
 - Dynamique des processus,
- Economie, gestion :
 - Modélisation macroéconomique,
 - Modèles financiers,
- ► Médecine :
 - Posologie, absorption et efficacité des traitements, croissance tumorales,



Quelques usages

Mécanique :

- Prédiction de trajectoire (planètes, comètes, astéroïdes, débris,...)
- Orbite des satellites, manœuvre spatiales, rentrée atmosphérique,
- Militaire (balistique),

▶ Biologie, environnement, :

- Evolution d'un écosystème, d'un population, croissance et extinction,
- Chimie, industrie :
 - Dynamique des processus,
- Economie, gestion :
 - Modélisation macroéconomique,
 - Modèles financiers,
- ► Médecine :
 - Posologie, absorption et efficacité des traitements, croissance tumorales,
 - Epidémiologie.



Une des premières utilisations célèbre





Mercury-Atlas 6 : comment calculer une trajectoire de rentrée atmosphérique ?



Comment viser le Kazakstan avec un Soyouz?

Pourquoi seule l'analyse numérique peut résoudre le problème ?

Quelles méthodes?

Quelle précision?



Constat d'échec

Important : en général, on ne sait pas résoudre explicitement une EDO!



Constat d'échec

Important : en général, on ne sait pas résoudre explicitement une EDO!

On s'en sort quand...



Constat d'échec

Important : en général, on ne sait pas résoudre explicitement une EDO!

On s'en sort quand...

l'équation différentielle est linéaire à coefficients explicitement intégrables,



Constat d'échec

Important : en général, on ne sait pas résoudre explicitement une EDO!

On s'en sort quand...

- l'équation différentielle est linéaire à coefficients explicitement intégrables,
- l'exercice est fait pour qu'on y arrive,



Constat d'échec

Important : en général, on ne sait pas résoudre explicitement une EDO!

On s'en sort quand...

- l'équation différentielle est linéaire à coefficients explicitement intégrables,
- l'exercice est fait pour qu'on y arrive,
- on a beaucoup de chance.



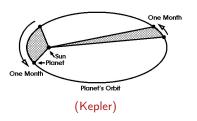
Constat d'échec

Important : en général, on ne sait pas résoudre explicitement une EDO!

On s'en sort quand...

- l'équation différentielle est linéaire à coefficients explicitement intégrables,
- l'exercice est fait pour qu'on y arrive,
- on a beaucoup de chance.

Exemple de coup de chance : les orbites de planètes sont des ellipses.



Équation gravitationelle :

$$x'' = -Gm_S \frac{x}{\left\|x\right\|^3}.$$



Modélisation rentrée atmosphérique





Modélisation du problème de rentrée atmosphérique :

Loi de Newton :

$$m\mathbf{x}'' = \text{ force de gravité} + \text{ force de frottement}$$

$$m\mathbf{x}'' = -\frac{Gmm_T}{\left\|\mathbf{x}\right\|^3}\mathbf{x} - \alpha(\mathbf{x})\left\|\mathbf{x}'\right\|\mathbf{x}'$$



Modélisation rentrée atmosphérique





Modélisation du problème de rentrée atmosphérique :

Loi de Newton :

$$m\mathbf{x}'' = \text{ force de gravit\'e} + \text{ force de frottement}$$

$$m\mathbf{x}'' = -\frac{Gmm_T}{\left\|\mathbf{x}\right\|^3}\mathbf{x} - \alpha(\mathbf{x})\left\|\mathbf{x}'\right\|\mathbf{x}'$$

▶ Inconnue : $\mathbf{x} : [0, T] \to \mathbb{R}^2$ ou \mathbb{R}^3



Modélisation rentrée atmosphérique





Modélisation du problème de rentrée atmosphérique :

► Loi de Newton :

$$m\mathbf{x}'' = \text{ force de gravité} + \text{ force de frottement}$$

$$m\mathbf{x}'' = -\frac{Gmm_T}{\|\mathbf{x}\|^3}\mathbf{x} - \alpha(\mathbf{x}) \|\mathbf{x}'\| \mathbf{x}'$$

- ▶ Inconnue : $\mathbf{x} : [0, T] \to \mathbb{R}^2$ ou \mathbb{R}^3
- ▶ Domaine : pour tout $t \in [0, T]$, $\mathbf{x}(t) \in \Omega := {\mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\| > R_T}$.



Modélisation rentrée atmosphérique





Modélisation du problème de rentrée atmosphérique :

► Loi de Newton :

$$m\mathbf{x}'' = \text{ force de gravité} + \text{ force de frottement}$$

$$m\mathbf{x}'' = -\frac{Gmm_T}{\left\|\mathbf{x}\right\|^3}\mathbf{x} - \alpha(\mathbf{x})\left\|\mathbf{x}'\right\|\mathbf{x}'$$

- ▶ Inconnue : $\mathbf{x} : [0, T] \to \mathbb{R}^2$ ou \mathbb{R}^3
- ▶ Domaine : pour tout $t \in [0, T]$, $\mathbf{x}(t) \in \Omega := {\mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\| > R_T}$.
- ▶ Données : G cste de gravité, m masse du vaisseau, m_T masse de la Terre, $\alpha(\mathbf{x})$ coefficient de frottement aérodynamique, R_T rayon de la Terre.



Problème de Cauchy

Méthode : on met l'équation sous forme d'un problème de Cauchy.

Problème de Cauchy. C'est un problème différentiel de la forme

$$(\mathcal{C}): egin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$



Problème de Cauchy

Méthode : on met l'équation sous forme d'un problème de Cauchy.

Problème de Cauchy. C'est un problème différentiel de la forme

$$(\mathcal{C}): egin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

▶ $\mathbf{u}:[0,T]\to\Omega\subset\mathbb{R}^d$ est la fonction inconnue ,



Problème de Cauchy

Méthode : on met l'équation sous forme d'un problème de Cauchy.

Problème de Cauchy. C'est un problème différentiel de la forme

$$(\mathcal{C}): egin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

- $\mathbf{u}: [0, T] \to \Omega \subset \mathbb{R}^d$ est la fonction inconnue,
- ▶ le domaine spatial Ω est une partie ouverte connexe de l'espace \mathbb{R}^d ,



Problème de Cauchy

Méthode : on met l'équation sous forme d'un problème de Cauchy.

Problème de Cauchy. C'est un problème différentiel de la forme

$$(\mathcal{C}): egin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

- $\mathbf{u}: [0, T] \to \Omega \subset \mathbb{R}^d$ est la fonction inconnue,
- ▶ le domaine spatial Ω est une partie ouverte connexe de l'espace $ℝ^d$,
- $\mathbf{f}: \Omega \times [0, T] \mapsto \mathbb{R}^d$ est la **dynamique** du problème,



58

Problème de Cauchy

Méthode : on met l'équation sous forme d'un problème de Cauchy.

Problème de Cauchy. C'est un problème différentiel de la forme

$$(\mathcal{C}): egin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

- $\mathbf{u}: [0, T] \to \Omega \subset \mathbb{R}^d$ est la fonction inconnue,
- ▶ le domaine spatial Ω est une partie ouverte connexe de l'espace $ℝ^d$,
- ▶ $\mathbf{f}: \Omega \times [0, T] \mapsto \mathbb{R}^d$ est la **dynamique** du problème,
- ▶ $\mathbf{u}_0 \in \Omega$ est l'état initial où condition initiale.

Problème de Cauchy pour la rentrée atmosphérique

Exemple: rentrée atmosphérique,

$$\mathbf{x}'' = -\frac{Gm_T}{\|\mathbf{x}\|^3} \mathbf{x} - \frac{\alpha(\mathbf{x})}{m} \|\mathbf{x}'\| \mathbf{x}'$$



Problème de Cauchy pour la rentrée atmosphérique

Exemple: rentrée atmosphérique,

$$\mathbf{x}'' = -\frac{Gm_T}{\|\mathbf{x}\|^3} \mathbf{x} - \frac{\alpha(\mathbf{x})}{m} \|\mathbf{x}'\| \mathbf{x}'$$
 on pose $\mathbf{u} := \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{bmatrix} := \begin{bmatrix} \mathbf{x} \\ \mathbf{x}' \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_1' \\ x_2' \end{bmatrix} \in \mathbb{R}^4 \text{ et on a}$
$$\mathbf{f}(\mathbf{u}) = \begin{bmatrix} u_3 \\ u_4 \\ -\frac{Gm_T}{\|(u_1,u_2)\|^3} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} - \frac{\alpha(u_1,u_2)}{m} \|(u_3,u_4)\| \begin{bmatrix} u_3 \\ u_4 \end{bmatrix} \end{bmatrix}.$$



Problème de Cauchy pour la rentrée atmosphérique

Exemple: rentrée atmosphérique,

$$\mathbf{x}'' = -\frac{Gm_T}{\|\mathbf{x}\|^3} \mathbf{x} - \frac{\alpha(\mathbf{x})}{m} \|\mathbf{x}'\| \mathbf{x}'$$
 on pose $\mathbf{u} := \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{bmatrix} := \begin{bmatrix} \mathbf{x} \\ \mathbf{x}' \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_1' \\ x_2' \end{bmatrix} \in \mathbb{R}^4 \text{ et on a}$
$$\mathbf{f}(\mathbf{u}) = \begin{bmatrix} u_3 \\ u_4 \\ -\frac{Gm_T}{\|(u_1,u_2)\|^3} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} - \frac{\alpha(u_1,u_2)}{m} \|(u_3,u_4)\| \begin{bmatrix} u_3 \\ u_4 \end{bmatrix} \end{bmatrix}.$$

Le problème de Cauchy s'écrit

$$\begin{cases} \mathbf{u}'(t) = \mathbf{f}(\mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \begin{bmatrix} \mathbf{x}_0 \\ \mathbf{v}_0 \end{bmatrix}. \end{cases}$$



Comment approcher numériquement la solution du problème de Cauchy?

Soit le problème de Cauchy

$$(\mathcal{C}): \begin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$



Comment approcher numériquement la solution du problème de Cauchy?

Soit le problème de Cauchy

$$(\mathcal{C}): \begin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

Existence et unicité de $u:[0,T^*[\to\Omega]$ dès que $\mathbf{f}\in\mathscr{C}^1$ [Cauchy-Lipschitz].



Comment approcher numériquement la solution du problème de Cauchy?

Soit le problème de Cauchy

$$(\mathcal{C}): \begin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

Existence et unicité de $u: [0, T^*[\to \Omega \text{ dès que } \mathbf{f} \in \mathscr{C}^1 \text{ [Cauchy-Lipschitz]}.$

Discrétisation du temps : soit h > 0, et par $t_n = nh$ avec $t_N = N_h h = T$.



Comment approcher numériquement la solution du problème de Cauchy?

Soit le problème de Cauchy

$$(\mathcal{C}): \begin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

Existence et unicité de $u: [0, T^*[\to \Omega \text{ dès que } \mathbf{f} \in \mathscr{C}^1 \text{ [Cauchy-Lipschitz]}.$

Discrétisation du temps: soit h > 0, et par $t_n = nh$ avec $t_N = N_h h = T$. Cette suite dépend de h, on la note donc aussi $(t_n^h)_{n=0,1,...,N}$.



Comment approcher numériquement la solution du problème de Cauchy?

Soit le problème de Cauchy

$$(\mathcal{C}): \begin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

Existence et unicité de $u: [0, T^*[\to \Omega \text{ dès que } f \in \mathscr{C}^1 \text{ [Cauchy-Lipschitz]}.$

Discrétisation du temps: soit h > 0, et par $t_n = nh$ avec $t_N = N_h h = T$. Cette suite dépend de h, on la note donc aussi $(t_n^h)_{n=0,1,\ldots,N}$.

Principe d'une méthode numérique : approcher les valeurs de la solution exacte $\mathbf{u}(t_n)$ par une suite $(\mathbf{U}_n)_{n=0,1,\ldots,N}$ de Ω .



Comment approcher numériquement la solution du problème de Cauchy?

Soit le problème de Cauchy

$$(\mathcal{C}): \begin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

Existence et unicité de $u: [0, T^*[\to \Omega \text{ dès que } f \in \mathscr{C}^1 \text{ [Cauchy-Lipschitz]}.$

Discrétisation du temps : soit h > 0, et par $t_n = nh$ avec $t_N = N_h h = T$. Cette suite dépend de h, on la note donc aussi $(t_n^h)_{n=0,1,\ldots,N}$.

Principe d'une méthode numérique : approcher les valeurs de la solution exacte $\mathbf{u}(t_n)$ par une suite $(\mathbf{U}_n)_{n=0,1,\ldots,N}$ de Ω .

Vocabulaire : Un **schéma numérique** est une relation de récurrence qui permet de définir la suite (\mathbf{U}_n) .



Comment approcher numériquement la solution du problème de Cauchy?

Soit le problème de Cauchy

$$(\mathcal{C}): \begin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

Existence et unicité de $u: [0, T^*[\to \Omega \text{ dès que } f \in \mathscr{C}^1 \text{ [Cauchy-Lipschitz]}.$

Discrétisation du temps: soit h > 0, et par $t_n = nh$ avec $t_N = N_h h = T$. Cette suite dépend de h, on la note donc aussi $(t_n^h)_{n=0,1,\ldots,N}$.

Principe d'une méthode numérique : approcher les valeurs de la solution exacte $\mathbf{u}(t_n)$ par une suite $(\mathbf{U}_n)_{n=0,1,\ldots,N}$ de Ω .

Vocabulaire : Un schéma numérique est une relation de récurrence qui permet de définir la suite (\mathbf{U}_n) .

Cette suite est dite solution numérique ou discrète.



La méthode d'Euler

Problème de Cauchy

$$(\mathcal{C}): \begin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

▶ On écrit, pour $n = 0, 1, ..., N_h - 1$,

$$u(t_{n+1}) = u(t_n) + \int_{t_n}^{t_{n+1}} u'(s) ds = u(t_n) + \int_{t_n}^{t_{n+1}} f(s, u(s)) ds.$$

Approximation de l'intégrale par la méthode des rectangles :

$$\int_{t_n}^{t_{n+1}} f(s, u(s)) ds \approx h f(t_n, u(t_n)).$$

D'où la récurrence (explicite mais inexacte) :

$$\mathbf{u}(t_{n+1}) \approx \mathbf{u}(t_n) + h \mathbf{f}(t_n, \mathbf{u}(t_n)).$$



Approximation des EDO Méthode d'Euler

On a obtenu:

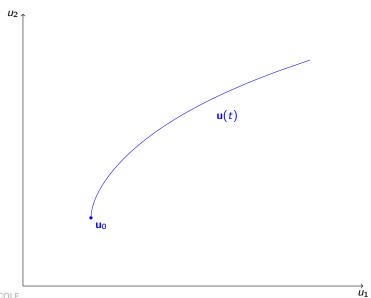
$$\mathbf{u}(t_{n+1}) \approx \mathbf{u}(t_n) + h \mathbf{f}(t_n, \mathbf{u}(t_n)).$$

Méthode d'Euler. C'est le système de récurrence/condition initiale suivant :

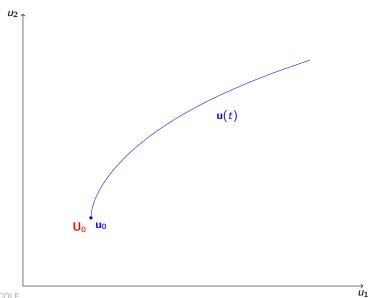
$$(E): egin{cases} \mathbf{U}_{n+1} = \mathbf{U}_n + h \ \mathbf{f}(t_n, \mathbf{U}_n), & n = 0, 1, \dots, N_h - 1 \ \mathbf{U}_0 = \mathbf{u}_0. \end{cases}$$

Remarques.

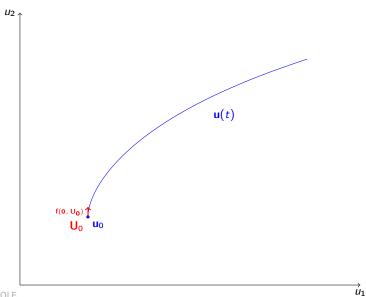
- ightharpoonup a priori, on n'a pas $\mathbf{u}(t_n) = U_n$.
- ightharpoonup il n'est pas évident non plus que $\mathbf{u}(t_n)$ approche U_n .



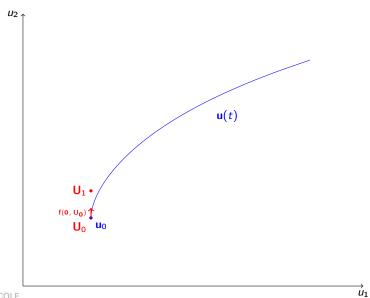




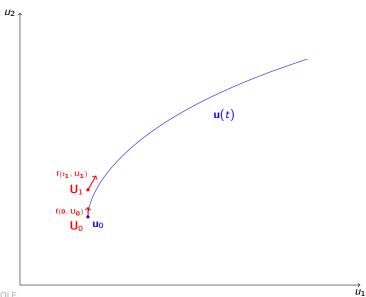




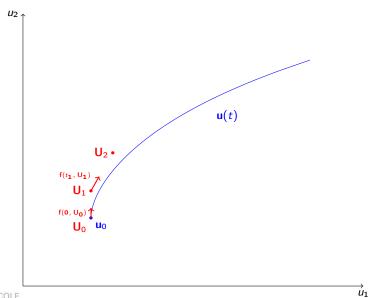




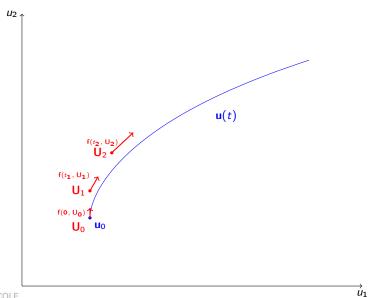




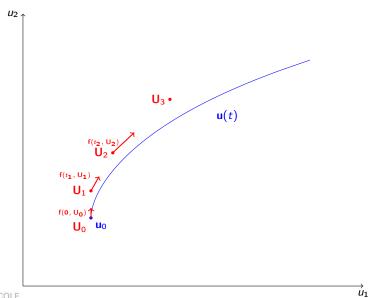




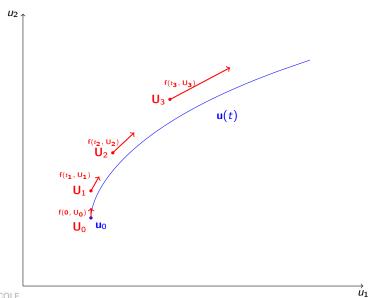




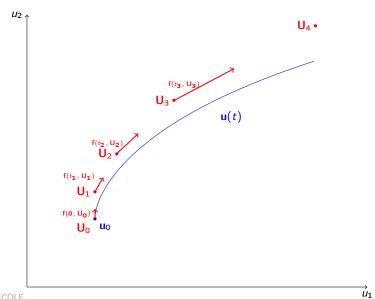






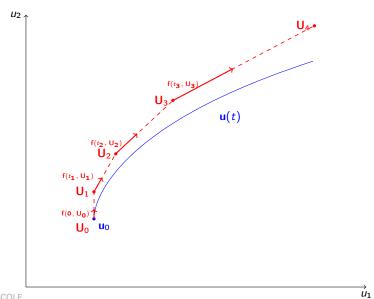








Méthode d'Euler





Méthode d'Euler : Exemple

Considérons le problème de Cauchy

$$\begin{cases} u'(t) = 2tu(t) & t \in [0, 1], \\ u(0) = 1, \end{cases}$$

dont la solution exacte est $u(t) = e^{t^2}$.



Méthode d'Euler : Exemple

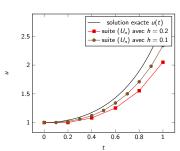
Considérons le problème de Cauchy

$$\begin{cases} u'(t) = 2tu(t) & t \in [0,1], \\ u(0) = 1, \end{cases}$$

dont la solution exacte est $u(t) = e^{t^2}$.

La méthode d'Euler s'écrit

$$\begin{cases} U_{n+1} = U_n + 2ht_nU_n & t \in [0, 2], \\ U_0 = 1. \end{cases}$$





Méthode d'Euler : Exemple

Considérons le problème de Cauchy

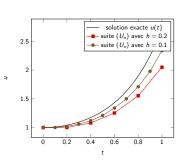
$$\begin{cases} u'(t) = 2tu(t) & t \in [0,1], \\ u(0) = 1, \end{cases}$$

dont la solution exacte est $u(t) = e^{t^2}$.

La méthode d'Euler s'écrit

$$\begin{cases} U_{n+1} = U_n + 2ht_nU_n & t \in [0, 2], \\ U_0 = 1. \end{cases}$$

La méthode d'Euler n'est pas exacte,





Méthode d'Euler : Exemple

Considérons le problème de Cauchy

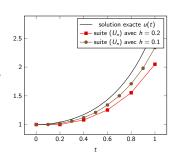
$$\begin{cases} u'(t) = 2tu(t) & t \in [0,1], \\ u(0) = 1, \end{cases}$$

dont la solution exacte est $u(t) = e^{t^2}$.

La méthode d'Euler s'écrit

$$\begin{cases} U_{n+1} = U_n + 2ht_nU_n & t \in [0,2], \\ U_0 = 1. \end{cases}$$

- La méthode d'Euler n'est pas exacte,
- ► l'erreur augmente avec le temps,





Méthode d'Euler : Exemple

Considérons le problème de Cauchy

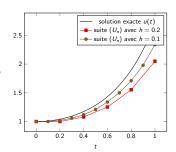
$$\begin{cases} u'(t) = 2tu(t) & t \in [0, 1], \\ u(0) = 1, \end{cases}$$

dont la solution exacte est $u(t) = e^{t^2}$.

La méthode d'Euler s'écrit

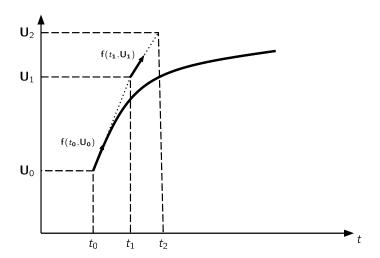
$$\begin{cases} U_{n+1} = U_n + 2ht_nU_n & t \in [0, 2], \\ U_0 = 1. \end{cases}$$

- La méthode d'Euler n'est pas exacte.
- ► l'erreur augmente avec le temps,
- ▶ l'erreur diminue si *h* diminue.



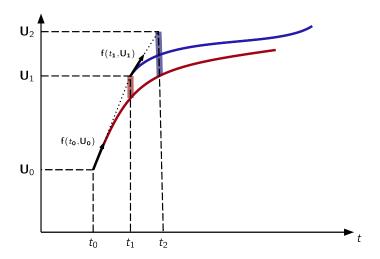


Analyse d'erreur





Analyse d'erreur





Consistance, stabilité, convergence

► Convergence : L'erreur globale du schéma tend vers 0

$$\max_{0\leqslant n\leqslant N_h}\|\mathbf{U}_n-\mathbf{u}(t_n)\|\xrightarrow[h\to 0]{?}0.$$



Consistance, stabilité, convergence

► Convergence : L'erreur globale du schéma tend vers 0

$$\max_{\mathbf{0}\leqslant n\leqslant N_h}\|\mathbf{U}_n-\mathbf{u}(t_n)\|\xrightarrow[h\to 0]{?}0.$$

- ► Analyse en 2 temps :
 - ► Consistance : analyse de l'erreur locale (commise à chaque pas).



Consistance, stabilité, convergence

► Convergence : L'erreur globale du schéma tend vers 0

$$\max_{\mathbf{0}\leqslant n\leqslant N_h}\|\mathbf{U}_n-\mathbf{u}(t_n)\|\xrightarrow[h\to 0]{?}0.$$

- ► Analyse en 2 temps :
 - ► Consistance : analyse de l'erreur locale (commise à chaque pas).
 - ▶ Stabilité : maîtrise de la propagation des erreurs locales.



Consistance, stabilité, convergence

► Convergence : L'erreur globale du schéma tend vers 0

$$\max_{0\leqslant n\leqslant N_h}\|\mathbf{U}_n-\mathbf{u}(t_n)\|\xrightarrow[h\to 0]{?}0.$$

- ► Analyse en 2 temps :
 - Consistance : analyse de l'erreur locale (commise à chaque pas).
 - ► Stabilité : maîtrise de la propagation des erreurs locales.

Théorème de Lax.

Consistance et Stabilité \Rightarrow Convergence.



Consistance

On note **u** la solution exacte du problème de Cauchy.

Définition. L'erreur de consistance du schéma d'Euler est la suite

$$arepsilon_n^h := rac{ \mathbf{\mathsf{u}}(t_{n+1}) - \mathbf{\mathsf{u}}(t_n)}{h} - \mathsf{f}(t_n, \mathbf{\mathsf{u}}(t_n)), \quad 0 \leqslant n \leqslant N_h - 1.$$



67

Consistance

On note **u** la solution exacte du problème de Cauchy.

Définition. L'erreur de consistance du schéma d'Euler est la suite

$$arepsilon_n^h := rac{ \mathbf{u}(t_{n+1}) - \mathbf{u}(t_n)}{h} - \mathbf{f}(t_n, \mathbf{u}(t_n)), \quad 0 \leqslant n \leqslant N_h - 1.$$

Remarque: Lien avec l'approximation de l'intégration de \mathbf{u}' sur $[t_n, t_n + h]$ par la méthode des rectangles à gauche:

$$arepsilon_n^h := \int_0^1 \mathbf{u}'(t_n + hs) \mathrm{d}s - \mathbf{u}'(t_n).$$



67

Consistance

Rappel: u désigne la solution exacte du problème de Cauchy.

Définition. L'erreur de consistance) du schéma d'Euler est la suite

$$arepsilon_n^h := rac{ extbf{u}(t_{n+1}) - extbf{u}(t_n)}{h} - extbf{f}(t_n, extbf{u}(t_n)), \quad 0 \leqslant n \leqslant N_h - 1.$$

Définition. Le schéma est consistant si l'erreur de consistance converge vers 0: $\max_{0\leqslant n\leqslant N_h-1}\left\|\varepsilon_n^h\right\|\to 0\quad \text{quand}\quad h\to 0.$ Si de plus il existe k>0 tel que $\max_{0\leqslant n\leqslant N_h-1}\left\|\varepsilon_n^h\right\|=\mathcal{O}(h^k),$ le schéma est dit consistant d'ordre k.

$$\max_{0 \leq n \leq N_t - 1} \left\| \varepsilon_n^h \right\| o 0$$
 quand $h o 0$

$$\max_{0 \le n \le N_h - 1} \left\| \varepsilon_n^h \right\| = \mathcal{O}(h^k),$$

Consistance pour Euler

Proposition. Si la solution exacte \mathbf{u} est de classe \mathscr{C}^2 , alors la méthode d'Euler est consistante d'ordre 1.



Consistance pour Euler

Proposition. Si la solution exacte \mathbf{u} est de classe \mathscr{C}^2 , alors la méthode d'Euler est consistante d'ordre 1.

Preuve: Par Taylor-Lagrange,

$$u(t_{n+1}) = u(t_n + h) = u(t_n) + hu'(t_n) + h^2R_n(h),$$

avec $\mathbf{R}_n(h)$ bornée. Plus précisément, $\|\mathbf{R}_n(h)\| \leqslant C := \frac{1}{2} \sup_{t \in [0,T]} \left\| \mathbf{u}''(t) \right\|$

$$arepsilon_n^h := rac{\mathsf{u}(t_{n+1}) - \mathsf{u}(t_n)}{h} - \mathsf{f}(t_n, \mathsf{u}(t_n)),$$
 $arepsilon_n^h := \mathsf{u}'(t_n) + h\mathsf{R}_n(h) - \mathsf{f}(t_n, \mathsf{u}(t_n)),$
 $arepsilon_n^h := h\mathsf{R}_n(h).$

Ainsi

$$\max_{0 \leqslant n \leqslant N_h - 1} \left\| \varepsilon_n^h \right\| \leqslant Ch.$$



Approximation des EDO Stabilité pour Euler

On considère le schéma d'Euler

$$(E): \begin{cases} \mathbf{U}_{n+1} = \mathbf{U}_n + h \ \mathbf{f}(t_n, \mathbf{U}_n), & n \geqslant 0, \\ \mathbf{U}_0 = \mathbf{u}_0. \end{cases}$$

et on le perturbe en ajoutant une suite d'erreurs μ_n



Stabilité pour Euler

On considère le schéma d'Euler

$$(E): \begin{cases} \mathbf{U}_{n+1} = \mathbf{U}_n + h \ \mathbf{f}(t_n, \mathbf{U}_n), & n \geqslant 0, \\ \mathbf{U}_0 = \mathbf{u}_0. \end{cases}$$

et on le perturbe en ajoutant une suite d'erreurs μ_n

$$(E_{pert}): \begin{cases} \mathbf{V}_{n+1} = \mathbf{V}_n + h \ \mathbf{f}(t_n, \mathbf{V}_n) + \boldsymbol{\mu}_n, & n \geqslant 0, \\ \mathbf{V}_0 = \mathbf{u}_0. \end{cases}$$



Stabilité pour Euler

On considère le schéma d'Euler

$$(E): \begin{cases} \mathbf{U}_{n+1} = \mathbf{U}_n + h \ \mathbf{f}(t_n, \mathbf{U}_n), & n \geqslant 0, \\ \mathbf{U}_0 = \mathbf{u}_0. \end{cases}$$

et on le perturbe en ajoutant une suite d'erreurs μ_n

$$(E_{\mathsf{pert}}): egin{cases} \mathbf{V}_{n+1} = \mathbf{V}_n + h \ \mathbf{f}(\mathbf{t}_n, \mathbf{V}_n) + \mu_n, & n \geqslant 0, \\ \mathbf{V}_0 = \mathbf{u}_0. \end{cases}$$

Définition. Le schéma est stable s'il existe une constante C > 0 telle que

$$\max_{\mathbf{0}\leqslant n\leqslant N_h}\|\mathbf{V}_n-\mathbf{U}_n\|\leqslant C\sum_{\mathbf{0}\leqslant n\leqslant N_h-\mathbf{1}}\|\mu_n\|\ .$$



Stabilité pour Euler

Proposition. Si f est L-lipschitzienne, la méthode d'Euler est stable et l'on a

max
$$|V_n-U_n|\leqslant e^{LT}\sum_{0\leqslant n\leqslant N_h-1}\|\mu_n\|$$
 .



71

Stabilité pour Euler

Proposition. Si f est L-lipschitzienne, la méthode d'Euler est stable et l'on a

$$\max_{0\leqslant n\leqslant N_h} |V_n - U_n| \leqslant e^{LT} \sum_{0\leqslant n\leqslant N_h - 1} \|\mu_n\|.$$

Preuve: Par soustraction

$$\begin{aligned} & \mathsf{V}_{n+1} - \mathsf{U}_{n+1} = \mathsf{V}_n - \mathsf{U}_n + h \big(\mathsf{f}(\mathsf{t}_n, \mathsf{V}_n) - \mathsf{f}(\mathsf{t}_n, \mathsf{U}_n) \big) + \mu_n, \\ & \| \mathsf{V}_{n+1} - \mathsf{U}_{n+1} \| \leqslant \| \mathsf{V}_n - \mathsf{U}_n \| + h \| \mathsf{f}(\mathsf{t}_n, \mathsf{V}_n) - \mathsf{f}(\mathsf{t}_n, \mathsf{U}_n) \| + \| \mu_n \| , \\ & \| \mathsf{V}_{n+1} - \mathsf{U}_{n+1} \| \leqslant (1 + hL) \| \mathsf{V}_n - \mathsf{U}_n \| + \| \mu_n \| . \end{aligned}$$

Posons $e_n = ||\mathbf{V}_n - \mathbf{U}_n||$, et divisons par $(1 + hL)^{n+1}$;

$$\frac{e_{n+1}}{(1+hL)^{n+1}} \leqslant \frac{e_n}{(1+hL)^n} + \frac{\|\mu_n\|}{(1+hL)^{n+1}},$$
$$\frac{e_{n+1}}{(1+hL)^{n+1}} \leqslant \frac{e_n}{(1+hL)^n} + \|\mu_n\|,$$



Stabilité pour Euler

Preuve (suite): On a

$$\frac{e_{n+1}}{(1+hL)^{n+1}} - \frac{e_n}{(1+hL)^n} \leqslant \|\mu_n\|,$$

En sommant, on obtient

$$\frac{e_n}{(1+hL)^n} \leqslant \sum_{n=0}^{n-1} \|\mu_k\|, \quad \forall n \in \mathbb{N},$$

$$e_n \leqslant (1+hL)^n \sum_{k=0}^{n-1} \|\mu_k\|, \quad \forall n \in \mathbb{N}.$$

Or pour tout $n \leq N_h$,

$$(1+hL)^n \leqslant e^{Lt_N} \leqslant e^{N_h hL} = e^{LT}.$$

Ainsi,

$$\|\mathbf{V}_n - \mathbf{U}_n\| \leqslant e^{LT} \sum_{k=0}^{n-1} \|\mu_k\|, \quad \forall n = 0, 1, \dots, N_h.$$

Convergence pour Euler

Définition. Le schéma converge si

definition. Le schéma converge si
$$\max_{0\leqslant n\leqslant N_h}|\mathbf{U}_n-\mathbf{u}(t_n)|\longrightarrow 0 \quad ext{ quand } \quad h\to 0.$$



Convergence pour Euler

Définition. Le schéma converge si
$$\max_{0\leqslant n\leqslant N_h} |\mathbf{U}_n - \mathbf{u}(t_n)| \longrightarrow 0 \quad \text{ quand } \quad h \to 0.$$

Théorème de Lax. Un schéma consistant et stable est convergent.



73

Convergence pour Euler

Définition. Le schéma converge si

$$\max_{0\leqslant n\leqslant N_h} |\mathsf{U}_n - \mathsf{u}(t_n)| \longrightarrow 0 \quad \text{ quand } \quad h \to 0.$$

Théorème de Lax. Un schéma consistant d'ordre *k* et stable est convergent d'ordre *k*.

Preuve (Pour Euler): Par définition,

$$\varepsilon_n^h := \frac{\mathsf{u}(t_{n+1}) - \mathsf{u}(t_n)}{h} - \mathsf{f}(t_n, \mathsf{u}(t_n)).$$

Donc la suite de valeurs exactes est solution du schéma perturbé :

$$\begin{cases} \mathsf{u}(t_{n+1}) = \mathsf{u}(t_n) + h\mathsf{f}(t_n, \mathsf{u}(t_n)) + h\varepsilon_n^h, \\ \mathsf{u}(t_0) = u_0. \end{cases}$$



Convergence pour Euler

Preuve (suite) : On utilise la stabilité du schéma :

$$\max_{0\leqslant n\leqslant N_h} |\mathsf{U}_n - \mathsf{u}(t_n)| \leqslant e^{LT} h \sum_{0\leqslant n\leqslant N_h-1} \left\| \varepsilon_n^h \right\|.$$

Or,

$$\sum_{0 \le n \le N_{t-1}} \left\| \varepsilon_{n}^{h} \right\| \leqslant N_{h} \max_{0 \leqslant n \leqslant N_{h}-1} \left\| \varepsilon_{n}^{h} \right\|,$$

Ainsi,

$$\max_{0 \leqslant n \leqslant N_h} |\mathsf{U}_n - \mathsf{u}(t_n)| \leqslant T e^{LT} \max_{0 \leqslant n \leqslant N_h - 1} \left\| \varepsilon_n^h \right\|.$$

Or, par consistance d'ordre 1,

$$\max_{0 \leqslant n \leqslant N_h - 1} \left\| \varepsilon_n^h \right\| \leqslant Ch.$$

Finalement,

$$\max_{0\leqslant n\leqslant N_h} |\mathbf{U}_n - \mathbf{u}(t_n)| \leqslant CTe^{LT}h.$$



Convergence pour Euler

Proposition. Le schéma d'Euler est convergent d'ordre 1 :

$$\max_{0 \le n \le N_t} |\mathbf{U}_n - \mathbf{u}(t_n)| \leqslant CT e^{LT} h$$

$$\max_{0\leqslant n\leqslant N_h} |\mathbf{U}_n - \mathbf{u}(t_n)| \leqslant CT \mathrm{e}^{LT} h,$$
 avec $C = \frac{1}{2} \sup_{t\in[0,T]} \left\|\mathbf{u}''(t)\right\|$.



Convergence pour Euler

Proposition. Le schéma d'Euler est convergent d'ordre 1 :

$$\max_{0\leqslant n\leqslant N_h} |\mathbf{U}_n - \mathbf{u}(t_n)| \leqslant CT e^{LT} h,$$
 avec $C = \frac{1}{2} \sup_{t\in[0,T]} \left\| \mathbf{u}''(t) \right\|$.

avec
$$C = \frac{1}{2} \sup_{t \in [0,T]} \left\| \mathbf{u}''(t) \right\|$$
.

Remarques:

- ▶ On a besoin que $\mathbf{u} \in \mathcal{C}^2([0, T])$.
- ► La constante *Te*^{LT} augmente exponentiellement
 - ▶ avec T (problèmes en temps long) ,
 - avec L (problèmes « raides »).

Schémas classiques

Définition. Un schéma à un pas (constant) pour approcher la solution du

$$(\mathcal{C}): egin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

Définition. Un schéma à un pas (constant) pour approcher la solution problème de Cauchy
$$(\mathcal{C}): \begin{cases} \mathbf{u}'(t) = \mathbf{f}(t,\mathbf{u}(t)) & t \in [0,T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$
 est une relation de récurrence de la forme
$$\begin{cases} \mathbf{U}_{n+1} = \mathbf{U}_n + h\mathbf{G}_h(t_n,\mathbf{U}_n,\mathbf{U}_{n+1}), & n=0,1,\ldots,N_h-1, \\ \mathbf{U}_0 = \mathbf{u}_0. \end{cases}$$
 où $h>0$ est le pas et \mathbf{G}_h est une fonction qui définit le schéma.



Schémas classiques

Définition. Un schéma à un pas (constant) pour approcher la solution du

$$(\mathcal{C}): egin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

Définition. Un schéma à un pas (constant) pour approcher la solution problème de Cauchy
$$(\mathcal{C}): \begin{cases} \mathbf{u}'(t) = \mathbf{f}(t,\mathbf{u}(t)) & t \in [0,T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$
 est une relation de récurrence de la forme
$$\begin{cases} \mathbf{U}_{n+1} = \mathbf{U}_n + h\mathbf{G}_h(t_n,\mathbf{U}_n,\mathbf{U}_{n+1}), & n=0,1,\ldots,N_h-1, \\ \mathbf{U}_0 = \mathbf{u}_0. \end{cases}$$
 où $h>0$ est le pas et \mathbf{G}_h est une fonction qui définit le schéma.

Vocabulaire : Si G_h dépend effectivement de U_{n+1} le schéma est dit implicite, sinon il est explicite.



Schémas classiques

Définition. Un schéma à un pas (constant) pour approcher la solution du

$$(\mathcal{C}): egin{cases} \mathbf{u}'(t) = \mathbf{f}(t, \mathbf{u}(t)) & t \in [0, T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$

Définition. Un schéma à un pas (constant) pour approcher la solution problème de Cauchy
$$(\mathcal{C}): \begin{cases} \mathbf{u}'(t) = \mathbf{f}(t,\mathbf{u}(t)) & t \in [0,T], \\ \mathbf{u}(0) = \mathbf{u}_0. \end{cases}$$
 est une relation de récurrence de la forme
$$\begin{cases} \mathbf{U}_{n+1} = \mathbf{U}_n + h\mathbf{G}_h(t_n,\mathbf{U}_n,\mathbf{U}_{n+1}), & n=0,1,\ldots,N_h-1, \\ \mathbf{U}_0 = \mathbf{u}_0. \end{cases}$$
 où $h>0$ est le pas et \mathbf{G}_h est une fonction qui définit le schéma.

Vocabulaire : Si G_h dépend effectivement de U_{n+1} le schéma est dit implicite, sinon il est explicite.

Exemple: Le schéma d'Euler $\mathbf{U}_{n+1} = \mathbf{U}_n + h\mathbf{f}(t_n, \mathbf{U}_n)$ correspond à

$$G_h(t_n, \mathbf{U}_n, \mathbf{U}_{n+1}) = \mathbf{f}(t_n, \mathbf{U}_n).$$



Consistance

Rappel: u désigne la solution exacte du problème de Cauchy.

Définition. L'erreur de consistance d'un schéma général à un pas est

Definition. Lerreur de consistance d'un schema general a un pas est
$$\varepsilon_n^h := \frac{\mathbf{u}(t_{n+1}) - \mathbf{u}(t_n)}{h} - G_h(t_n, \mathbf{u}(t_n), \mathbf{u}(t_{n+1})), \quad 0 \leqslant n \leqslant N_h.$$

Remarque : Lien avec l'erreur de l'approximation de l'intégration de \mathbf{u}' sur $[t_n, t_n + h]$:

$$\varepsilon_n^h := \int_0^1 \mathbf{u}'(t_n + hs) ds - G_h(t_n, \mathbf{u}(t_n), \mathbf{u}(t_{n+1})).$$



Approximation des EDO Stabilité

On considère un schéma général

$$(E): \begin{cases} \mathbf{U}_{n+1} = \mathbf{U}_n + h\mathbf{G}_h(t_n, \mathbf{U}_n, \mathbf{U}_{n+1}), & n = 0, 1, \dots, N_h - 1, \\ \mathbf{U}_0 = \mathbf{u}_0. \end{cases}$$

et on le perturbe en ajoutant une suite d'erreurs μ_n :

$$(E_{\mathsf{pert}}): egin{cases} \mathbf{V}_{n+1} = \mathbf{V}_n + h \mathbf{G}_h(t_n, \mathbf{V}_n, \mathbf{V}_{n+1}) + \mu_n, & n = 0, 1, \dots, N_h - 1, \\ \mathbf{V}_0 = \mathbf{u}_0. \end{cases}$$

Définition. Le schéma est stable s'il existe C>0 t.q.

$$\max_{0\leqslant n\leqslant N_h}\|\mathsf{V}_n-\mathsf{U}_n\|\leqslant C\sum_{0\leqslant n\leqslant N_h-1}\|\mu_n\|\ .$$



Euler, Euler implicite

$$\int_{t_n}^{t_{n+1}} u'(t) dt = h \int_0^1 u'(t_n + hs) ds$$

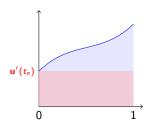
Rectangle à gauche :

$$I(h) = \int_0^1 \mathbf{u}'(t_n + hs) ds$$

 $\approx \mathbf{u}'(t_n) = \mathbf{f}(t_n, \mathbf{u}(t_n))$

conduit à poser

Euler: $U_{n+1} = U_n + hf(t_n, U_n)$.



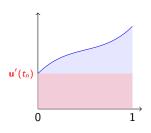
Euler, Euler implicite

Rectangle à gauche :

$$I(h) = \int_0^1 \mathbf{u}'(t_n + hs) ds$$
$$\approx \mathbf{u}'(t_n) = \mathbf{f}(t_n, \mathbf{u}(t_n))$$

conduit à poser

Euler: $U_{n+1} = U_n + hf(t_n, U_n)$.



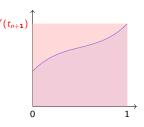
Rectangle à droite :

$$I(h) = \int_0^1 \mathbf{u}'(t_n + hs) ds$$

 $\approx \mathbf{u}'(t_{n+1}) = \mathbf{f}(t_{n+1}, \mathbf{u}(t_{n+1}))$

conduit à poser

Euler (implicite) : $U_{n+1} = U_n + hf(t_{n+1}, U_{n+1})$



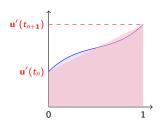
Crank-Nicolson, Heun

Trapèze:

$$I(h) \approx \frac{\mathbf{u}'(t_n) + \mathbf{u}'(t_{n+1})}{2}$$

Crank-Nicolson (implicite)

$$\mathbf{U}_{n+1} = \mathbf{U}_n + \frac{h}{2} \left(\mathbf{f}(t_n, \mathbf{U}_n) + f(t_{n+1}, \mathbf{U}_{n+1}) \right)$$



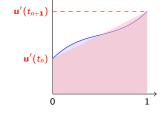
Crank-Nicolson, Heun

Trapèze:

$$I(h) \approx \frac{\mathbf{u}'(t_n) + \mathbf{u}'(t_{n+1})}{2}$$

Crank-Nicolson (implicite)

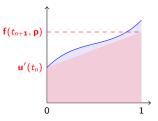
$$U_{n+1} = U_n + \frac{h}{2} (f(t_n, U_n) + f(t_{n+1}, U_{n+1}))$$



Trapèze approché : Le vecteur **p** approche $\mathbf{u}(t_{n+1})$ par Euler :

Heun (explicite)

$$\begin{vmatrix} \mathbf{p} = \mathbf{U}_n + h\mathbf{f}(t_n, \mathbf{U}_n). \\ \mathbf{U}_{n+1} = \mathbf{U}_n + \frac{h}{2} \left(\mathbf{f}(t_n, \mathbf{U}_n) + f(t_{n+1}, \mathbf{p}) \right) \end{vmatrix}$$



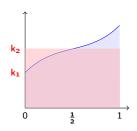
Runge-Kutta

Point milieu approché : RK-2 (explicite)

$$k_{1} = f(t_{n}, U_{n}).$$

$$k_{2} = f\left(t_{n} + \frac{h}{2}, U_{n} + \frac{h}{2}k_{1}\right).$$

$$U_{n+1} = U_{n} + hk_{2}.$$



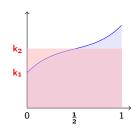
Runge-Kutta

Point milieu approché : RK-2 (explicite)

$$\mathbf{k}_{1} = \mathbf{f}(t_{n}, \mathbf{U}_{n}).$$

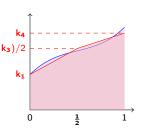
$$\mathbf{k}_{2} = \mathbf{f}\left(t_{n} + \frac{h}{2}, \mathbf{U}_{n} + \frac{h}{2}k_{1}\right).$$

$$\mathbf{U}_{n+1} = \mathbf{U}_{n} + hk_{2}.$$



Méthode de Simpson approchée : RK-4 (explicite)

$$\begin{aligned} \mathbf{k}_1 &= \mathbf{f} \left(t_n, \mathbf{U}_n \right). \\ \mathbf{k}_2 &= \mathbf{f} \left(t_n + \frac{h}{2}, \mathbf{U}_n + \frac{h}{2} k_1 \right). \\ \mathbf{k}_3 &= \mathbf{f} \left(t_n + \frac{h}{2}, \mathbf{U}_n + \frac{h}{2} k_2 \right). \\ \mathbf{k}_4 &= \mathbf{f} \left(t_n + h, \mathbf{U}_n + h \mathbf{k}_3 \right). \\ \mathbf{U}_{n+1} &= \mathbf{U}_n + \frac{h}{6} \left(k_1 + 2k_2 + 2k_3 + k_4 \right). \end{aligned}$$



Pourquoi les méthodes implicites?

On considère le problème de Cauchy :

$$\begin{cases} y'(t) = -1000y(t), \\ y(0) = 1. \end{cases}$$



Pourquoi les méthodes implicites?

On considère le problème de Cauchy :

$$\begin{cases} y'(t) = -1000y(t), \\ y(0) = 1. \end{cases}$$

- Problème raide : f(y) = -1000y est 1000-lipschitzienne.
- ▶ Solution exacte : $y(t) = e^{-1000t}$ (positive, converge vers 0 en $+\infty$).
- ► Constante d'erreur pour Euler proportionnelle à *Te*¹⁰⁰⁰*T*.



Pourquoi les méthodes implicites?

On considère le problème de Cauchy :

$$\begin{cases} y'(t) = -1000y(t), \\ y(0) = 1. \end{cases}$$

- Problème raide : f(y) = -1000y est 1000-lipschitzienne.
- ▶ Solution exacte : $y(t) = e^{-1000t}$ (positive, converge vers 0 en $+\infty$).
- Constante d'erreur pour Euler proportionnelle à Te^{1000T} .

Méthode d'Euler :

$$\left\{ \begin{aligned} Y_{n+1} &= Y_n - 1000hY_n, \\ Y_0 &= 1, \end{aligned} \right.$$
 soit $Y_n = (1-1000h)^n.$

Pourquoi les méthodes implicites?

On considère le problème de Cauchy :

$$\begin{cases} y'(t) = -1000y(t), \\ y(0) = 1. \end{cases}$$

- Problème raide : f(y) = -1000y est 1000-lipschitzienne.
- ▶ Solution exacte : $y(t) = e^{-1000t}$ (positive, converge vers 0 en $+\infty$).
- Constante d'erreur pour Euler proportionnelle à Te^{1000T} .

Méthode d'Euler :

$$\begin{cases} Y_{n+1} = Y_n - 1000hY_n, \\ Y_0 = 1, \end{cases}$$
 soit $Y_n = (1 - 1000h)^n$.

soit
$$Y_n = (1 - 1000h)^n$$
.

La solution numérique reste positive ssi $h \leq 10^{-3}$



Pourquoi les méthodes implicites?

On considère le problème de Cauchy :

$$\begin{cases} y'(t) = -1000y(t), \\ y(0) = 1. \end{cases}$$

- Problème raide : f(y) = -1000y est 1000-lipschitzienne.
- ▶ Solution exacte : $y(t) = e^{-1000t}$ (positive, converge vers 0 en $+\infty$).
- Constante d'erreur pour Euler proportionnelle à Te^{1000T} .

Méthode d'Euler :

$$\begin{cases} Y_{n+1} = Y_n - 1000hY_n, \\ Y_0 = 1, \end{cases}$$
 soit $Y_n = (1 - 1000h)^n$.

soit
$$Y_n = (1 - 1000h)^n$$
.

La solution numérique tend vers 0 ssi $h < 2 \cdot 10^{-3}$



Stabilité absolue/asymptotique

Euler implicite :
$$\begin{cases} Y_{n+1}=Y_n-1000hY_{n+1},\\ Y_0=1, \end{cases}$$
 soit $Y_n=(1+1000h)^{-n}.$

Stabilité absolue/asymptotique

Euler implicite :
$$\begin{cases} Y_{n+1}=Y_n-1000hY_{n+1},\\ Y_0=1, \end{cases}$$
 soit $Y_n=(1+1000h)^{-n}.$

La solution numérique reste positive et converge vers 0 pour tout h > 0.

Pas de contrainte sur h.



Stabilité absolue/asymptotique

Euler implicite

$$\begin{cases} Y_{n+1} = Y_n - 1000hY_{n+1}, \\ Y_0 = 1, \end{cases}$$
 soit $Y_n = (1 + 1000h)^{-n}$.

La solution numérique reste positive et converge vers $\mathbf{0}$ pour tout h > 0.

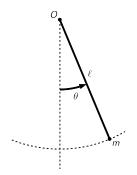
- Méthode implicite plus coûteuse (pour les problèmes non-linéaires).
- Mais utile pour les problèmes raides en temps long.



Conservation d'énergie

Modélisation du pendule sans frottement :

$$\left\{ \begin{array}{l} \theta^{\prime\prime}(t) = -\frac{\mathbf{g}}{\ell}\sin\left(\theta(t)\right) \\ \theta(0), \theta^\prime(0) \text{ donnés.} \end{array} \right.$$



Énergie du système :

$$E(t) = \frac{1}{2}\theta'(t)^2 - \frac{g}{\ell}\cos(\theta(t)).$$

Résultat : E'(t) = 0.

Question : est-ce respecté par les schémas ?





Méthode d'Euler explicite (ordre 1) pour h = 0.1.





Méthode d'Euler implicite (ordre 1) pour h = 0.1.





Méthode RK4 (ordre 4) pour h = 0.1.





Méthode d'Euler symplectique (ordre 1) pour h = 0.1.

