

Méthodes numériques

Résoudre les problèmes numériques rencontrés
en Physique, Chimie et SI

Germain Gondor

LYCÉE CARNOT - DIJON, 2021 - 2022

Sommaire

Première partie I

Problèmes stationnaires de dimension 1

Sommaire

- 1 Introduction
 - Objectifs
 - Support de l'étude
 - Linéarité
 - Non linéarité
- 2 Convergence
- 3 Rappel sur la méthode par dichotomie
- 4 Méthode de la corde (de Lagrange)
- 5 Méthode de Newton
- 6 Dilemme robustesse/rapidité
- 7 Résolution avec Scipy

Objectifs

L'objectif de cette partie est de traiter des problèmes stationnaires (constant au cours du temps), linéaires ou non, conduisant à la résolution approchée d'une équation algébrique ou transcendante.

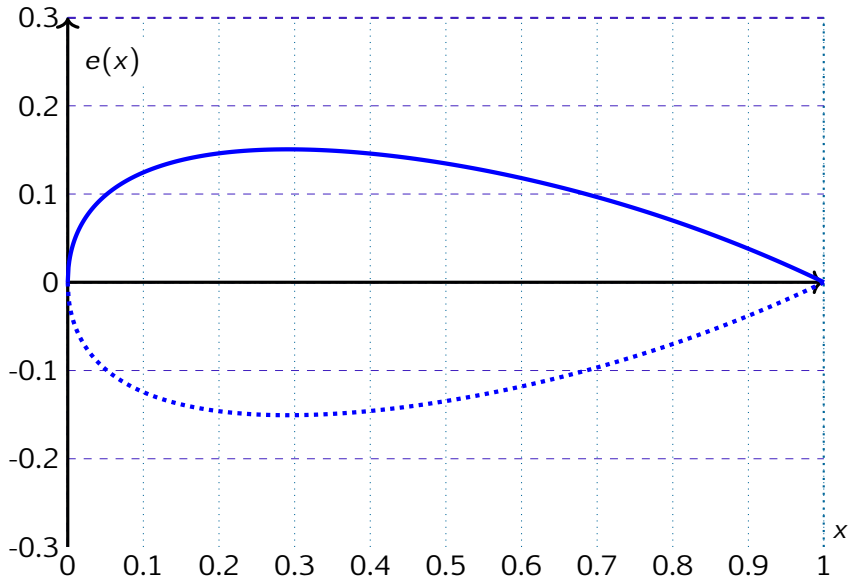
Support de l'étude

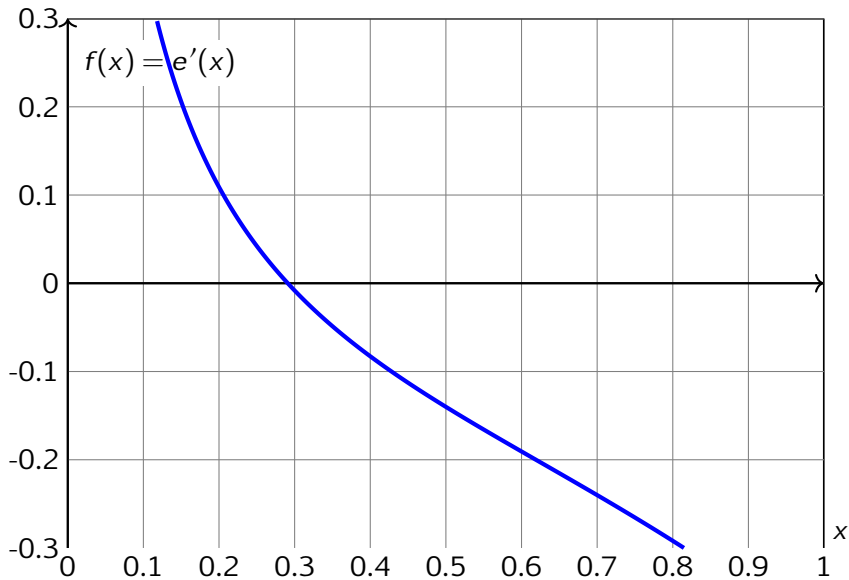
Comme support de cours, nous prenons un profil d'aile d'avion dont la demi-épaisseur $e(x)$ est approchée par l'équation:

$$e(x) = 0,15 \times (3,7 \times \sqrt{x} - 3,4 \times x - 0,3 \times x^4), \forall x \in [0, 1]$$

Rechercher la valeur de x pour laquelle le profil est le plus épais, revient à résoudre l'équation :

$$f(x) = 0,15 \times \left(\frac{3,7}{2 \cdot \sqrt{x}} - 3,4 - 1,2 \times x^3 \right) = 0$$





Linéarité

Soit f une fonction définie de E , un espace vectoriel, dans K , un corps commutatif. f est linéaire, si elle vérifie la propriété suivante:

$$\forall (x, y) \in E^2, \forall (\lambda, \mu) \in K^2, f(\lambda.x + \mu.y) = \lambda.f(x) + \mu.f(y)$$

EXEMPLE :

En 1D: Equation de droite

$$\begin{aligned} \mathbb{R} &\rightarrow \mathbb{R} \\ f: x &\mapsto a.x + b \end{aligned}$$

$$\text{avec } (a, b) \in \mathbb{R}^2$$

En 2D: Définition d'un plan

$$f(x, y, z) = a.x + b.y + c.z + d$$

La forme linéaire $f(x, y, z) = 0$ définit un plan où le vecteur (a, b, c) est un vecteur normal à ce plan.

Résoudre une équation linéaire 1D dans \mathbb{R} ou \mathbb{C} est immédiat:

$$\forall a \neq 0, a.x + b = 0 \Rightarrow x = -\frac{b}{a}$$

EXEMPLE : $3.i.x - 12.i + 6 = 0 \Rightarrow x = 4 + 2.i$

Résoudre un système linéaire de n équations à n inconnues suppose la résolution d'un système matriciel du type:

$$\begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} \Leftrightarrow [A] \cdot [X] = [B] \stackrel{\det(A) \neq 0}{\Rightarrow} [X] = [A^{-1}] \cdot [B]$$

Nous verrons dans un chapitre suivant comment le résoudre par un pivot de Gauss.

Non linéarité

Les cas non-linéaires sont plus difficiles à résoudre. Parfois une **résolution analytique** est possible comme dans la résolution d'une équation du second degré :

EXEMPLE :

$$x^2 - x - 6 = 0 \Rightarrow \Delta = (-1)^2 - 4 \times 1 \times (-6) = 25 = 5^2$$

$$\Rightarrow \begin{cases} x_+ &= \frac{-(-1) + 5}{2} = 2 \\ x_- &= \frac{-(-1) - 5}{2} = -3 \end{cases}$$

Dans d'autre cas, la solution analytique n'est pas connue (exemple : profil d'aile d'avion). On utilise alors une **approche numérique** pour la calculer. C'est l'objet de ce cours.

Sommaire

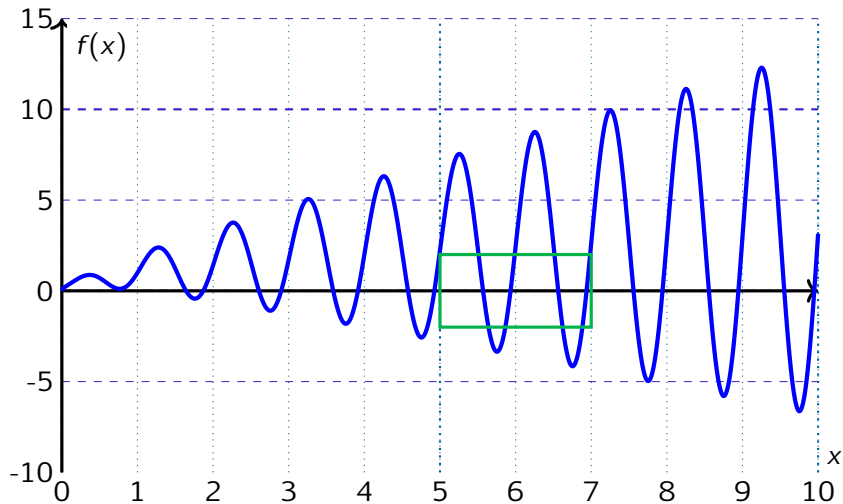
- 1 Introduction
- 2 **Convergence**
 - Séparation des racines
 - Représentations graphiques
 - Critères de convergence
 - Ordre de convergence
- 3 Rappel sur la méthode par dichotomie
- 4 Méthode de la corde (de Lagrange)
- 5 Méthode de Newton
- 6 Dilemme robustesse/rapidité
- 7 Résolution avec Scipy

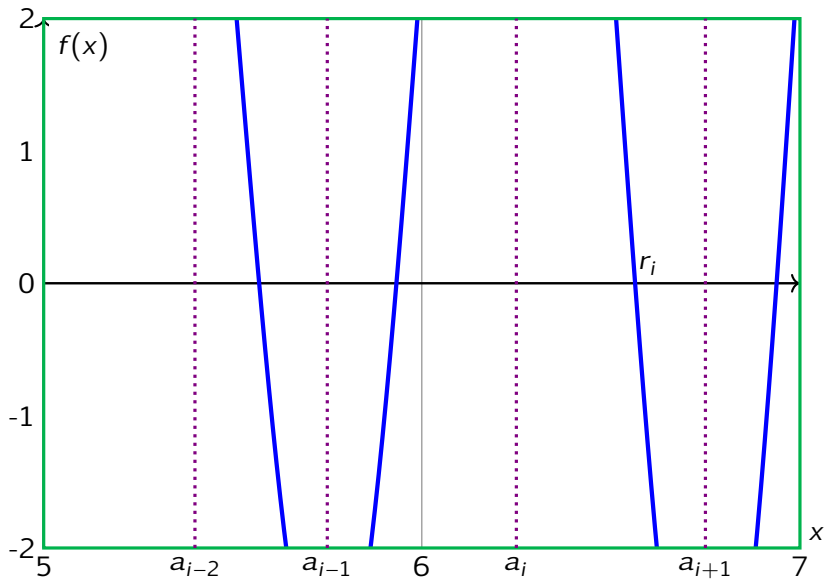
Séparation des racines

La plupart des méthodes de recherche de racines d'une fonction, se comportent bien si une seule racine r (valeur pour laquelle la fonction s'annule) est présente dans l'intervalle d'étude.

Séparer la racine r_i (ième annulation de la fonction sur l'intervalle de départ) revient à trouver un intervalle $]a_i, a_{i+1}[$ où cette racine est unique.

EXEMPLE : $f(x) = \sqrt{x} + x \cdot \sin(2\pi \cdot x)$





Il convient alors de faire un test aux bornes de l'intervalle d'étude pour éliminer les cas d'inadéquation de la comparaison à zéro, i.e, les cas sans racines sur l'intervalle. Il est en effet, possible d'élaborer une procédure permettant de déterminer la parité du nombre de racines sur l'intervalle :

- Posons $p_i = f(a_i).f(a_{i+1})$
- En tenant compte de l'ordre de multiplicité des racines, nous avons:
 - si $p_i < 0$, il existe $2.n + 1$ racines dans $]a_i, a_{i+1}[$, $n \in \mathbb{N}$
 - si $p_i > 0$, il existe $2.n$ racines dans $]a_i, a_{i+1}[$, $n \in \mathbb{N}$
 - si $p_i = 0$, alors a_i ou a_{i+1} est racine (voire les deux).

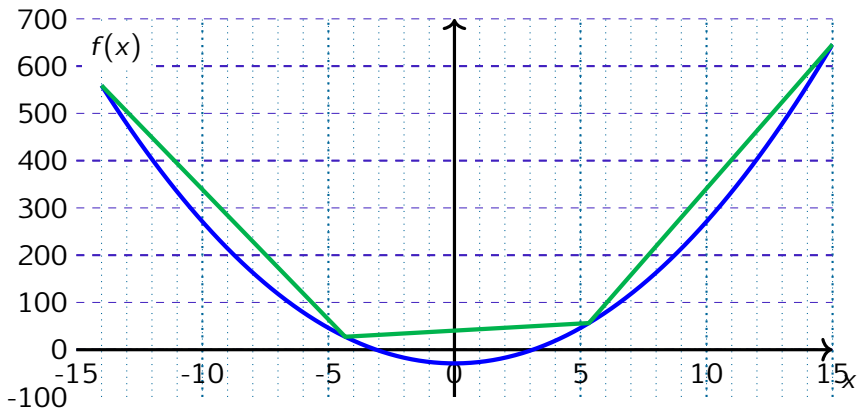
Evidemment, cet algorithme ne donne pas la valeur de n . Nous verrons plus loin l'utilité de ce test, notamment dans la recherche de la racine par dichotomie.

Représentations graphiques

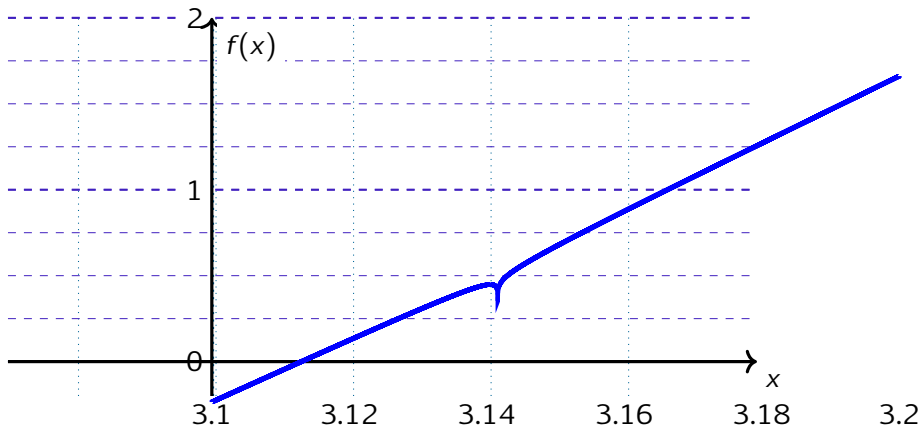
Les représentations graphiques sont utiles parfois pour avoir une idée des lieux de convergence. Cela permet d'obtenir à vu d'œil un premier intervalle ou un premier candidat pour démarrer une méthode numérique.

Cependant, il convient aussi de se méfier des représentations graphiques. Pour cela, il est préférable d'observer attentivement l'expression de la fonction et son domaine de définition.

EXEMPLE : $f(x) = 3.x^3 + \frac{1}{\pi^4} \cdot \ln[(\pi - x)^2] - 29$



Courbes tracées entre -14 et 15. La courbe verte n'ayant que 4 points, on ne constate pas d'intersection avec l'axe des abscisses.



Avec le zoom sur l'intervalle $[3, 1; 3, 2]$, il semble encore qu'il n'y ait qu'une racine sur l'intervalle. Or $\lim_{x \rightarrow \pi} f(x) = -\infty$.

Par continuité sur les intervalle $[3, 14; \pi[$ et $[\pi; 3, 15]$, la fonction f admet donc 3 racines dans l'intervalle $[3, 1; 3, 2]$.

Critères de convergence

Les méthodes numériques ne permettent pas d'obtenir une réponse formelle à un problème.

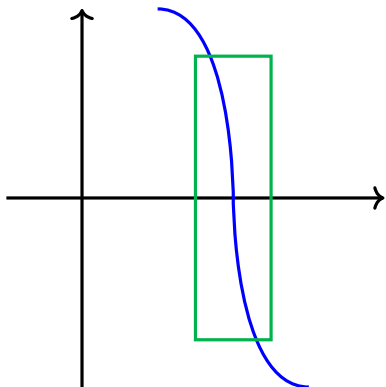
En revanche, elles permettent d'approcher la solution d'un problème avec une précision fonction du type de stockage des nombres (cf premier semestre) et du calcul numérique nécessaire pour obtenir cette solution (cf la sensibilité de certains calculs aux erreurs d'arrondis).

Pour une erreur normée en abscisse (fonction de l'"ordre de grandeur" de x), une tolérance de 10^{-8} est raisonnable. En revanche, chercher à obtenir 10^{-16} est a priori non nécessaire et impossible.

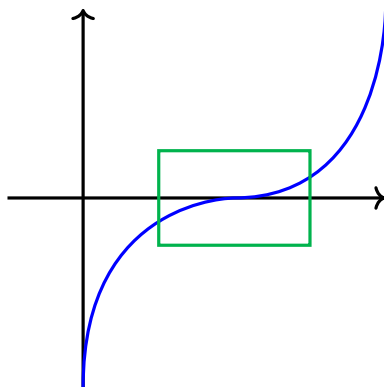
Dans le cas de la méthode par dichotomie, il est possible de faire un test de convergence sur la valeur de x et/ou sur la valeur de la fonction:

$$\frac{|b_n - a_n|}{|a_n| + |b_n|} < \varepsilon_x \quad \text{et/ou} \quad |f(c_n)| < \varepsilon_f$$

Suivant le gradient (pente) de la fonction, un des deux tests sera privilégié. Pour un fort (faible) gradient, il convient de suivre la valeur de f (de x).



Fort gradient



Faible gradient

En absence d'information, on peut suivre les deux (si la méthode donne une information sur l'erreur en x).

Ordre de convergence

Soit x_n la n ème valeur calculée pour résoudre l'équation $f(x) = 0$ sur l'intervalle donné. Soit r la solution exacte de cette équation. On peut alors définir ε_n comme l'erreur en abscisse à l'itération n .

La méthode numérique est qualifiée de méthode à l'ordre p si :

$$\exists K \in \mathbb{R}^+ / \lim_{n \rightarrow \infty} \frac{\varepsilon_{n+1}}{\varepsilon_n^p} = K$$

- Si $p = 1$ et $K = 1$ la méthode est dite sous-linéaire.
- Si $p = 1$ et $K \in]0, 1[$ la méthode est dite linéaire.
- Si $p > 1$, la méthode est qualifiée de super-linéaire.
- Si $p = 2$ la méthode est quadratique.

Sommaire

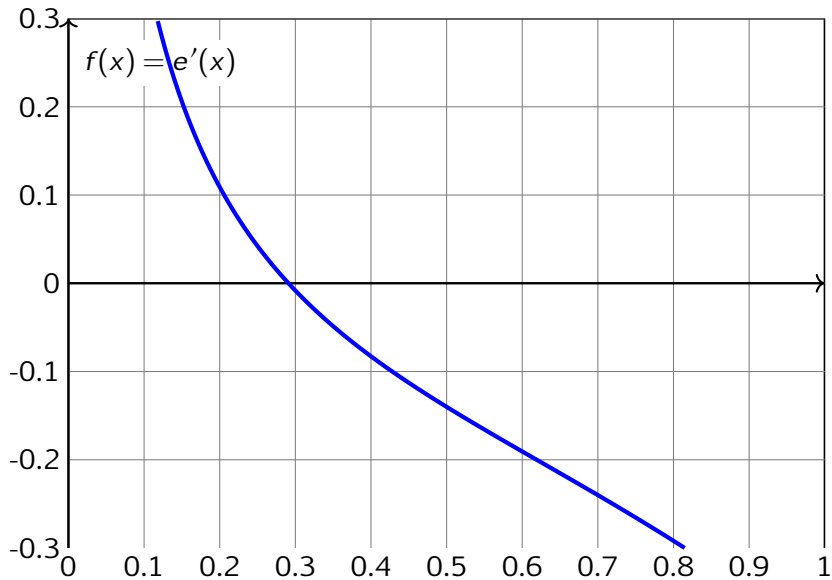
- 1 Introduction
- 2 Convergence
- 3 Rappel sur la méthode par dichotomie**
 - Principe
 - Algorithme
 - Convergence de la méthode
- 4 Méthode de la corde (de Lagrange)
- 5 Méthode de Newton
- 6 Dilemme robustesse/rapidité
- 7 Résolution avec Scipy

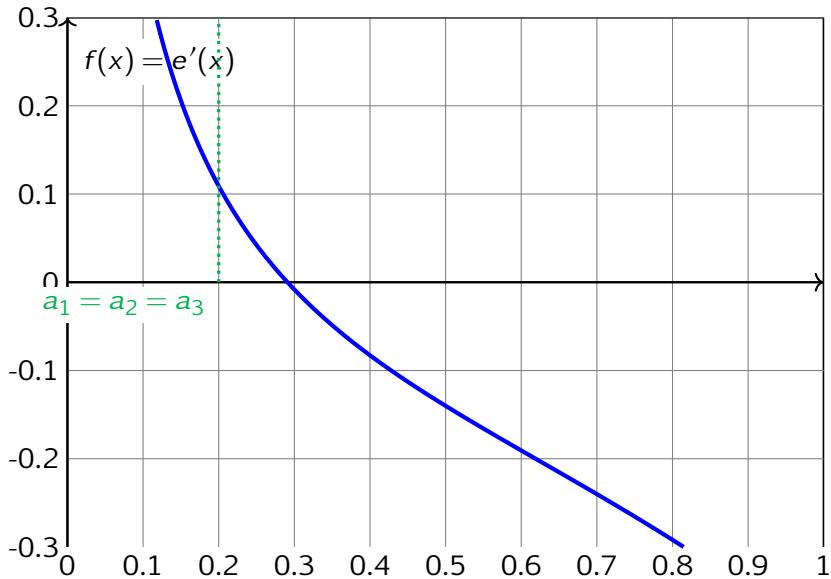
Principe

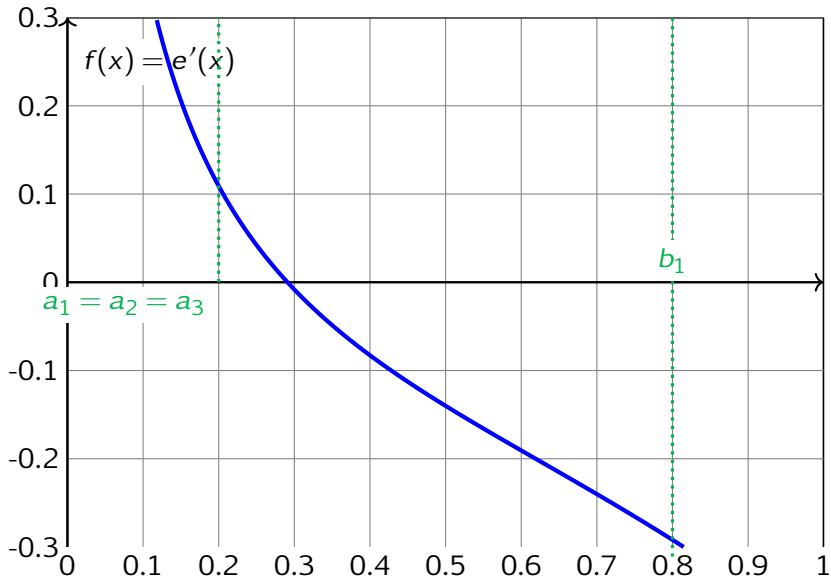
Le principe est de diviser l'intervalle en deux parts égales, de conserver la partie contenant la racine et d'y reproduire l'opération jusqu'à ce que le critère de convergence soit satisfait.

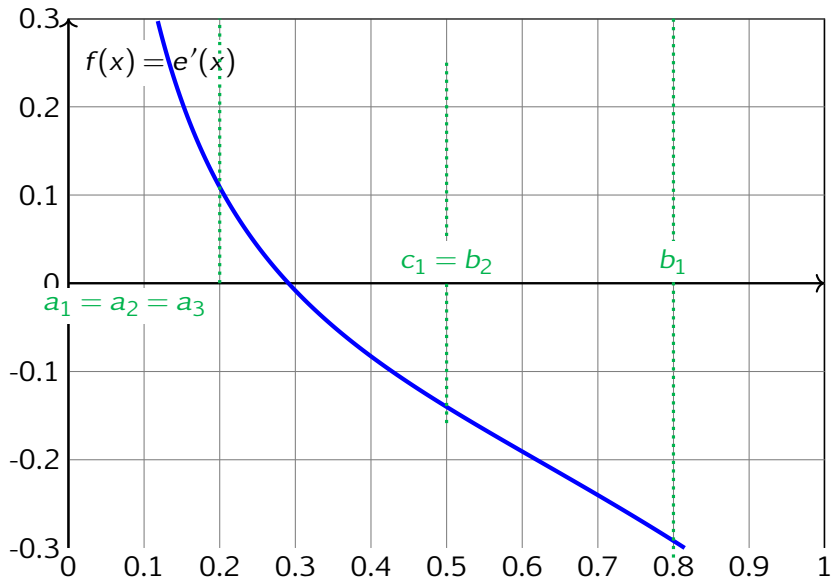
A partir d'un intervalle donné $[a, b]$, encadrant une racine de la fonction f étudiée :

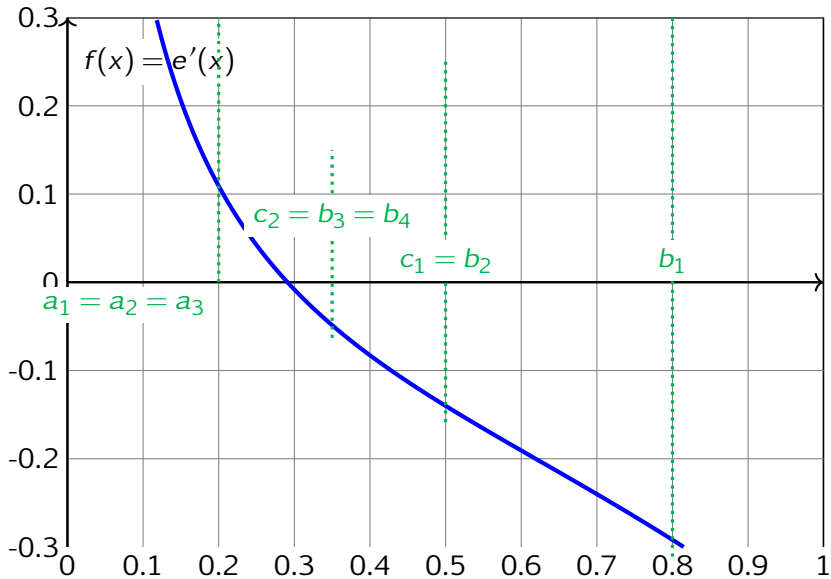
- Calculer le point c milieu de l'intervalle : $c = \frac{a+b}{2}$
- Evaluer $p = f(a).f(c)$ puis **test** :
 - si $p > 0$, il n'y a pas de racine dans l'intervalle $[a, c]$. La racine est donc dans l'intervalle $[c, b]$. On donne alors à a la valeur de c
 - si $p < 0$, la racine est dans $[a, c]$. On donne alors à b la valeur de c
 - si $p = 0$, alors la racine est c . Ô Miracle,...
- **test d'arrêt** : Evaluer le critère de convergence
- Recommencer si le critère de convergence n'est pas satisfait

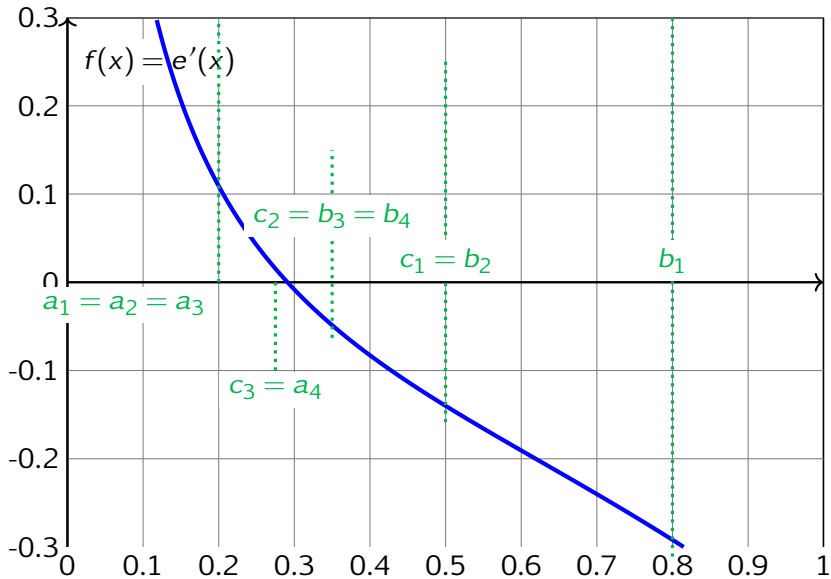


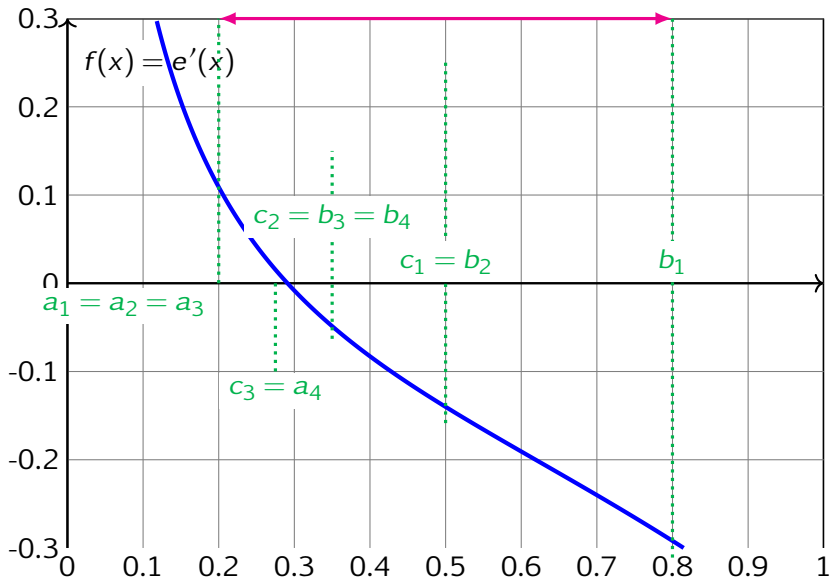


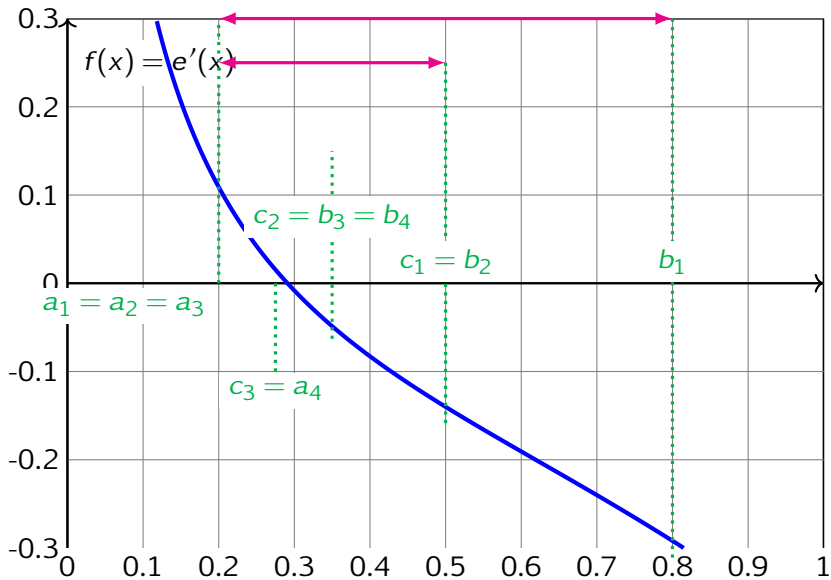


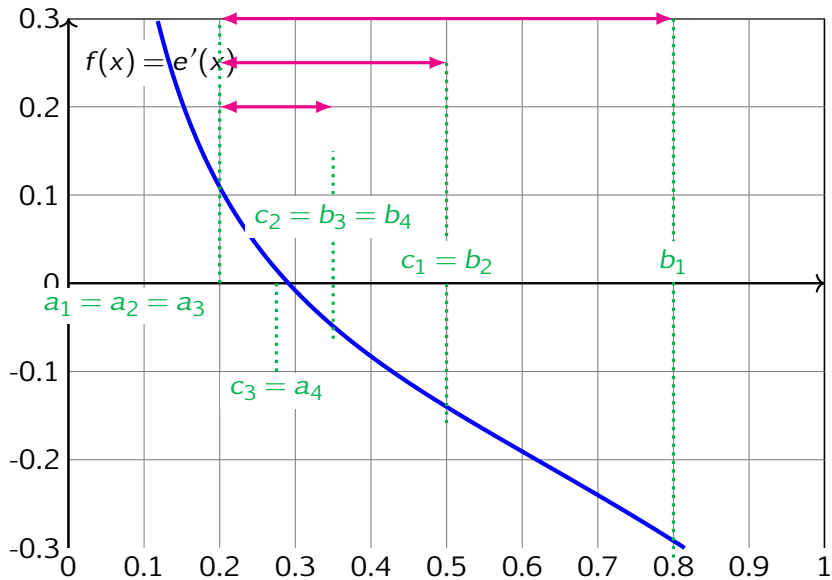


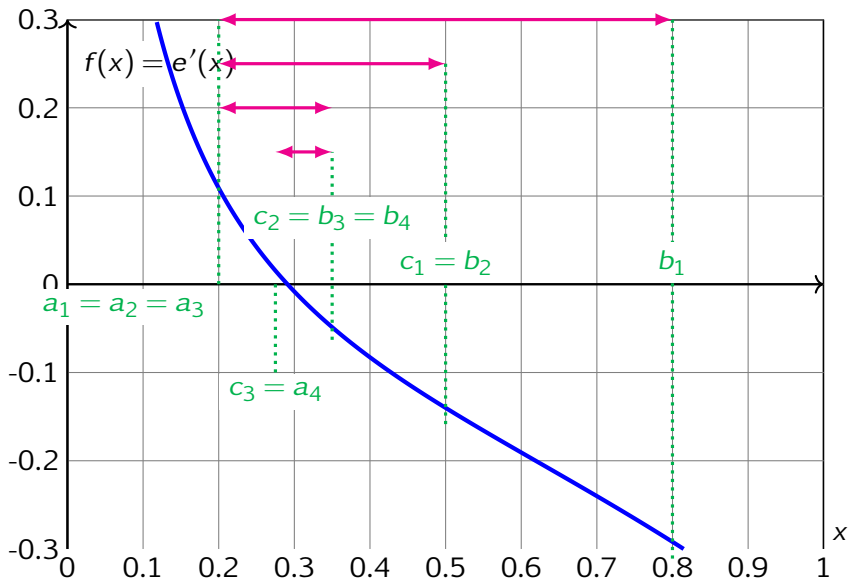












Algorithme

Méthode par dichotomie pour la fonction f sur $[a, b]$ avec une tolérance ε :

entrée: f , une fonction, a et b un encadrement de la racine cherchée et ε une tolérance

résultat: $(b_i - a_i)/2$ tel que $b_i - a_i \geq 2.\varepsilon$

Algorithm 1 Dichotomie

Dichotomie(f, a, b, ε)

$a_i, b_i \leftarrow a, b$

$f_a \leftarrow f(a_i)$

1: **tant que** $b_i - a_i > 2 \cdot \varepsilon$ **faire**

2: $c_i \leftarrow (a_i + b_i)/2$

3: $f_c \leftarrow f(c_i)$

4: **si** $f_a \cdot f_c \leq 0$ **alors**

5: $b_i \leftarrow c_i$

6: **sinon**

7: $a_i \leftarrow c_i$

8: $f_a \leftarrow f_c$

9: **fin si**

10: **fin tant que**

11: **renvoi:** $\frac{a_i + b_i}{2}$

Convergence de la méthode

À l'itération n , l'erreur ϵ_n entre la solution exacte et la solution numérique calculée peut être majorée par $|\epsilon_n| < \frac{|b_n - a_n|}{2}$. Ainsi, la racine r de la fonction peut être encadrée :

$$c_n - \frac{b_n - a_n}{2} \leq r \leq c_n + \frac{b_n - a_n}{2}$$

Soit $[a_n, b_n]$ l'intervalle à l'itération n . Définissons ϵ_n par $\epsilon_n = b_n - a_n$. Comme l'intervalle est divisé par deux à chaque itération, $|\epsilon_{n+1}| \leq \epsilon_{n+1} = \frac{\epsilon_n}{2} = \frac{\epsilon_1}{2^n}$. La convergence est donc linéaire.

Pour test d'arrêt, il est alors possible de prendre :

- soit un encadrement de la solution exacte : $b_n - a_n \leq \epsilon$
- soit la valeur de la fonction au point calculé : $|f(c_n)| < \epsilon$

Le nombre d'itération n pour obtenir un erreur inférieur à ε est tel que :

$$\varepsilon \leq \frac{b-a}{2^n} \Rightarrow 2^n \leq \frac{b-a}{\varepsilon}$$

$$\Rightarrow n \cdot \ln(2) \leq \ln\left(\frac{b-a}{\varepsilon}\right) \Rightarrow n \geq \frac{1}{\ln 2} \cdot \ln\left(\frac{b-a}{\varepsilon}\right)$$

Avec $\varepsilon = 10^{-p}$ (resp. $\varepsilon = 2^{-p}$) pour une approche décimal (resp. binaire) :

$$n \geq \frac{\ln(b-a) + p \cdot \ln(10)}{\ln 2} \quad \left(\text{resp. } n \geq p + \frac{\ln(b-a)}{\ln 2} \right)$$

La vitesse de convergence est lente mais la méthode est robuste. Si:

- la fonction f est définie et continue sur l'intervalle $[a, b]$
- la fonction n'admet qu'une seule racine sur l'intervalle $[a, b]$

alors la méthode converge.

Sommaire

- 1 Introduction
- 2 Convergence
- 3 Rappel sur la méthode par dichotomie
- 4 Méthode de la corde (de Lagrange)**
 - Principe
 - Convergence de la méthode
 - Algorithme
- 5 Méthode de Newton
- 6 Dilemme robustesse/rapidité
- 7 Résolution avec Scipy

Principe

Le principe de la méthode de la corde est de diviser l'intervalle en deux, de conserver la partie contenant la racine et d'y reproduire l'opération jusqu'à ce que le critère de convergence soit satisfait. Au lieu de diviser en deux parts égales, on cherche le point c (diviseur de l'intervalle intersection de l'axe des abscisses et de la droite passant par les points $A = (a, f(a))$ et $B = (b, f(b))$).

L'équation de la corde est donnée par:

$$y = \frac{f(b) - f(a)}{b - a} \cdot (x - a) + f(a) \quad \Rightarrow \quad 0 = \frac{f(b) - f(a)}{b - a} \cdot (c - a) + f(a)$$

Principe

Le principe de la méthode de la corde est de diviser l'intervalle en deux, de conserver la partie contenant la racine et d'y reproduire l'opération jusqu'à ce que le critère de convergence soit satisfait. Au lieu de diviser en deux parts égales, on cherche le point c (diviseur de l'intervalle intersection de l'axe des abscisses et de la droite passant par les points $A = (a, f(a))$ et $B = (b, f(b))$).

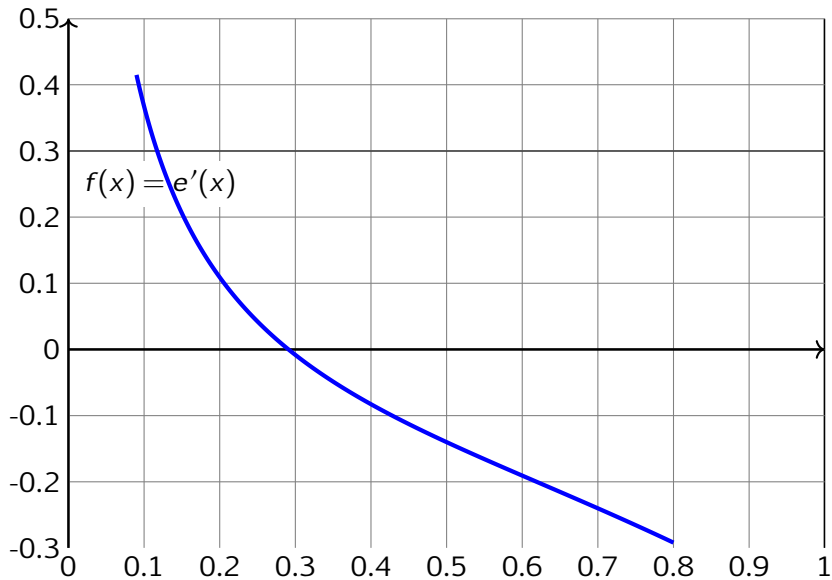
L'équation de la corde est donnée par:

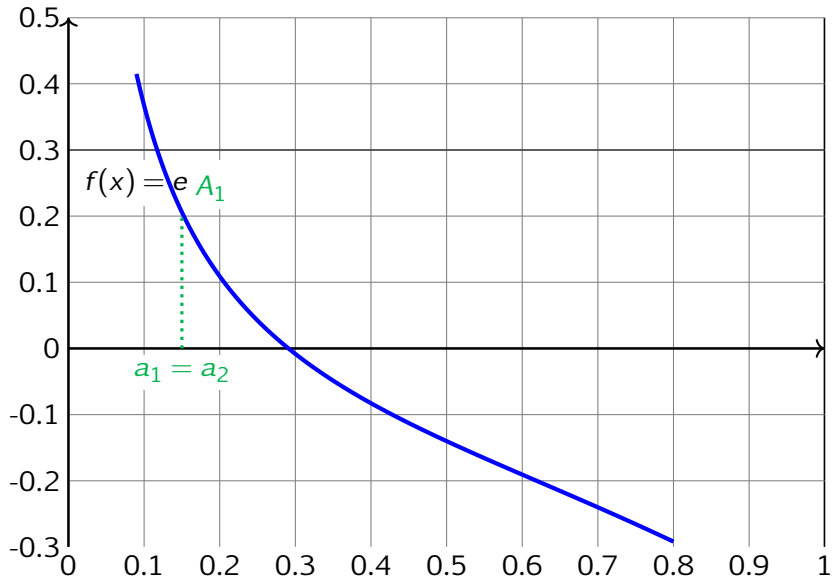
$$y = \frac{f(b) - f(a)}{b - a} \cdot (x - a) + f(a) \quad \Rightarrow \quad 0 = \frac{f(b) - f(a)}{b - a} \cdot (c - a) + f(a)$$

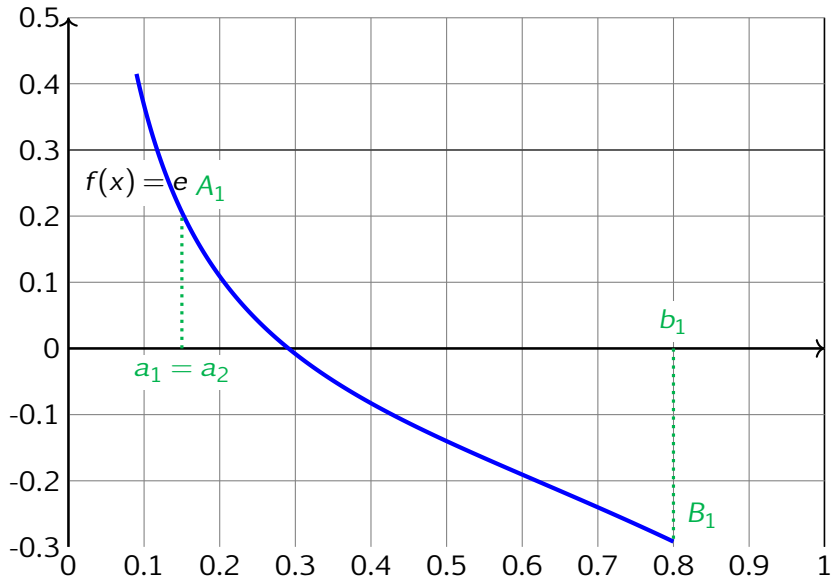
$$\Rightarrow \quad c = \frac{a \cdot f(b) - b \cdot f(a)}{f(b) - f(a)}$$

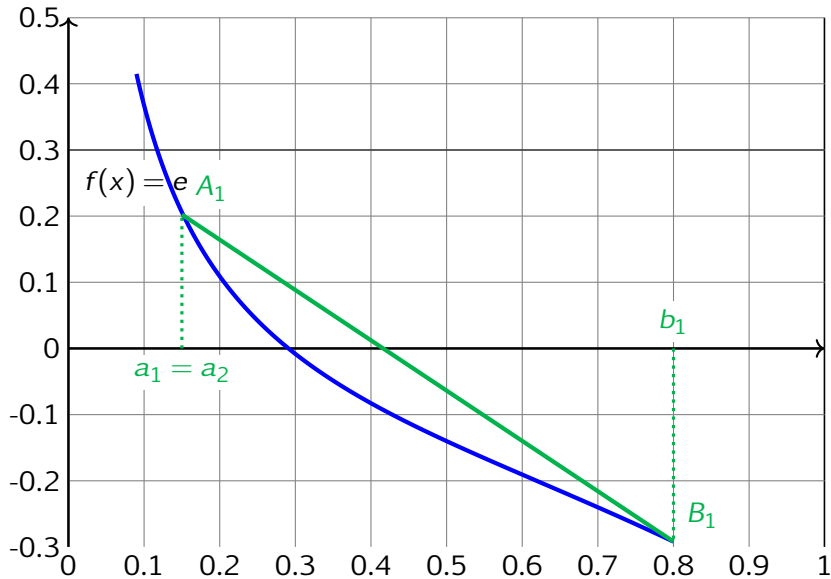
A partir d'un intervalle donné $[a, b]$, encadrant une racine de la fonction f étudiée :

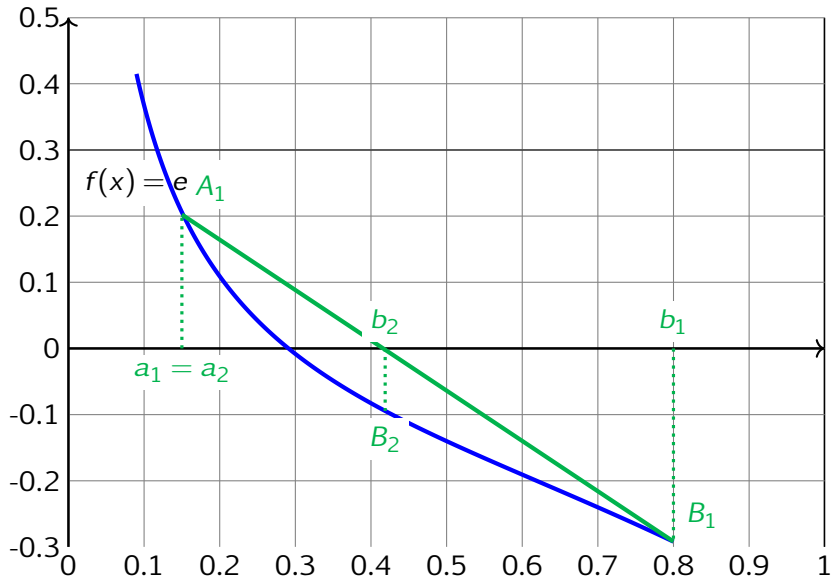
- Calculer le point c sur la corde et sur l'axe des abscisses : $c = \frac{a.f(b) - b.f(a)}{f(b) - f(a)}$
- Evaluer $p = f(a).f(c)$ puis **test** :
 - si $p > 0$, il n'y a pas de racine dans l'intervalle $[a, c]$. La racine est donc dans l'intervalle $[c, b]$. On donne alors à a la valeur de c
 - si $p < 0$, la racine est dans $[a, c]$. On donne alors à b la valeur de c
 - si $p = 0$, alors la racine est c . Ô Miracle,...
- **test d'arrêt** : Evaluer le critère de convergence
- Recommencer si le critère de convergence n'est pas satisfait

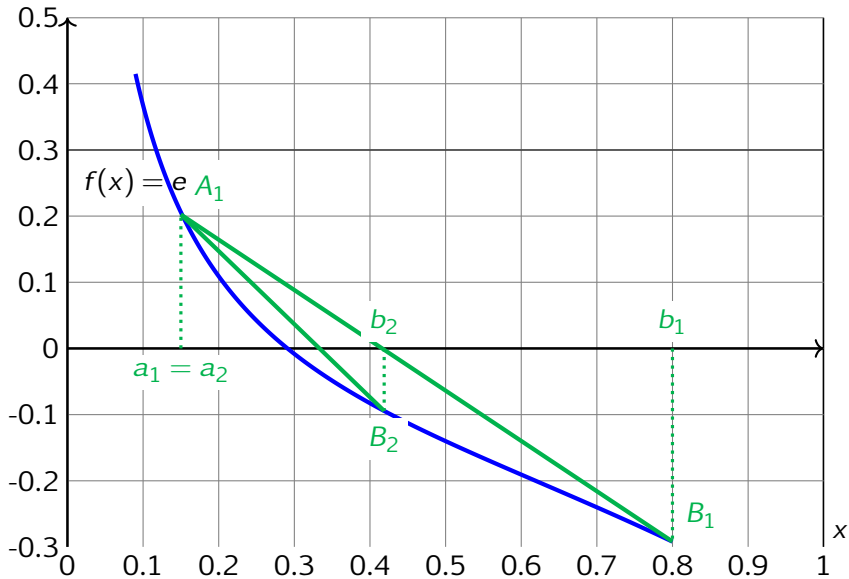












Convergence de la méthode

Comme nous pouvons le voir sur la figure précédente, dans le cas d'une fonction convexe sur $[a, b]$ telle que $f(a) > 0$ et $f(b) < 0, \forall n \in \mathbb{N}, a_n = a_1$.

Il ne sera donc pas possible d'encadrer la solution exacte par intervalle dont la longueur tendrait vers 0 avec n . En effet :

A l'itération n , r étant la solution exacte du problème $f(x) = 0$ sur l'intervalle d'étude :

$$c_n = \frac{a_n \cdot f(b_n) - b_n \cdot f(a_n)}{f(b_n) - f(a_n)} \Rightarrow \begin{cases} a_n &= c_n + \frac{b_n - a_n}{f(b_n) - f(a_n)} \cdot f(a_n) \\ b_n &= c_n + \frac{b_n - a_n}{f(b_n) - f(a_n)} \cdot f(b_n) \end{cases}$$

or $a_n < r < b_n$

$$\Rightarrow c_n + \frac{b_n - a_n}{f(b_n) - f(a_n)} \cdot f(a_n) < r < c_n + \frac{b_n - a_n}{f(b_n) - f(a_n)} \cdot f(b_n)$$

$$\text{d'où } |c_n - r| < \frac{\max(|f(a_n)|, |f(b_n)|)}{|f(b_n) - f(a_n)|} \cdot |b_n - a_n|$$

Il convient alors de prendre pour test d'arrêt:

- la valeur de la fonction au point calculé : $|f(c_n)| < \varepsilon$
- l'évolution de l'algorithme : $|c_n - c_{n-1}| < \varepsilon$

Algorithme

Algorithm 2 Méthode de la corde de Lagrange

```
1:  $a_i, b_i \leftarrow a, b$ 
2:  $f_a, f_b \leftarrow f(a_i), f(b_i)$ 
3:  $f_c \leftarrow f_a$ 
4: tant que  $|f_c| > \varepsilon$  faire
5:    $c_i \leftarrow (a_i.f_b - b_i.f_a)/(f_b - f_a)$ 
6:    $f_c \leftarrow f(c_i)$ 
7:   si  $f_a.f_c \geq 0$  alors
8:      $b_i \leftarrow c_i$ 
9:      $f_b \leftarrow f_c$ 
10:  sinon  $a_i \leftarrow c_i$   $f_a \leftarrow f_c$ 
11:  fin si
12: fin tant que
13: renvoi:  $\frac{a_i.f_b - b_i.f_a}{f_b - f_a}$ 
```

Sommaire

- 1 Introduction
- 2 Convergence
- 3 Rappel sur la méthode par dichotomie
- 4 Méthode de la corde (de Lagrange)
- 5 Méthode de Newton**
 - Principe
 - Convergence de la méthode
 - Algorithme
 - Limites et précautions
 - Fausse position
 - Fausse position
 - Méthode de Newton-Raphson

Principe

Les méthodes de Schröder sont basés sur les n premières dérivées d'une fonction f pour une méthode d'ordre $n + 1$. Dans le cas de la méthode de Newton, on n'utilise que la dérivée première.

RAPPEL Si la fonction est de classe C^1 sur l'intervalle $[a, b]$ alors, le développement de Taylor à l'ordre 1 donne:

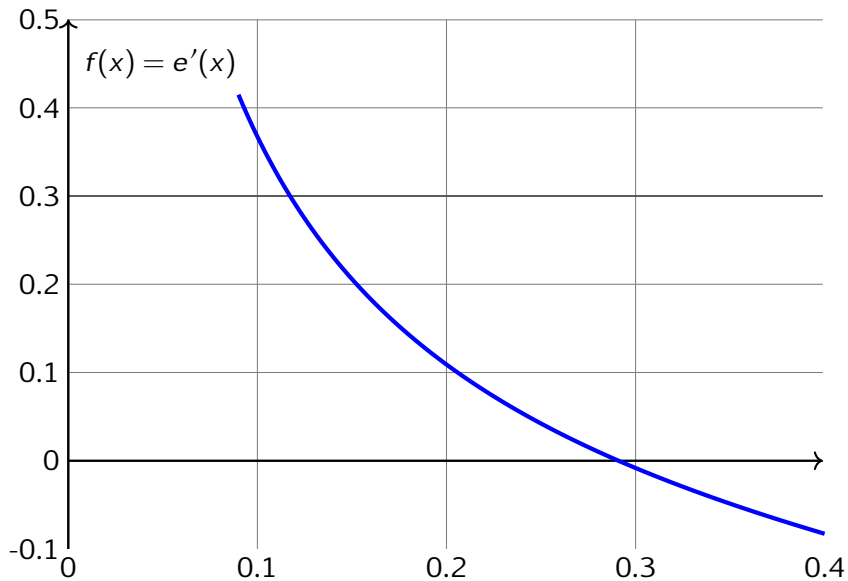
$$f(b) = f(a) + f'(a).(b - a) + o(b - a)$$

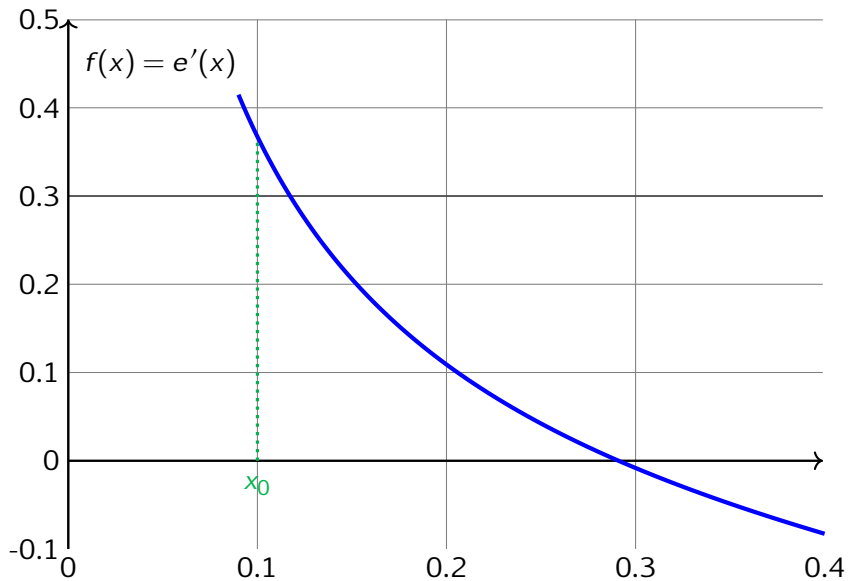
A partir d'un point x_0 , la méthode consiste à rechercher le point suivant x_1 en le supposant racine de la fonction et en négligeant le terme $o(b - a)$ dans le développement de Taylor à l'ordre 1. Ainsi:

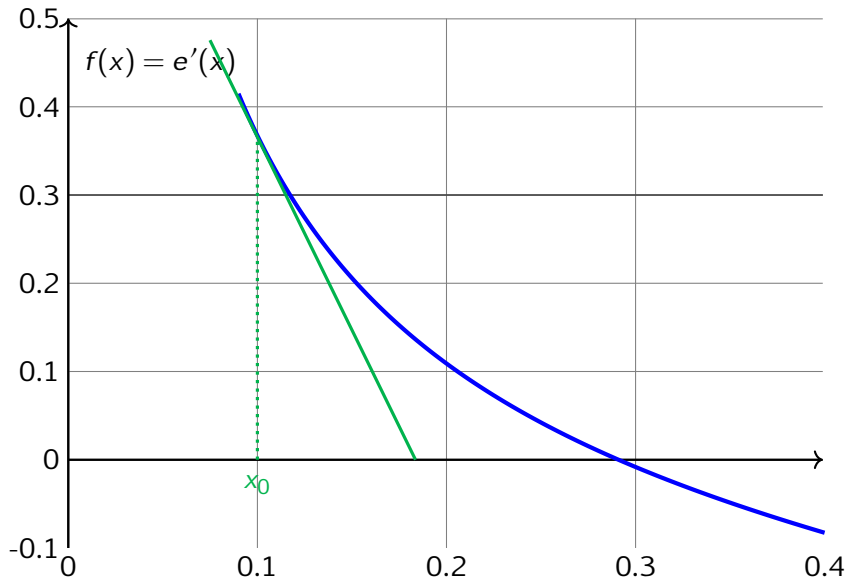
$$0 = f(x_1) = f(x_0) + f'(x_0).(x_1 - x_0) \Rightarrow x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

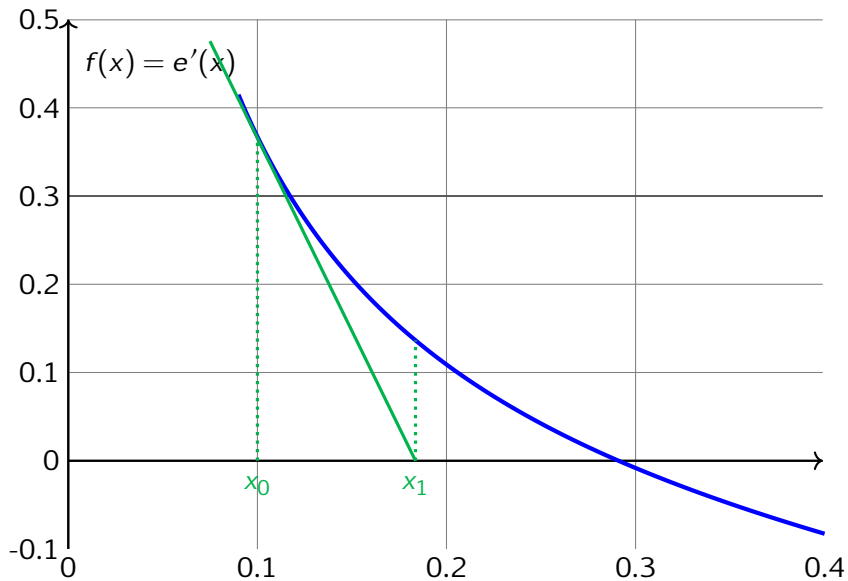
Soit f de classe C^1 sur $I = [a, b]$:

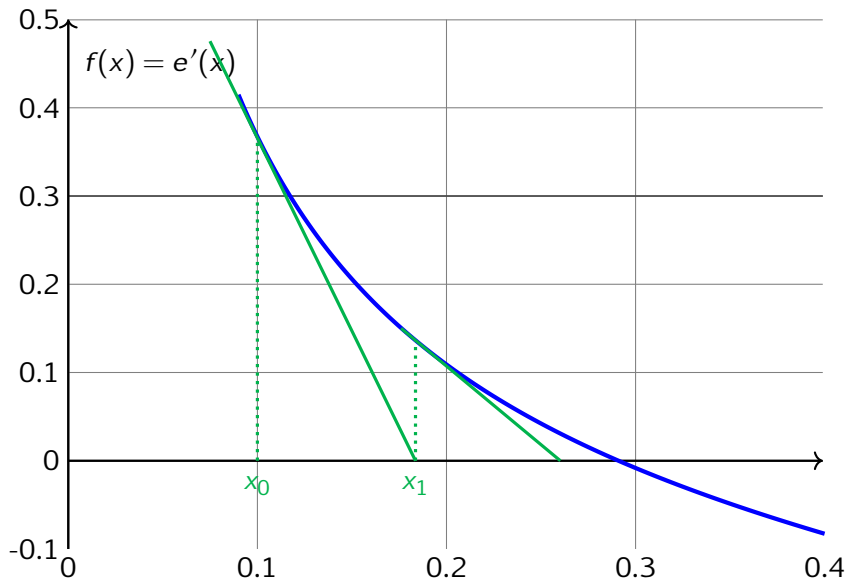
- Choisir un premier candidat x_0
- **Test d'arrêt** : Tester le critère de convergence $|f(x_0)| < \epsilon$ et $x_0 \notin I$
- Si le critère n'est pas atteint, calculer le nouveau candidat $x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$
- Reprendre depuis le test d'arrêt

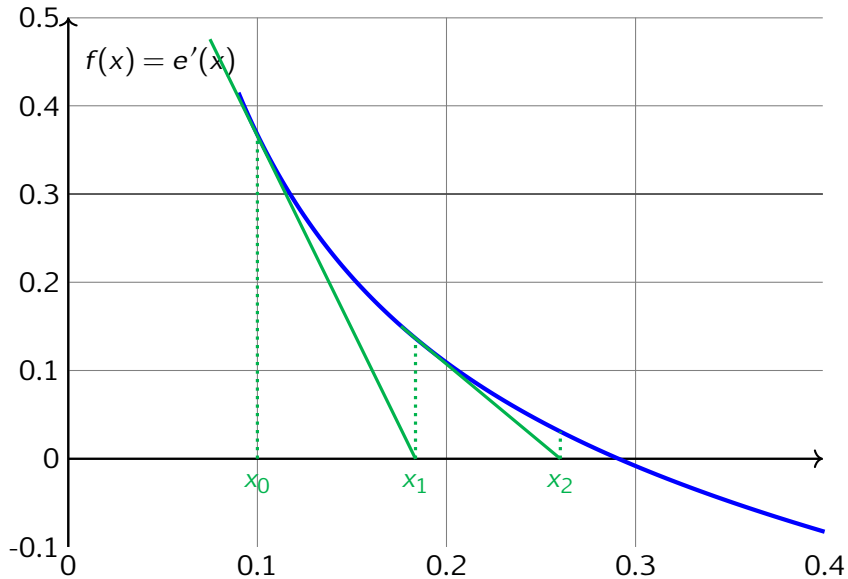


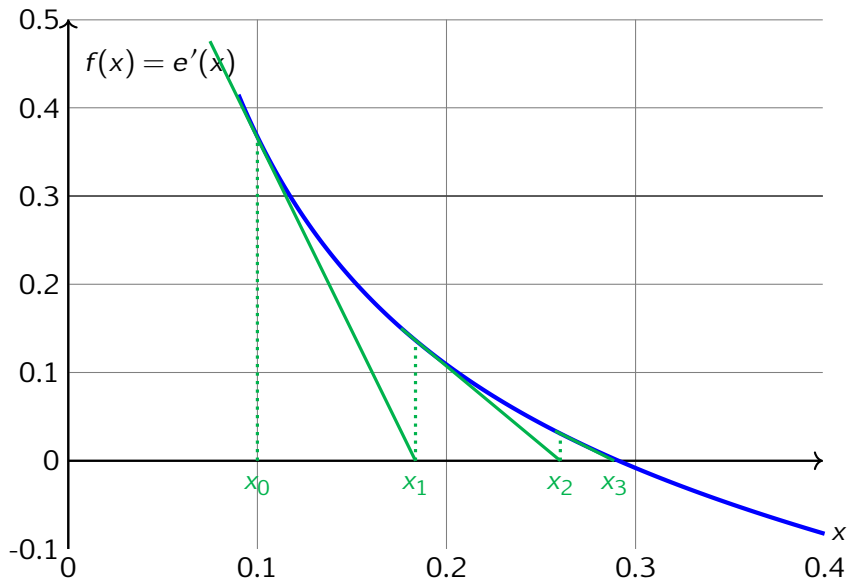












Convergence de la méthode

Soit r la racine de la fonction f sur l'intervalle d'étude et x_n , la n ème valeur calculée par itération de Newton. Notons alors $\varepsilon_n = x_n - r$. Le développement de Taylor à l'ordre 2 au voisinage de r donne:

$$f(x_n) = \cancel{f(r)} + f'(r) \cdot \varepsilon_n + f''(r) \cdot \frac{\varepsilon_n^2}{2} + o(\varepsilon_n^2)$$

$$f'(x_n) = f'(r) + f''(r) \cdot \varepsilon_n + o(\varepsilon_n)$$

$$\begin{aligned} \text{alors } \varepsilon_{n+1} &= x_{n+1} - r = x_n - r - \frac{f(x_n)}{f'(x_n)} \\ &= \varepsilon_n - \frac{f'(r) \cdot \varepsilon_n + f''(r) \cdot \frac{\varepsilon_n^2}{2} + o(\varepsilon_n^2)}{f'(r) + f''(r) \cdot \varepsilon_n + o(\varepsilon_n)} = \frac{\varepsilon_n^2}{2} \cdot \frac{f''(r)}{f'(r)} + o(\varepsilon_n^2) \end{aligned}$$

La convergence est donc quadratique.

Algorithme

Méthode de Newton pour la fonction f sur $[a, b]$ avec une tolérance ε :

Algorithm 3 Méthode de Newton

```
1:  $x_i \leftarrow x_0$ 
2:  $f_x \leftarrow f(x_i)$ 
3: tant que  $|f_x| > \varepsilon$  faire
4:    $fp_x \leftarrow f'(x_i)$ 
5:   si  $fp_x = 0$  alors
6:     break
7:   fin si
8:    $x_i \leftarrow x_i - f_x / fp_x$ 
9:    $f_x \leftarrow f(x_i)$ 
10: fin tant que
11: renvoi:  $x_i$ 
```

Limites et précautions

La méthode est bien adapté si la valeur de x_0 est proche de la racine r de la fonction f et si la dérivée n'est pas trop faible (pente trop horizontale).

Limites et précautions

La méthode est bien adaptée si la valeur de x_0 est proche de la racine r de la fonction f et si la dérivée n'est pas trop faible (pente trop horizontale).

Soit la fonction $f(x) = x - 2.\sin(x)$. La dérivée $f'(x) = 1 - 2.\cos(x)$ s'annule pour $x = \frac{\pi}{3} \approx 1,0472$.

Limites et précautions

La méthode est bien adaptée si la valeur de x_0 est proche de la racine r de la fonction f et si la dérivée n'est pas trop faible (pente trop horizontale).

Soit la fonction $f(x) = x - 2.\sin(x)$. La dérivée $f'(x) = 1 - 2.\cos(x)$ s'annule pour $x = \frac{\pi}{3} \approx 1,0472$.

Prendre comme valeur $x_0 = \frac{\pi}{3}$ conduit tout de suite au crash. La dérivée étant nulle, x_1 ne pourra pas être calculé. Prendre $0 < x_0 < \frac{\pi}{3}$ conduit à la racine $r_0 = 0$. Pour trouver la deuxième racine $r_1 \approx 1,9$, il convient de prendre une valeur de départ $x_0 > \frac{\pi}{3}$.

n	x_n	$f(x_n)$
0	1,1	-0,6824147
1	8,452992	6,801205
2	5,256414	6,967679
3	203,384184	201,922795
\vdots	\vdots	\vdots
57	-0,380713	0,362452
58	0,042317	-0,042292
59	-0,000051	0,000051
60	0,000000	-0,000000

Cependant, prendre pour valeur $x_0 = 1,1$ ou $x_0 = 1,2$ conduit dans le premier cas à r_0 et dans le deuxième à r_1 , même si la fonction est monotone sur $\left[\frac{\pi}{3}, \pi\right]$ et que $\frac{\pi}{3} < 1,1 < 1,2 < \pi$.

Notons dans les deux cas, la convergence quadratique à proximité de r .

n	x_n	$f(x_n)$
0	1,2	-0,6640782
1	3,612334	4,519429
2	1,988080	0,159694
3	1,899879	0,007200
4	1,895505	0,000018
5	1,895494	0,000000

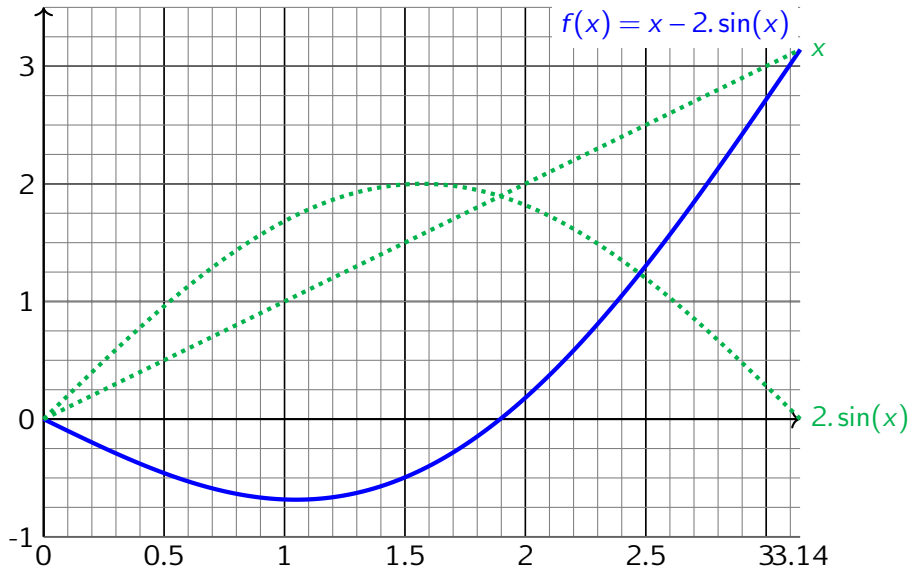
n	x_n	$f(x_n)$
0	1,2	-0,6640782
1	3,612334	4,519429
2	1,988080	0,159694
3	1,899879	0,007200
4	1,895505	0,000018
5	1,895494	0,000000

Il est donc préférable que la fonction soit monotone sur l'intervalle d'étude et que la première valeur x_0 soit proche de r .

n	x_n	$f(x_n)$
0	1,2	-0,6640782
1	3,612334	4,519429
2	1,988080	0,159694
3	1,899879	0,007200
4	1,895505	0,000018
5	1,895494	0,000000

Il est donc préférable que la fonction soit monotone sur l'intervalle d'étude et que la première valeur x_0 soit proche de r .

Lorsque l'algorithme patine ou que l'itération conduit à une valeur hors du domaine de définition, il est possible d'utiliser une méthode par dichotomie et reprendre le schéma de Newton près de la racine.



Fausse position

Principe

La méthode de Newton étant basée sur un développement de Taylor à l'ordre 1, il est nécessaire de calculer la dérivée première.

La méthode de la fausse position est basée sur une estimation de la dérivée à partir des deux points précédents. Ce n'est pas la valeur exacte comme dans le schéma de Newton. Nous avons alors :

$$f'(x_n) \approx \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}$$

$$\begin{aligned}\Rightarrow x_{n+1} &= x_n - \frac{f(x_n)}{f'(x_n)} = x_n - f(x_n) \cdot \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} \\ &= \frac{x_{n-1} \cdot f(x_n) - x_n \cdot f(x_{n-1})}{f(x_n) - f(x_{n-1})}\end{aligned}$$

Il est donc préférable que la fonction soit monotone sur l'intervalle d'étude et que la première valeur x_0 soit proche de r .

Lorsque l'algorithme patine ou que l'itération conduit à une valeur hors du domaine de définition, il est possible d'utiliser une méthode par dichotomie et reprendre le schéma de Newton près de la racine.

Soit f définie sur $I = [a, b]$:

- Choisir un premier couple candidat $(x_0, x_1) \in I^2$
- **Test d'arrêt** : Tester le critère de convergence $|f(x_1)| < \varepsilon$ et $x_1 \notin I$
- Si le critère n'est pas atteint :
 - calculer $f(x_1)$
 - si $f(x_1) = f(x_0)$ soit **arrêt**, soit choisir une autre valeur de x_0
 - calculer le nouveau candidat $x_2 = x_1 - f(x_1) \cdot \frac{x_1 - x_0}{f(x_1) - f(x_0)}$
- Reprendre depuis le test d'arrêt avec : si $f(x_n) = f(x_{n-1})$ soit **arrêt**, soit prendre x_{n_2} à la place de x_{n_1}

Fausse position

Convergence de la méthode

Soit r la racine de la fonction f sur l'intervalle d'étude et x_n , la n ème valeur calculée par la méthode de la fausse position. Notons alors $\varepsilon_n = x_n - r$. Le développement de Taylor à l'ordre 2 au voisinage de r donne:

$$f(x_n) = \cancel{f(r)} + f'(r) \cdot \varepsilon_n + f''(r) \cdot \frac{\varepsilon_n^2}{2} + o(\varepsilon_n^2)$$

$$f(x_{n-1}) = \cancel{f(r)} + f'(r) \cdot \varepsilon_{n-1} + f''(r) \cdot \frac{\varepsilon_{n-1}^2}{2} + o(\varepsilon_{n-1}^2)$$

alors

$$\begin{aligned}
 \varepsilon_{n+1} &= x_{n+1} - r \\
 &= \frac{x_{n-1} \cdot f(x_n) - x_n \cdot f(x_{n-1})}{f(x_n) - f(x_{n-1})} - r \\
 &= \frac{\varepsilon_{n-1} \cdot f(x_n) - \varepsilon_n \cdot f(x_{n-1})}{f(x_n) - f(x_{n-1})} \\
 &= \frac{f''(r) \cdot \left(\frac{\varepsilon_n^2}{2} \cdot \varepsilon_{n-1} - \frac{\varepsilon_{n-1}^2}{2} \cdot \varepsilon_n \right) + o(\varepsilon_n^2) - o(\varepsilon_{n-1}^2)}{f'(r) \cdot (\varepsilon_n - \varepsilon_{n-1}) + o(\varepsilon_n) - o(\varepsilon_{n-1})} \\
 &\approx \frac{f''(r) \cdot \varepsilon_n \cdot \varepsilon_{n-1} \cdot \left(\frac{\varepsilon_n}{2} - \frac{\varepsilon_{n-1}}{2} \right)}{f'(r) \cdot (\varepsilon_n - \varepsilon_{n-1})} = \alpha \cdot \varepsilon_n \cdot \varepsilon_{n-1}
 \end{aligned}$$

Nous cherchons à connaître l'ordre p de convergence, i.e, p tel que $\varepsilon_{n+1} = \varepsilon_n^p$. Or:

$$\varepsilon_{n+1} \equiv \varepsilon_n \cdot \varepsilon_{n-1} \Rightarrow \varepsilon_n^p \equiv \varepsilon_n \cdot \varepsilon_n^{\frac{1}{p}} \Rightarrow p = 1 + \frac{1}{p}$$

$$\Rightarrow p^2 - p - 1 = 0 \Rightarrow p = \frac{1 + \sqrt{5}}{2} \approx 1,62 < 2$$

La convergence est donc super-linéaire ($p > 1$) mais inférieure à la convergence quadratique ($p = 2$) de la méthode de Newton. En revanche, elle ne nécessite pas le calcul de la dérivée première qui peut parfois être très lourd.

Méthode de Newton-Raphson

La méthode de Newton-Raphson est une généralisation de la méthode de Newton pour les fonctions de plusieurs variables. Il s'agit de résoudre un système de n équations à n inconnues :

$$\vec{F}(\vec{X}) = \vec{0} \Leftrightarrow \begin{bmatrix} F_1(x_1, \dots, x_i, \dots, x_n) \\ \vdots \\ F_i(x_1, \dots, x_i, \dots, x_n) \\ \vdots \\ F_n(x_1, \dots, x_i, \dots, x_n) \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

La formule de Taylor à l'ordre 1 donne pour chaque F_i :

$$F_i(\vec{x} + \vec{\delta x}) = F_i(\vec{x}) + \sum_{j=1}^n \frac{\partial F_i(\vec{x})}{\partial x_j} \cdot \delta x_j + o(\vec{\delta x}), \forall i \in \llbracket 1, n \rrbracket$$

En introduisant alors la matrice jacobienne \mathbb{J} définie par

$$[\mathbb{J}]_{(i,j)} = J_{i,j} = \frac{\partial F_i(\vec{x})}{\partial x_j} \quad \Leftrightarrow \quad \mathbb{J} = \begin{bmatrix} \frac{\partial F_1(\vec{x})}{\partial x_1} & \cdots & \frac{\partial F_1(\vec{x})}{\partial x_j} & \cdots & \frac{\partial F_1(\vec{x})}{\partial x_n} \\ \vdots & \ddots & \vdots & & \vdots \\ \frac{\partial F_i(\vec{x})}{\partial x_1} & \cdots & \frac{\partial F_i(\vec{x})}{\partial x_j} & \cdots & \frac{\partial F_i(\vec{x})}{\partial x_n} \\ \vdots & & \vdots & \ddots & \vdots \\ \frac{\partial F_n(\vec{x})}{\partial x_1} & \cdots & \frac{\partial F_n(\vec{x})}{\partial x_j} & \cdots & \frac{\partial F_n(\vec{x})}{\partial x_n} \end{bmatrix}$$

la formule générale du développement limité à l'ordre 1 de \vec{F} au voisinage de \vec{x} devient:

$$\vec{F}(\vec{x} + \delta\vec{x}) = \vec{F}(\vec{x}) + \mathbb{J}(\vec{x}) \cdot \delta\vec{x} + o(\delta\vec{x})$$

L'algorithme de Newton - Raphson s'obtient en cherchant $\overrightarrow{X_{n+1}}$ solution du problème $\vec{F}(\vec{x}) = \vec{0}$. Le développement de Taylor, conduit à :

$$\vec{F}(\overrightarrow{X_{n+1}}) = \vec{F}(\overrightarrow{X_n}) + \mathbb{J}(\overrightarrow{X_n}) \cdot (\overrightarrow{X_{n+1}} - \overrightarrow{X_n}) + o(\overrightarrow{X_{n+1}} - \overrightarrow{X_n})$$

$$\text{d'où } \overrightarrow{X_{n+1}} = \overrightarrow{X_n} - \mathbb{J}^{-1}(\overrightarrow{X_n}) \cdot \vec{F}(\overrightarrow{X_n})$$

Chaque itération demande alors le calcul de l'inverse de la matrice Jacobienne...

Sommaire

- 1 Introduction
- 2 Convergence
- 3 Rappel sur la méthode par dichotomie
- 4 Méthode de la corde (de Lagrange)
- 5 Méthode de Newton
- 6 Dilemme robustesse/rapidité**
- 7 Résolution avec Scipy

Dilemme robustesse/rapidité

Nous avons vu différentes méthodes de résolution de l'équation $f(x) = 0$. Que choisir ?

Dilemme robustesse/rapidité

Robustesse

La méthode la plus naïve est la dichotomie.

A partir d'une fonction continue sur un intervalle $[a, b]$, avec $f(a).f(b) \leq 0$, la recherche de x tel que $f(x) = 0$ est simple et robuste mais lente (convergence linéaire).

Si l'on souhaite un programme simple et robuste, la dichotomie est une solution.

Dilemme robustesse/rapidité

Rapidité

Si on cherche une valeur avec précision, une convergence quadratique sera préférée (méthode de Newton).

A ce moment là, il convient d'avoir une idée approximative de la solution pour proposer un premier candidat au voisinage de cette solution de sorte que la fonction y présente de bonnes propriétés.

Si on ne souhaite pas calculer la dérivée de la fonction, on peut alors utiliser la méthode de la fausse position.

Dilemme robustesse/rapidité

Efficacité - compromis

Dans le cas où l'on cherche rapidité et stabilité, on peut utiliser la méthode par dichotomie dans un premier temps pour localiser le zéro de la fonction, puis appliquer un algorithme de Newton.

En cas d'instabilité de l'algorithme de Newton, il est toujours possible de réutiliser une méthode par dichotomie.

Enfin, mentionnons qu'il existe aussi d'autres types d'algorithmes, comme les algorithmes génétiques pour trouver les différentes solutions du problème.

Sommaire

- 1 Introduction
- 2 Convergence
- 3 Rappel sur la méthode par dichotomie
- 4 Méthode de la corde (de Lagrange)
- 5 Méthode de Newton
- 6 Dilemme robustesse/rapidité
- 7 Résolution avec Scipy**
 - En dimension 1
 - En dimension n

Résolution avec Scipy

En dimension 1

Pour résoudre un problème du type $f(x) = 0$, le plus simple est d'utiliser le module `scipy.optimize` avec la méthode `newton`.

Ainsi, `scipy.optimize.newton(f, x0)` permet de déterminer un zéro de la fonction f en partant de x_0 .

Il n'est pas nécessaire de donner la dérivée f_p de f mais il est possible de le faire `scipy.optimize.newton(f, x0, fp)`.

```
def f(x):  
    return x**2/2 - 1  
  
def fp(x):  
    return x  
  
sol = scipy.optimize.newton(f, 3, fp)
```

Résolution avec Scipy

Fonction lambda

RAPPEL il n'est pas nécessaire de déclarer les fonctions avant d'utiliser `newton`. La fonction `lambda` permet de faire lors de l'appel.

```
>>> import scipy.optimize as sciop
>>> sciop.newton(lambda x: x**2/2 - 1, 3)
1.4142135623730971
```

Résolution avec Scipy

En dimension n

```
>>> import scipy.optimize as sciop
>>> f = lambda X : [X[0]**2-2, X[1]**2-9]
>>> sciop.newton(f, [2, 2])
array([1.41421356, 3.          ])
```

Deuxième partie II

Interpolation, intégration et dérivation numérique

Sommaire

- 8 Interpolation et approximation de fonctions
 - Définitions
 - Interpolation polynomiale
 - Interpolation par morceaux
 - Approximation polynomiale par moindres carrés
 - Approximation polynomiale avec Numpy
- 9 Intégration numérique
- 10 Dérivation numérique

DÉFINITION : Interpolation

|| *Méthode qui consiste à déterminer une fonction (dans un ensemble donné), passant par un certain nombre de points imposés.*

DÉFINITION : Interpolation

|| *Méthode qui consiste à déterminer une fonction (dans un ensemble donné), passant par un certain nombre de points imposés.*

L'interpolation est utile lorsqu'une loi est donnée à partir d'une liste de points et qu'il est nécessaire d'évaluer le résultat en des points intermédiaires. Nous verrons par la suite que les fonctions d'interpolations sont aussi à la base de l'intégration et de la dérivation numérique.

DÉFINITION : Interpolation

|| *Méthode qui consiste à déterminer une fonction (dans un ensemble donné), passant par un certain nombre de points imposés.*

L'interpolation est utile lorsqu'une loi est donnée à partir d'une liste de points et qu'il est nécessaire d'évaluer le résultat en des points intermédiaires. Nous verrons par la suite que les fonctions d'interpolations sont aussi à la base de l'intégration et de la dérivation numérique.

DÉFINITION : Approximation

|| *Méthode qui consiste à déterminer une fonction passant "au mieux" à proximité des points donnés.*

DÉFINITION : Interpolation

|| *Méthode qui consiste à déterminer une fonction (dans un ensemble donné), passant par un certain nombre de points imposés.*

L'interpolation est utile lorsqu'une loi est donnée à partir d'une liste de points et qu'il est nécessaire d'évaluer le résultat en des points intermédiaires. Nous verrons par la suite que les fonctions d'interpolations sont aussi à la base de l'intégration et de la dérivation numérique.

DÉFINITION : Approximation

|| *Méthode qui consiste à déterminer une fonction passant "au mieux" à proximité des points donnés.*

Une approximation est utile lorsque une loi théorique est recherchée à partir de points de mesure (nombreux, mais entachés de bruit de mesure).

Interpolation polynomiale

DÉFINITION : Interpolation polynomiale

|| *Interpolation où la fonction est recherchée dans l'ensemble des polynômes.*

Interpolation polynomiale

DÉFINITION : Interpolation polynomiale

Interpolation où la fonction est recherchée dans l'ensemble des polynômes.

En exprimant le polynôme dans une base, par exemple dans la base canonique $(1, x, x^2, \dots, x^{n-1})$ (mais ce n'est pas le seul choix possible) :

$$P(x) = a_0 + a_1 \cdot x + a_2 \cdot x^2 + \dots + a_{n-1} \cdot x^{n-1}$$

L'opération consiste à déterminer les n composantes a_i du polynôme passant par les points imposés. Chaque point de passage constituant une condition scalaire sur les coefficients, il existe une solution unique si la dimension n de la base correspond au nombre de points.

Les composantes a_i sont solutions du système de n équations, pour n points imposés (x_i, y_i) :

$$y_i = a_0 + a_1.x_i + a_2.x_i^2 + \cdots + a_{n-1}.x_i^{n-1} \quad \text{avec} \quad i \in \llbracket 1..n \rrbracket$$

La **base polynomiale de Lagrange** est plus pratique pour l'expression directe du polynôme interpolant à partir des points, mais s'avère plus lourde à évaluer pour l'ordinateur. Le polynôme interpolant n points, s'écrit directement dans la base de Lagrange sous la forme :

La **base polynomiale de Lagrange** est plus pratique pour l'expression directe du polynôme interpolant à partir des points, mais s'avère plus lourde à évaluer pour l'ordinateur. Le polynôme interpolant n points, s'écrit directement dans la base de Lagrange sous la forme :

$$P(x) = \sum_{i=1}^n \prod_{j=1 \& j \neq i}^n y_j \cdot \frac{x - x_j}{x_i - x_j}$$

La **base polynomiale de Lagrange** est plus pratique pour l'expression directe du polynôme interpolant à partir des points, mais s'avère plus lourde à évaluer pour l'ordinateur. Le polynôme interpolant n points, s'écrit directement dans la base de Lagrange sous la forme :

$$P(x) = \sum_{i=1}^n \prod_{j=1 \& j \neq i}^n y_j \cdot \frac{x - x_j}{x_i - x_j}$$

EXEMPLE : Pour 3 points $\{(x_1, y_1), (x_2, y_2), (x_3, y_3)\}$, le polynôme interpolant s'écrit :

$$P(x) = y_1 \cdot \frac{(x - x_2) \cdot (x - x_3)}{(x_1 - x_2) \cdot (x_1 - x_3)} + y_2 \cdot \frac{(x - x_1) \cdot (x - x_3)}{(x_2 - x_1) \cdot (x_2 - x_3)} + y_3 \cdot \frac{(x - x_1) \cdot (x - x_2)}{(x_3 - x_1) \cdot (x_3 - x_2)}$$

L'expression se simplifie pour chaque (x_i, y_i) , ce qui permet de vérifier que le polynôme P passe bien par les points.

L'interpolation polynomiale n'est cependant pas idéale dès que le nombre de point augmente : le polynôme interpolant peut alors présenter des oscillations entre les points (phénomène de Runge, FIG 1). L'interpolation peut alors être localement très éloignée des points. Il est souvent mieux adapté d'interpoler par morceaux.

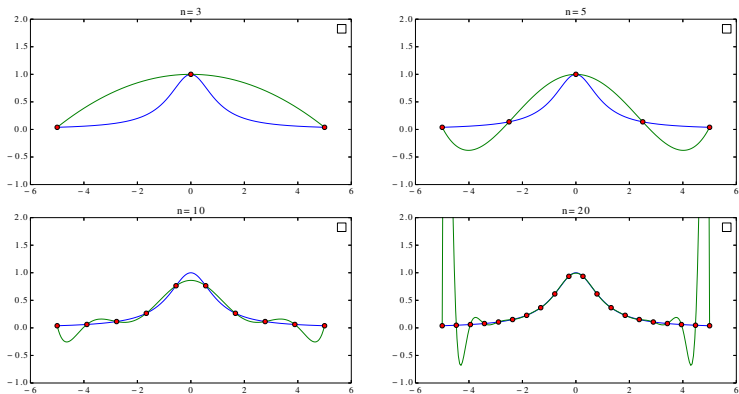


FIGURE 1 – Interpolation lagrangienne de degré n , $n \in \{3, 5, 10, 20\}$, de la fonction $f(x) = \frac{1}{x^2 + 1}$.

Interpolation par morceaux

Pour éviter les oscillations sur certaines fonctions, il est plus satisfaisant de réaliser une interpolation polynomiale de faible degré mais par morceaux.

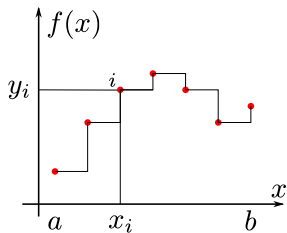


FIGURE 2 – Interpolation de degré 0.

La fonction interpolatrice n'est pas continue bien entendu.

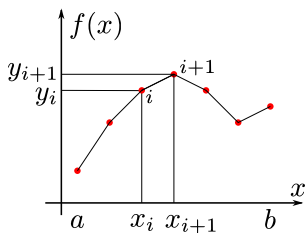


FIGURE 3 – Interpolation de degré 1.

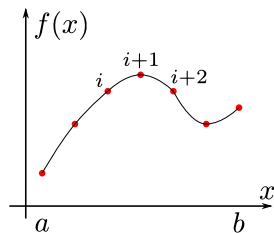


FIGURE 4 – Interpolation de degré 2.

Interpolation par morceaux de degrés 0

Il s'agit de considérer qu'entre deux points, la valeur de la fonction vaut une constante, égale à la valeur du point précédent, du point suivant ou encore égale à la moyenne des valeurs des points encadrant.

Cette interpolation est très rudimentaire mais elle peut être suffisante si le nombre de point est très important.

Interpolation linéaire par morceaux (degré 1)

Une loi affine ($a.x + b$) est adoptée entre deux points successifs, passant évidemment par les deux points :

$$\forall x \in [x_i, x_{i+1}], \quad y = y_i + (x - x_i) \cdot \frac{y_{i+1} - y_i}{x_{i+1} - x_i}$$

La fonction obtenue est continue mais sa dérivée ne l'est pas.

Interpolation quadratique par morceaux (degré 2)

Une loi parabolique ($a.x^2 + b.x + c$) est adoptée sur chaque intervalle regroupant 3 points successifs, passant par les 3 points.

La fonction obtenue est continue mais sa dérivée ne l'est pas car elle présente des discontinuités de la pente entre chaque portion de parabole. L'interpolation spline résout ce problème et assure une continuité C_1 .

Autres méthodes par morceaux

Il existe bien d'autres méthodes pour interpoler un ensemble de points, pour une fonction ou plus généralement pour une courbe du plan, ou encore pour des surfaces, des champs scalaires ou vectoriels (plans ou volumiques).

Pour les courbes dans le plan, citons l'interpolation d'Hermite cubique, qui est une interpolation de degrés 3 par morceau assurant la continuité de la dérivée aux extrémités des morceaux, ou encore les splines cubiques, qui est une interpolation cubique (degrés 3) par morceaux avec des conditions de continuité des deux premières dérivées aux extrémités des morceaux, ce qui conduit à une fonction de classe C_2 .

De même, les courbes de Bezier permettent d'interpoler des courbes ou des surfaces par des expressions polynomiales.

De même, les courbes de Bezier permettent d'interpoler des courbes ou des surfaces par des expressions polynomiales.

On retrouve alors un équivalent des cerces (Fig 5) utilisées en menuiserie, carrosserie...

De même, les courbes de Bezier permettent d'interpoler des courbes ou des surfaces par des expressions polynomiales.

On retrouve alors un équivalent des cerces (Fig 5) utilisées en menuiserie, carrosserie...

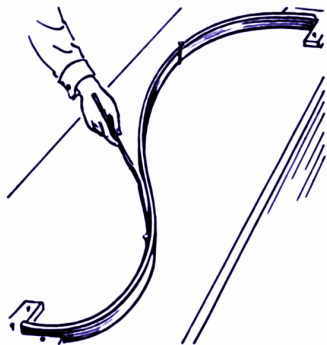


FIGURE 5 – Interpolation à l'aide d'une cerce en menuiserie

Choix du degré du polynôme d'interpolation

Pour limiter les oscillations (phénomène de Runge), nous avons déjà indiqué qu'il faut éviter les polynômes de degrés trop élevés (FIG 1). Par ailleurs, la qualité de la courbe devient très pauvre pour des degrés trop faibles, ce qui pousse à utiliser préférentiellement les degrés 1, 2 ou 3.

Néanmoins, soulignons que ce choix dépend aussi de la régularité de la courbe à interpoler car le choix d'un degré 2 ou 3 suppose une continuité de la dérivée sur la courbe d'origine. Si la courbe présente des discontinuités ou des ruptures de pentes, le degré 1 est à privilégier.

La FIG 6 montre l'interpolation par morceaux d'un signal sinusoïdal (C_∞) et d'un signal carré (discontinu). Le choix d'un degré 2 est bien approprié pour la première courbe mais pas pour la seconde !

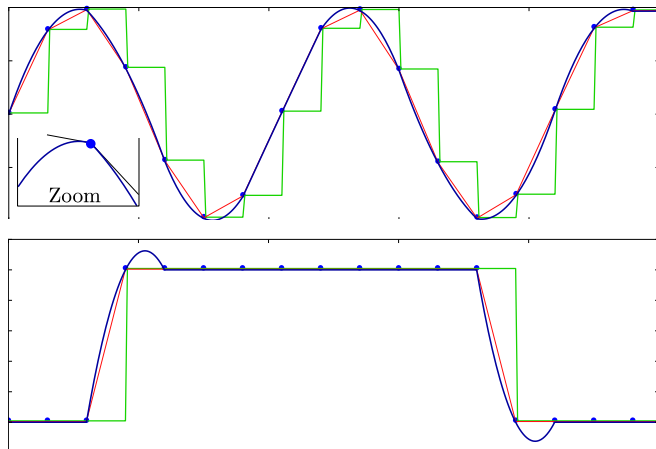


FIGURE 6 – Interpolations d'ordres 0, 1 et 2 pour une sinusoïde (fonction régulière) et un signal carré (fonction discontinue).

Approximation polynomiale par moindres carrés

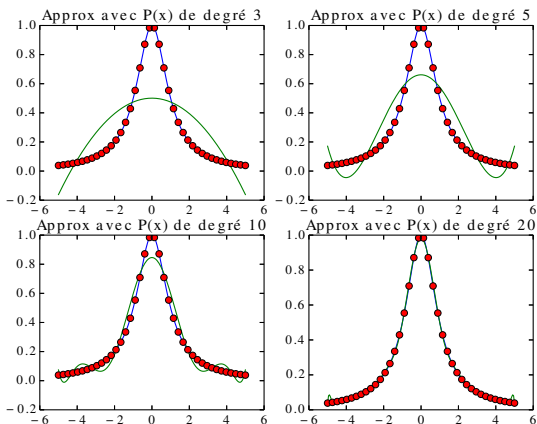


FIGURE 7 – Approximation avec un polynôme de degré n , $n \in \{3, 5, 10, 20\}$, de la fonction $f(x) = \frac{1}{x^2 + 1}$ discrétisée en 40 points.

Dans le cas de l'approximation polynomiale, on cherche à minimiser la distance (en norme 2) entre un polynôme de degré m et les n points imposés $(x_i, y_i)_1^n$.

$$\text{Si } P(x) = a_0 + a_1 \cdot x + \dots + a_m \cdot x^m$$

$$\text{alors } E_{rr}(a_0, \dots, a_m) = \frac{1}{2} \sum_{i=1}^n (P(x_i) - y_i)^2$$

Minimiser E_{rr} suivant les paramètres $(a_j)_{j=0}^m$, revient à résoudre le système suivant :

$$\begin{bmatrix} \frac{\partial E_{rr}}{\partial a_0} \\ \vdots \\ \frac{\partial E_{rr}}{\partial a_m} \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\Rightarrow \forall j \in \llbracket 0, m \rrbracket \sum_{i=1}^n \left(P(x_i) - y_i \right) \cdot \frac{\partial P(x_i)}{\partial a_j} = \sum_{i=1}^n \left(P(x_i) - y_i \right) \cdot x_i^j = 0$$

La minimisation de D revient alors à résoudre le système suivant:

$$\underbrace{\sum_{i=1}^n \begin{bmatrix} x_i^0 \cdot x_i^0 & \dots & x_i^m \cdot x_i^0 \\ \vdots & \ddots & \vdots \\ x_i^0 \cdot x_i^m & \dots & x_i^m \cdot x_i^m \end{bmatrix}}_M \cdot \underbrace{\begin{bmatrix} a_0 \\ \vdots \\ a_m \end{bmatrix}}_a = \underbrace{\sum_{i=1}^n y_i \cdot \begin{bmatrix} x_i^0 \\ \vdots \\ x_i^m \end{bmatrix}}_b \Rightarrow a = M^{-1} \cdot b$$

Minimiser E_{rr} suivant les paramètres $(a_j)_{j=0}^m$, revient à résoudre le système suivant :

$$\begin{bmatrix} \frac{\partial E_{rr}}{\partial a_0} \\ \vdots \\ \frac{\partial E_{rr}}{\partial a_m} \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\Rightarrow \forall j \in \llbracket 0, m \rrbracket \sum_{i=1}^n \left(P(x_i) - y_i \right) \cdot \frac{\partial P(x_i)}{\partial a_j} = \sum_{i=1}^n \left(P(x_i) - y_i \right) \cdot x_i^j = 0$$

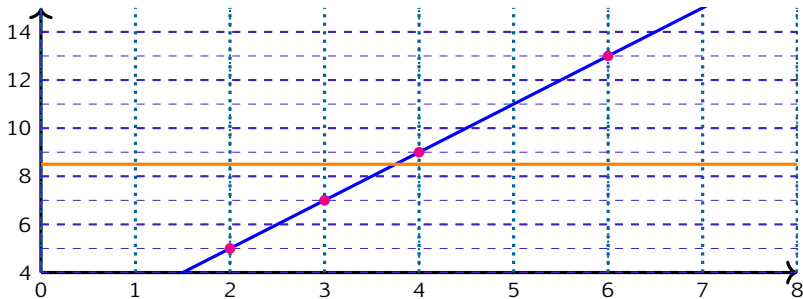
La minimisation de D revient alors à résoudre le système suivant:

$$\underbrace{\sum_{i=1}^n \begin{bmatrix} x_i^0 \cdot x_i^0 & \dots & x_i^m \cdot x_i^0 \\ \vdots & \ddots & \vdots \\ x_i^0 \cdot x_i^m & \dots & x_i^m \cdot x_i^m \end{bmatrix}}_M \cdot \underbrace{\begin{bmatrix} a_0 \\ \vdots \\ a_m \end{bmatrix}}_a = \underbrace{\sum_{i=1}^n y_i \cdot \begin{bmatrix} x_i^0 \\ \vdots \\ x_i^m \end{bmatrix}}_b \Rightarrow a = M^{-1} \cdot b$$

si M est inversible.

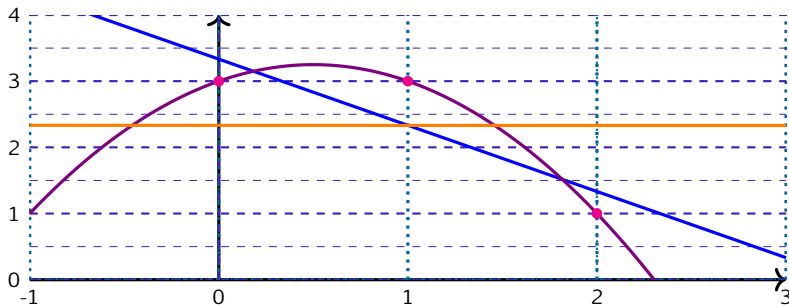
Approximation polynomiale avec Numpy

```
>>> X = [2, 3, 4, 6]
>>> Y = [5, 7, 9, 13]
>>> np.polyfit(X, Y, 0)
array([8.5])
>>> np.polyfit(X, Y, 1)
array([2., 1.])
>>> np.polyfit(X, Y, 2)
array([7.06708144e-16, 2.00000000e+00, 1.00000000e+00])
```



Approximation polynomiale avec Numpy

```
>>> X = [0, 1, 2]
>>> Y = [3, 3, 1]
>>> np.polyfit(X, Y, 0)
array([2.33333333])
>>> np.polyfit(X, Y, 1)
array([-1.          ,  3.33333333])
>>> np.polyfit(X, Y, 2)
array([-1.,  1.,  3.])
```



Sommaire

8 Interpolation et approximation de fonctions

9 Intégration numérique

- Principe, degré et ordre de la méthode
- Méthodes d'intégration composée
- Intégrer une fonction ou une liste de points avec Scipy

10 Dérivation numérique

Principe, degré et ordre de la méthode

L'**intégration numérique (ou quadrature)** consiste à intégrer (de façon approchée) une fonction sur un intervalle borné $[a, b]$, c'est-à-dire **calculer l'aire** sous la courbe représentant la fonction, à partir d'un calcul ou d'une mesure en un nombre fini n de points.

La répartition des points en abscisse est généralement uniforme (échantillonnage à pas constant $h = \frac{b-a}{n}$) mais il existe des méthodes à pas variable, ou encore à pas adaptatif.

L'intégration des polynômes étant très simple, l'opération consiste généralement à **construire une interpolation polynomiale** (de degré plus ou moins élevé) par morceaux (intégration composée) puis d'**intégrer le polynôme sur chaque morceau**.

Les méthodes de quadrature élémentaires composées sont de la forme:

$$\int_a^b f(x).dx = \sum_{i=0}^{n-1} \int_{x_i}^{x_{i+1}} f(x).dx \approx \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} \sum_{j=0}^m \omega_j \cdot f(\lambda(i,j)) = I(f, n, m)$$

$$\text{avec } x_i = a + i \cdot \frac{b-a}{n} \quad ; \quad \lambda(i,j) \in [x_i, x_{i+1}] \quad \text{et} \quad \sum_{j=0}^m \omega_j = 1$$

Les méthodes de quadrature élémentaires composées sont de la forme:

$$\int_a^b f(x).dx = \sum_{i=0}^{n-1} \int_{x_i}^{x_{i+1}} f(x).dx \approx \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} \sum_{j=0}^m \omega_j \cdot f(\lambda(i,j)) = I(f, n, m)$$

avec $x_i = a + i \cdot \frac{b-a}{n}$; $\lambda(i,j) \in [x_i, x_{i+1}]$ et $\sum_{j=0}^m \omega_j = 1$

DÉFINITION : Méthode de degré N

Méthode pour laquelle la formule approchée est exacte pour tout polynôme de degré au plus N et inexacte pour au moins un polynôme de degré $N + 1$.

Si on appelle $E_{rr}(f, n, N)$ la différence entre la solution exacte de l'intégrale et sa valeur approchée par une méthode d'ordre N sur n segments

$$\begin{aligned} E_{rr}(f, n, N) &= \left| \int_a^b f(x).dx - I(f, n, m) \right| \\ &= \left| \int_a^b f(x).dx - \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} \sum_{j=0}^m \omega_j \cdot f(\lambda(i, j)) \right| \end{aligned}$$

alors on peut montrer, moyennant une régularité suffisante de f , qu'il existe $K \in \mathbb{R}$ tel que :

$$E_{rr}(f, n, N) \leq \frac{K}{n^{N+1}}$$

Une méthode de degré N est donc d'ordre au moins $N + 1$.

La précision de l'intégration numérique peut ainsi s'améliorer en augmentant le nombre de points n (en diminuant le pas d'échantillonnage h) ou en augmentant le degré de l'interpolation polynomiale (sous réserve de bonnes propriétés de continuité de la courbe).

Méthodes d'intégration composée

Méthode des rectangles (degré 0)

L'intégration la plus simple est celle de degré 0 où le polynôme interpolateur sur chaque segment est une constante ($m = 0$ et $\omega_0 = 1$) prise soit à gauche, soit à droite de l'intervalle d'intégration. L'intégrale est donc approchée par des rectangles pour calculer l'aire sous la courbe :

$$\int_a^b f(x).dx \approx \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} f(\lambda(i,0))$$

Méthode des rectangles (degré 0)

L'intégration la plus simple est celle de degré 0 où le polynôme interpolateur sur chaque segment est une constante ($m = 0$ et $\omega_0 = 1$) prise soit à gauche, soit à droite de l'intervalle d'intégration. L'intégrale est donc approchée par des rectangles pour calculer l'aire sous la courbe :

$$\int_a^b f(x).dx \approx \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} f(\lambda(i,0))$$

$$\Rightarrow E_{rr}(f, n, 0) \leq \frac{K}{n} \quad \text{avec} \quad K = \frac{\sup |f'|}{2} \cdot |b-a|^2$$

Rectangles à gauche (degré 0)

Dans le cas de la formule des rectangles à gauche, on pose $\lambda(i, 0) = x_i$:

$$\int_a^b f(x).dx \approx \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} f(x_i) = \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} y_i$$

Rectangles à droite (degré 0)

Dans le cas de la formule des rectangles à droite, on pose $\lambda(i, 0) = x_{i+1}$:

$$\int_a^b f(x).dx \approx \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} f(x_{i+1}) = \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} y_{i+1}$$

Méthode point milieu (degré 1)

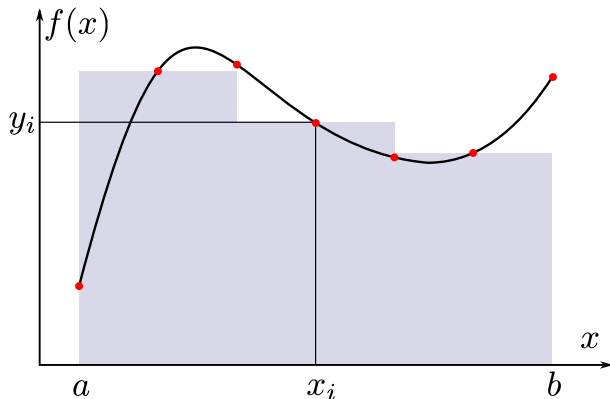


FIGURE 8 – Intégration au point milieu.

La méthode du point milieu consiste à considérer la fonction interpolante constante sur chaque intervalle $[x_i, x_{i+1}]$ et égale à la valeur prise par le point au milieu de l'intervalle $\lambda(i, 0) = \frac{x_i + x_{i+1}}{2}$ (FIG 8). La valeur approchée de l'intégrale s'écrit alors :

$$I(f, n, 0) = \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} f\left(\frac{x_i + x_{i+1}}{2}\right)$$

$$\Rightarrow E_{rr}(f, n, 0) \leq \frac{K}{n^2} \quad \text{avec} \quad K = \frac{\sup |f''|}{24} \cdot |b-a|^3$$

La méthode est d'ordre 2. En effet, la méthode du point milieu est exacte pour les polynômes de degré 1.

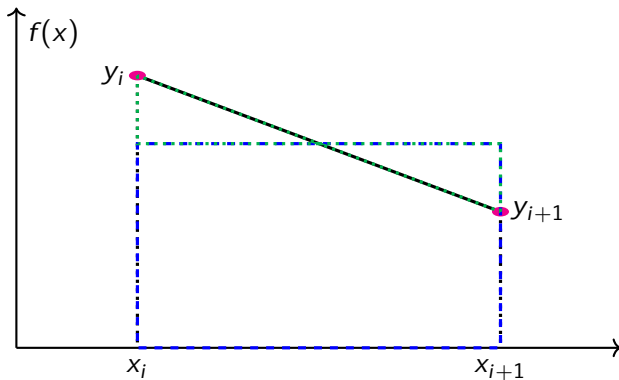


FIGURE 9 – Méthode point milieu - degré 1

Méthode trapèze (degré 1)

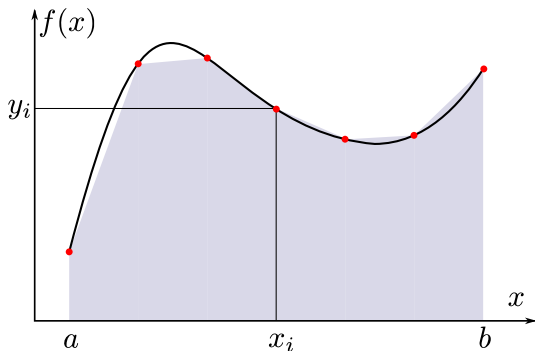


FIGURE 10 – Intégration par la méthode trapèze.

La méthode du trapèze s'appuie sur une interpolation linéaire entre chaque point (Fig 10). La valeur approchée de l'intégrale s'écrit alors :

$$I(f, n, 1) = \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} \sum_{j=0}^1 \omega_j \cdot f(\lambda(i, j))$$

$$\text{avec } \omega_0 = \omega_1 = \frac{1}{2} \quad \text{et} \quad \lambda(i, 0) = x_i; \lambda(i, 1) = x_{i+1}$$

$$\frac{b-a}{n} \cdot \sum_{i=0}^{n-1} \frac{y_i + y_{i+1}}{2} = \frac{b-a}{n} \cdot \left(\frac{y_0 + y_n}{2} + \left(\sum_{i=2}^{n-1} y_i \right) \right)$$

$$\Rightarrow E_{rr}(f, n, 1) \leq \frac{K}{n^2} \quad \text{avec} \quad K = \frac{\sup |f''|}{12} \cdot |b - a|^3$$

Cette expression est à nouveau similaire au calcul approché réalisé au paragraphe précédent.

La méthode du point milieu est au final plus précise que celle des trapèzes.

Méthode de Simpson (degré 2)

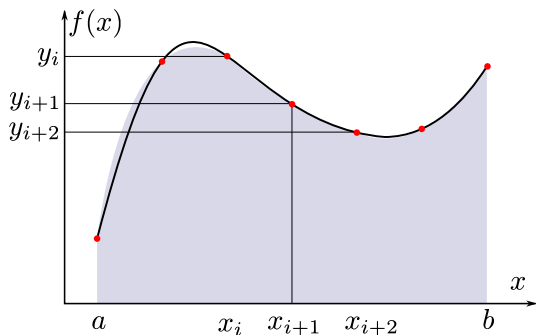


FIGURE 11 – Intégration par la méthode de Simpson.

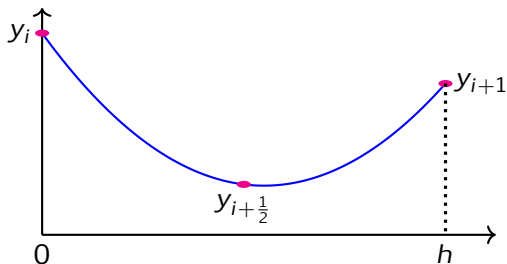
La méthode de Simpson s'appuie sur une interpolation quadratique ($a.x^2 + b.x + c$) sur chaque intervalle $[x_i, x_{i+1}]$: (FIG 11).

$$\int_a^b f(x).dx \approx \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} \sum_{j=0}^2 \omega_j \cdot f(\lambda(i,j)) = I(f, n, 2)$$

$$\text{avec } \lambda(i,0) = x_i; \lambda(i,1) = x_{i+\frac{1}{2}} = \frac{x_i + x_{i+1}}{2} \quad \text{et} \quad \lambda(i,2) = x_{i+1}$$

En appliquant une translation de l'intervalle $[x_i, x_{i+1}]$ sur l'intervalle $[0, h]$, l'interpolation s'écrit :

$$\begin{cases} y_i = c \\ y_{i+\frac{1}{2}} = a \cdot \frac{h^2}{4} + b \cdot \frac{h}{2} + c \\ y_{i+1} = a \cdot h^2 + b \cdot h + c \end{cases} \Rightarrow \begin{cases} a = 2 \cdot \frac{y_i + y_{i+1} - 2 \cdot y_{i+\frac{1}{2}}}{h^2} \\ b = \frac{4 \cdot y_{i+\frac{1}{2}} - y_{i+1} - 3 \cdot y_i}{h} \\ c = y_i \end{cases}$$



L'intégrale sur $[0, h]$ vaut alors :

L'intégrale sur $[0, h]$ vaut alors :

$$\begin{aligned}\int_0^h (a.x^2 + b.x + c).dx &= \left[a.\frac{x^3}{3} + b.\frac{x^2}{2} + c \right]_0^h \\ &= \left(a.\frac{h^3}{3} + b.\frac{h^2}{2} + c.h \right) \\ &= h.\frac{(4-9+6).y_i + (-8+12).y_{i+\frac{1}{2}} + (4-3).y_{i+1}}{6} \\ &= \frac{h}{6} \cdot \left(y_i + 4.y_{i+\frac{1}{2}} + y_{i+1} \right)\end{aligned}$$

$$\text{d'où } \omega_0 = \omega_2 = \frac{1}{6} \quad \text{et} \quad \omega_1 = \frac{2}{3}$$

La méthode est de degré 3 et d'ordre 4. En effet :

$$\text{d'où } \omega_0 = \omega_2 = \frac{1}{6} \quad \text{et} \quad \omega_1 = \frac{2}{3}$$

La méthode est de degré 3 et d'ordre 4. En effet :

$$\int_a^b f(x).dx \approx \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} \frac{y_i + 4.y_{i+\frac{1}{2}} + y_{i+1}}{6}$$

$$\text{d'où } \omega_0 = \omega_2 = \frac{1}{6} \quad \text{et} \quad \omega_1 = \frac{2}{3}$$

La méthode est de degré 3 et d'ordre 4. En effet :

$$\int_a^b f(x).dx \approx \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} \frac{y_i + 4.y_{i+\frac{1}{2}} + y_{i+1}}{6}$$

$$\Rightarrow E_{rr}(f, n, 2) \leq \frac{K}{n^4} \quad \text{avec} \quad K = \frac{\sup |f^{(4)}|}{180} \cdot |b-a|^5$$

L'interpolation de Simpson est plus précise que l'interpolation trapèze lorsque la fonction à intégrer est raisonnablement continue. Elle se justifie beaucoup moins lorsque la fonction présente des discontinuités.

Lors d'une intégration "temps réel", elle introduit par ailleurs un décalage temporel de $2.h$ (pour éviter $y_{i+\frac{1}{2}} \dots$) qui est souvent plus pénalisant que l'erreur d'intégration.

Autres méthodes

Méthode de	m	degré	pois
Boole-Villarceau	4	5	$\omega_0 = \omega_4 = \frac{7}{90}$; $\omega_1 = \omega_3 = \frac{16}{45}$; $\omega_2 = \frac{2}{15}$
Weddle-Hardy	6	7	$\omega_0 = \omega_6 = \frac{41}{840}$; $\omega_1 = \omega_5 = \frac{9}{35}$; $\omega_2 = \omega_4 = \frac{9}{280}$; $\omega_3 = \frac{34}{105}$

Autres méthodes

Méthode de	m	degré	poids
Boole-Villarceau	4	5	$\omega_0 = \omega_4 = \frac{7}{90}; \omega_1 = \omega_3 = \frac{16}{45}; \omega_2 = \frac{2}{15}$
Weddle-Hardy	6	7	$\omega_0 = \omega_6 = \frac{41}{840}; \omega_1 = \omega_5 = \frac{9}{35}; \omega_2 = \omega_4 = \frac{9}{280}; \omega_3 = \frac{34}{105}$

Méthode de	m	degré	points
Boole-Villarceau	4	5	$\lambda(i, j) = x_i + j \cdot \frac{x_{i+1} - x_i}{4}$
Weddle-Hardy	6	7	$\lambda(i, j) = x_i + j \cdot \frac{x_{i+1} - x_i}{6}$

Autres méthodes

Méthode de	m	degré	poids
Boole-Villarceau	4	5	$\omega_0 = \omega_4 = \frac{7}{90}; \omega_1 = \omega_3 = \frac{16}{45}; \omega_2 = \frac{2}{15}$
Weddle-Hardy	6	7	$\omega_0 = \omega_6 = \frac{41}{840}; \omega_1 = \omega_5 = \frac{9}{35}; \omega_2 = \omega_4 = \frac{9}{280}; \omega_3 = \frac{34}{105}$

Méthode de	m	degré	points
Boole-Villarceau	4	5	$\lambda(i, j) = x_i + j \cdot \frac{x_{i+1} - x_i}{4}$
Weddle-Hardy	6	7	$\lambda(i, j) = x_i + j \cdot \frac{x_{i+1} - x_i}{6}$

Voir aussi les méthodes de Newton-Cotes.

Intégrer une fonction ou une liste de points avec Scipy

OBJECTIF : obtenir une approximation de $\int_a^b f(x).dx$

Pour obtenir une approximation d'une intégrale simple, double ou triple d'une fonction, le module `scipy.integrate` propose les fonctions `quad`, `dblquad` et `tplquad`.

Ainsi `quad(f, a, b)` renvoie un tuple contenant une approximation de $\int_a^b f(x).dx$ et une estimation de l'erreur commise.

```
>>> scint.quad(lambda x : x**3, 0, 2)
(4.0, 4.440892098500626e-14)
```

Intégrale double

$$I = \int_{y_g}^{y_d} \left(\int_{x_g}^{x_d} f(y, x) \cdot dx \right) \cdot dy$$

\approx `scint.dblquad(f, xg, xd, lambda x: yg, lambda x, yd)`

Aire d'une sphère de rayon unitaire.

```
def jac2(theta, phi):  
    return m.sin(theta)  
  
A = scint.dblquad(jac2, 0, 2*m.pi, lambda x: 0, lambda x: m.pi)
```

Intégrale triple

Volume d'une boule de rayon unitaire.

```
def jac3(r, theta, phi):  
    return r**2*m.sin(theta)  
  
V = scint.tplquad(jac3, 0, 2*m.pi, lambda x: 0, lambda x: m.pi,  
                 lambda x, y: 0, lambda x, y : 1)
```

Utiliser `cumtrapz` du module `scipy.integrate` permet d'intégrer une liste de valeurs, par exemple pour passer des mesures d'accélération à la vitesse en chaque point de mesure.

Ainsi `scipy.integrate.cumtrapz(val, x=t, initial=i0)` permet d'obtenir l'évolution de l'aire sous la courbe formée par les éléments de `val`, avec `t`, la liste de leurs abscisses et `i0` la valeur initiale.

REMARQUE : la fonction ne semble pas tenir compte de la valeur initiale. Le troisième exemple serait peut-être la solution à adopter pour obtenir le résultat recherché :

```
>>> scint.cumtrapz([2, 3, 4], x=[0, 1, 2], initial=3)
array([3. , 2.5, 6. ])
>>> scint.cumtrapz([2, 3, 4], x=[0, 1, 2], initial=0)
array([0. , 2.5, 6. ])
>>> scint.cumtrapz([2, 3, 4], x=[0, 1, 2], initial=0) +3
array([3. , 5.5, 9. ])
```

Sommaire

8 Interpolation et approximation de fonctions

9 Intégration numérique

10 Dérivation numérique

- Dérivée première
- Dérivée seconde
- Influence du bruit de mesure

La dérivation numérique consiste à dériver (de façon approchée) une fonction sur un intervalle borné $[a, b]$, c'est-à-dire calculer la pente de la courbe représentant la fonction, à partir d'un calcul ou d'une mesure en un nombre fini de points.

La répartition des points en abscisse est généralement uniforme (pas d'échantillonnage constant h) mais il existe des méthodes à pas variable, ou encore à pas adaptatif.

La dérivation numérique consiste à dériver (de façon approchée) une fonction sur un intervalle borné $[a, b]$, c'est-à-dire calculer la pente de la courbe représentant la fonction, à partir d'un calcul ou d'une mesure en un nombre fini de points.

La répartition des points en abscisse est généralement uniforme (pas d'échantillonnage constant h) mais il existe des méthodes à pas variable, ou encore à pas adaptatif.

La dérivation des polynômes étant très simple, l'opération consiste généralement à construire une interpolation polynomiale par morceaux (de degré plus ou moins élevé) puis de dériver le polynôme sur chaque morceau.

La précision de l'intégration numérique peut s'améliorer en augmentant le nombre de points n (en diminuant le pas d'échantillonnage h) ou en augmentant le degré de l'interpolation polynomiale (sous réserve de bonnes propriétés de continuité de la courbe).

Dérivée première

Méthode à 1 pas

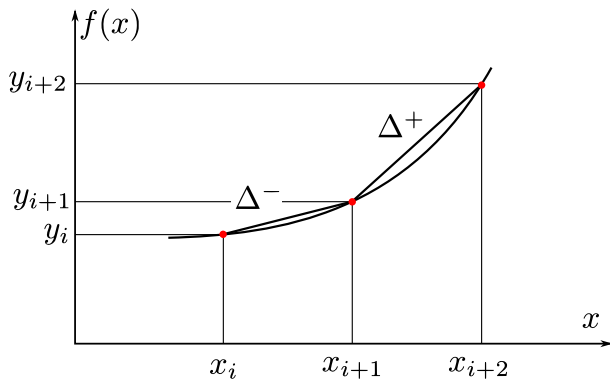


FIGURE 12 – Dérivation à 1 pas.

L'estimation de la dérivée la plus simple consiste à calculer la pente à partir du point courant et du point précédent ou suivant (FIG 12). L'estimation de la dérivée au point i peut alors se calculer par :

- différence avant : $D_i = \frac{1}{h}(y_{i+1} - y_i)$ (pente de la droite Δ^+),
- différence arrière : $D_i = \frac{1}{h}(y_i - y_{i-1})$ (pente de la droite Δ^-),

L'estimation de la dérivée la plus simple consiste à calculer la pente à partir du point courant et du point précédent ou suivant (FIG 12). L'estimation de la dérivée au point i peut alors se calculer par :

- différence avant : $D_i = \frac{1}{h}(y_{i+1} - y_i)$ (pente de la droite Δ^+),
- différence arrière : $D_i = \frac{1}{h}(y_i - y_{i-1})$ (pente de la droite Δ^-),

Évidemment, lorsqu'il s'agit de dériver une fonction temporelle "en temps réel", le point suivant n'est pas encore connu donc seule la différence arrière peut être calculée.

Notons aussi que le calcul de la dérivée à partir de n points, conduit à un tableau de valeur de dimension $n - 1$.

Méthode à 2 pas

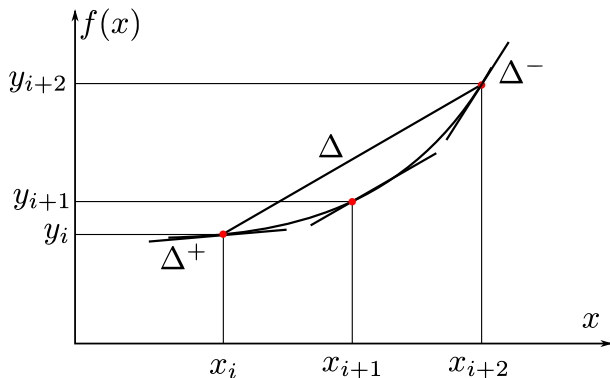


FIGURE 13 – Dérivation à 2 pas.

Une méthode à 2 pas consiste à utiliser 3 points pour une meilleure estimation de la dérivée (FIG 13).

Sur un intervalle $[x_i, x_{i+2}]$ de longueur $H = 2.h$, la courbe est interpolée par un polynôme d'ordre 2 : $y = a_i.x^2 + b_i.x + c_i$. Les coefficients a_i , b_i et c_i ont déjà été calculés au paragraphe 2 en fonction des y_i . Il suffit de dériver le polynôme respectivement en x_i , x_{i+1} et x_{i+2} pour obtenir respectivement les différences avant en x_i , centrée en x_{i+1} et arrière en x_{i+2} .

Après translation des abscisses, il s'agit en réalité de dériver en 0, $\frac{H}{2} = h$ et $H = 2.h$.

- différence avant :

$$D_i = b_i = \frac{4.y_{i+1} - y_{i+2} - 3.y_i}{2.h} \text{ (pente de la droite } \Delta^+)$$

- différence avant :

$$D_i = b_i = \frac{4.y_{i+1} - y_{i+2} - 3.y_i}{2.h} \text{ (pente de la droite } \Delta^+)$$

- différence centrée :

$$\begin{aligned} D_{i+1} &= 2.a_j.h + b_i = 2.\frac{2.y_i + 2.y_{i+2} - 4.y_{i+1}}{4.h^2}.h + \frac{4.y_{i+1} - y_{i+2} - 3.y_i}{2.h} \\ &= \frac{y_{i+2} - y_i}{2.h} \text{ (pente de la droite } \Delta) \end{aligned}$$

- différence avant :

$$D_i = b_i = \frac{4.y_{i+1} - y_{i+2} - 3.y_i}{2.h} \text{ (pente de la droite } \Delta^+)$$

- différence centrée :

$$\begin{aligned} D_{i+1} &= 2.a_j.h + b_i = 2.\frac{2.y_i + 2.y_{i+2} - 4.y_{i+1}}{4.h^2}.h + \frac{4.y_{i+1} - y_{i+2} - 3.y_i}{2.h} \\ &= \frac{y_{i+2} - y_i}{2.h} \text{ (pente de la droite } \Delta) \end{aligned}$$

- différence arrière :

$$\begin{aligned} D_{i+2} &= 2.a_j.2.h + b_i = 4.\frac{2.y_i + 2.y_{i+2} - 4.y_{i+1}}{4.h^2}.h + \frac{4.y_{i+1} - y_{i+2} - 3.y_i}{2.h} \\ &= \frac{y_i - 4.y_{i+1} + 3.y_{i+2}}{2.h} \text{ (pente de la droite } \Delta^-) \end{aligned}$$

Il faut noter que la différence centrée est aussi simple à calculer que dans le cas d'une méthode à 1 pas, et correspond à la pente entre les deux points de part et d'autre du point courant. La précision étant d'ordre 2, elle constitue un très bon compromis.

Évidemment, lorsqu'il s'agit de dériver une fonction temporelle "en temps réel", le ou les points suivants ne sont pas encore connus donc seule la différence arrière peut être calculée.

Dérivée seconde

Pour dériver deux fois une courbe donnée par une liste de points, le premier réflexe consiste à appliquer deux fois une dérivée simple (par une méthode à 1 ou 2 pas).

Il est cependant plus rapide de calculer directement la dérivée seconde à partir du polynôme d'interpolation, qui doit bien évidemment être au moins de degré 2 pour obtenir un résultat non nul.

En utilisant une méthode à 2 pas, le polynôme d'interpolation s'écrit $y = a.x^2 + b.x + c$. Les coefficients a , b et c ont déjà été calculés au paragraphe 2 en fonction des y_i , ce qui conduit à une dérivée seconde constante sur tout l'intervalle $[x_i, x_{i+2}]$:

$$f''(x) = 2.a = 2. \frac{2.y_i + 2.y_{i+2} - 4.y_{i+1}}{H^2} = \frac{y_i + y_{i+2} - 2.y_{i+1}}{h^2}$$

Selon si cette valeur de la dérivée seconde est adoptée en i , $i + 1$ ou $i + 2$, il s'agit de la différence seconde avant, centrée ou arrière. On peut montrer que cela revient à calculer deux dérivées simples à 1 pas.

Influence du bruit de mesure

Lorsque la courbe est issue d'une mesure, elle est généralement entachée d'un léger bruit, qui peut devenir catastrophique pour l'évaluation de la dérivée (FIG 14).

En effet, si les points de mesure restent "en moyenne" au voisinage de la valeur mesurée, il existe des fluctuations rapides (c'est-à-dire à la fréquence d'échantillonnage) entre les points successifs (voir zoom de la FIG 14).

Le calcul de la dérivée conduit à déterminer la pente entre deux points successifs, ce qui perturbe fortement le signal dérivé et cache les évolutions lentes du signal (lentes devant la période d'échantillonnage).

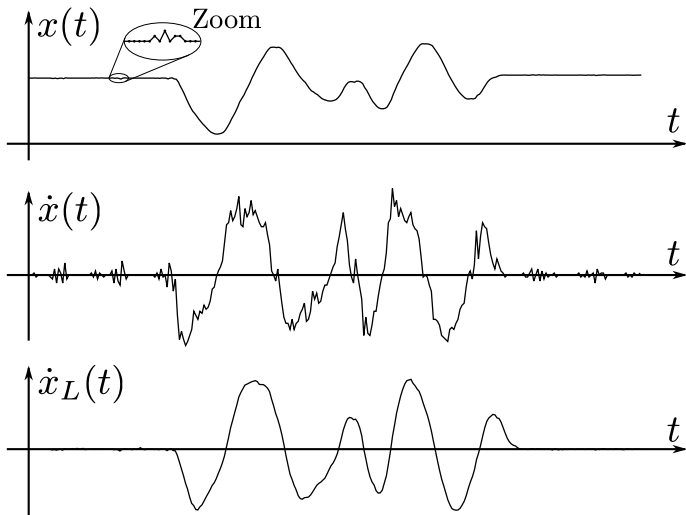


FIGURE 14 – Mesure d’une position au cours du temps $x(t)$ par un capteur potentiométrique, dérivée à 1 pas $\dot{x}(t)$ et dérivée à 1 pas lissée en effectuant la différence sur 10 pas.

Deux solutions sont possibles :

- filtrer (ou lisser) le signal d'origine pour supprimer l'essentiel du bruit, puis dériver,
- calculer la dérivée sur un temps plus long que le temps d'échantillonnage, par exemple pour une méthode à 1 pas en calculant la pente entre deux points espacés de k pas (solution adoptée pour $\dot{x}_L(t)$ sur la FIG 14, avec $k = 10$).

Dans les deux cas, le signal dérivé sera entaché d'un retard sur le signal d'origine, ce qui oblige à trouver un compromis entre la qualité du signal dérivé et le retard. k est généralement de l'ordre de 5 à 20 pas selon le bruit, le pas d'échantillonnage, les fréquences à observer dans le signal et le retard admissible.

Pour un lissage par moyenne mobile, on peut montrer que les deux méthodes s'avèrent identiques.

Troisième partie III

Intégration des équations différentielles

Sommaire

- 11 Mise en place du problème
 - Problème de Cauchy
 - Existence et unicité de la solution
 - Résolution numérique
- 12 Annexe mathématique
- 13 Méthodes à un pas
- 14 Mise en forme des systèmes d'équations différentielles
- 15 Résolution avec Scipy

Equations différentielles non linéaires

Une grande part des problèmes scientifiques se modélisent par une équation différentielle dont on cherche la solution pour dimensionner ou comprendre le phénomène.

Quand les équations sont simples (linéaires d'ordre 1 ou 2) la résolution analytique est aisée, mais nombre de modélisations conduisent à des équations non linéaires. Ce cours a pour objectif de proposer des méthodes pour déterminer une solution approchée.

Problème de Cauchy

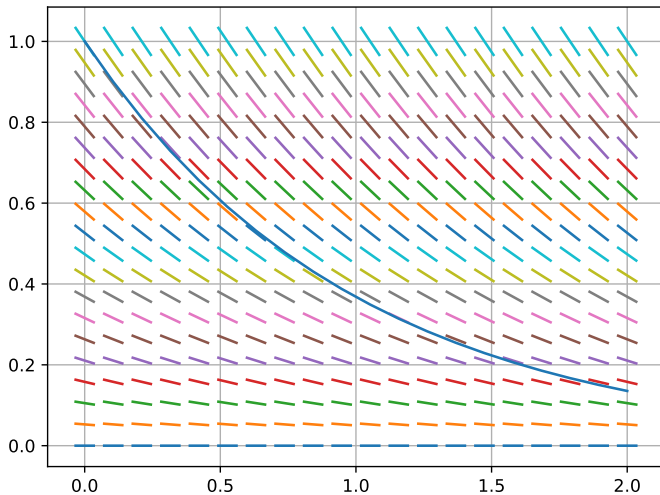
Le problème de Cauchy consiste à trouver les fonctions

$$Y \text{ de } [0, T] \rightarrow \mathbb{R}^N, \text{ telles que : } \begin{cases} \frac{dY}{dt} = F(Y, t) \\ Y(t_0) = Y_0 \end{cases}$$

où $t_0 \in [0, T]$ et $Y_0 \in \mathbb{R}^N$ sont des données.

La plupart des systèmes d'équations différentielles de tout ordre peuvent se mettre sous cette forme de système d'équations différentielles du premier ordre^a.

a. à l'exception des équations différentielles implicites



Existence et unicité de la solution

Théorème de Cauchy-Lipschitz : soit F une fonction de $\mathbb{R}^N \times [0, T] \rightarrow \mathbb{R}^N$ continue et lipschitzienne en Y .

Alors, $\forall t_0 \in [0, T]$ et $\forall Y_0 \in \mathbb{R}^N$, le problème de Cauchy admet une solution unique définie sur $[0, T]$.

Rappel : F lipschitzienne en Y :

$\exists k > 0$ tel que, $\forall Y \in \mathbb{R}^N, \forall Z \in \mathbb{R}^N, \forall t \in [0, T]$,

$$\|F(Y, t) - F(Z, t)\| \leq k \cdot \|Y - Z\|$$

Résolution numérique

OBJECTIF : Obtenir une solution approchée au problème de Cauchy pour une discrétisation temporelle de l'intervalle donné.

Résolution numérique

OBJECTIF : Obtenir une solution approchée au problème de Cauchy pour une discrétisation temporelle de l'intervalle donné.

La démarche s'appuie sur la forme initiale de l'équation différentielle où il est facile de voir que la fonction $F(Y, t)$ traduit l'évolution de Y , c'est-à-dire la pente de la courbe $Y(t)$.

Résolution numérique

OBJECTIF : Obtenir une solution approchée au problème de Cauchy pour une discrétisation temporelle de l'intervalle donné.

La démarche s'appuie sur la forme initiale de l'équation différentielle où il est facile de voir que la fonction $F(Y, t)$ traduit l'évolution de Y , c'est-à-dire la pente de la courbe $Y(t)$.

Les schémas d'intégration numériques exploitent $F(Y, t)$ pour traduire l'évolution sur un pas de temps.

En notant h le pas de temps et N le nombre de pas de temps, l'intervalle est discrétisé par $t_0 = 0, t_1 = h, \dots, t_N = N.h$.

En notant h le pas de temps et N le nombre de pas de temps, l'intervalle est discrétisé par $t_0 = 0, t_1 = h, \dots, t_N = N.h$.

Sur un sous-intervalle donné, on cherche à déterminer la solution sous la forme :

$$\int_{t_i}^{t_{i+1}} \frac{dY}{dt} dt = \int_{t_i}^{t_{i+1}} F(Y, t) dt$$

En notant h le pas de temps et N le nombre de pas de temps, l'intervalle est discrétisé par $t_0 = 0, t_1 = h, \dots, t_N = N.h$.

Sur un sous-intervalle donné, on cherche à déterminer la solution sous la forme :

$$\int_{t_i}^{t_{i+1}} \frac{dY}{dt} dt = \int_{t_i}^{t_{i+1}} F(Y, t) dt$$

On en déduit que :

$$Y_{i+1} = Y_i + \int_{t_i}^{t_{i+1}} F(Y, t) dt$$

Les méthodes de résolution numérique des équations différentielles sont basées sur des techniques d'estimation de l'intégrale de la fonction F .

Les méthodes de résolution numérique des équations différentielles sont basées sur des techniques d'estimation de l'intégrale de la fonction F .

On s'intéressera dans ce cours aux méthodes de résolution à un pas qui se mettent sous la forme générale :

Les méthodes de résolution numérique des équations différentielles sont basées sur des techniques d'estimation de l'intégrale de la fonction F .

On s'intéressera dans ce cours aux méthodes de résolution à un pas qui se mettent sous la forme générale :

$$\left\{ \begin{array}{l} Y_{i+1} = Y_i + h \cdot \vec{\Phi}(Y_j, t_j, h) \\ Y(t_0) = Y_0 \end{array} \right.$$

Les méthodes de résolution numérique des équations différentielles sont basées sur des techniques d'estimation de l'intégrale de la fonction F .

On s'intéressera dans ce cours aux méthodes de résolution à un pas qui se mettent sous la forme générale :

$$\left\{ \begin{array}{l} Y_{i+1} = Y_i + h \cdot \vec{\Phi}(Y_j, t_j, h) \\ Y(t_0) = Y_0 \end{array} \right.$$

où il faudra choisir $\vec{\Phi}$.

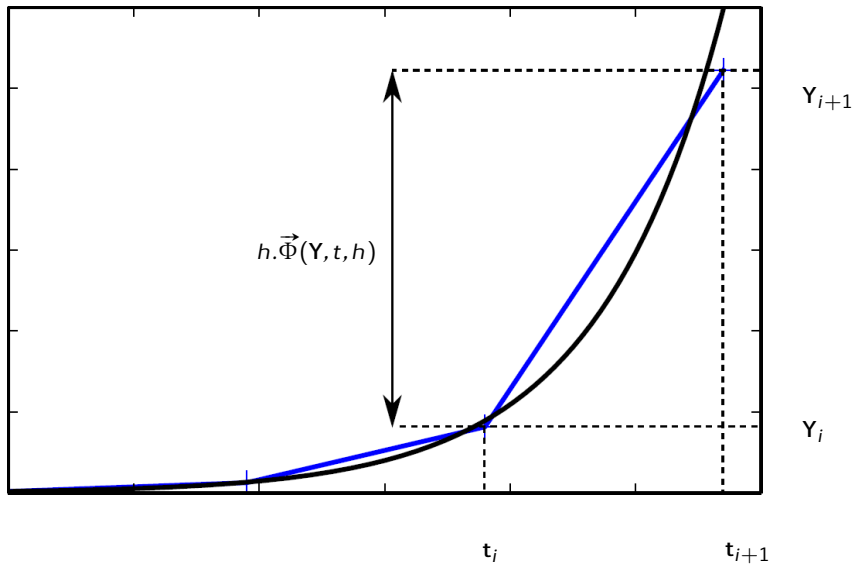
Les méthodes de résolution numérique des équations différentielles sont basées sur des techniques d'estimation de l'intégrale de la fonction F .

On s'intéressera dans ce cours aux méthodes de résolution à un pas qui se mettent sous la forme générale :

$$\begin{cases} Y_{i+1} = Y_i + h \cdot \vec{\Phi}(Y_j, t_j, h) \\ Y(t_0) = Y_0 \end{cases}$$

où il faudra choisir $\vec{\Phi}$.

La fonction $\vec{\Phi}$ représente la pente de la droite permettant de passer de Y_i à Y_{i+1} .



Sommaire

- 11 Mise en place du problème
- 12 Annexe mathématique
 - Définitions
 - Propriétés
- 13 Méthodes à un pas
- 14 Mise en forme des systèmes d'équations différentielles
- 15 Résolution avec Scipy

Définitions

On note $Y_{ex}(t)$ la solution exacte.

Erreur de consistance au pas i

DÉFINITION : Erreur de consistance au pas i

On note c_i l'erreur de consistance au pas i telle que

$$c_i = Y_{\text{ex}}(t_{i+1}) - Y_{\text{ex}}(t_i) - h \cdot \vec{\Phi}(Y_{\text{ex}}(t_j), t_j, h)$$

Erreur de consistance au pas i

DÉFINITION : Erreur de consistance au pas i

On note c_i l'erreur de consistance au pas i telle que

$$c_i = Y_{\text{ex}}(t_{i+1}) - Y_{\text{ex}}(t_i) - h \cdot \vec{\Phi}(Y_{\text{ex}}(t_j), t_j, h)$$

Elle permet de définir l'erreur réalisée entre l'approximation et la solution exacte au pas de temps i .

Méthode consistante

DÉFINITION : Méthode consistante

La méthode est dite consistante si $\lim_{h \rightarrow 0} \sum_{i=1}^N |c_i| = 0$

Méthode stable

Soit le schéma choisi et le schéma perturbé

$$\begin{cases} Y_{i+1} = Y_i + h \cdot \vec{\Phi}(Y_j, t_j, h) \\ Y_{\text{ex}}(t_0) = Y_0 \end{cases} \quad \begin{cases} Y_{i+1}^* = Y_i^* + h \cdot \vec{\Phi}(Y_j^*, t_j, h) + \zeta_i \\ Y_{\text{ex}}^*(t_0) = Y_0^* \end{cases}$$

DÉFINITION : Méthode stable

La méthode est stable si

$$\exists M > 0 / \sup_{0 \leq i \leq N} |Y_i^* - Y_i| \leq M \cdot \left(|Y_0^* - Y_0| + \sum_{i=0}^N |\zeta_i| \right)$$

Méthode convergente

DÉFINITION : Méthode convergente

La méthode est convergente si

$$\lim_{h \rightarrow 0} \left[\sup_{0 \leq i \leq N} |Y_{ex}(t_i) - Y_i| \right] = 0$$

Ordre de convergence

DÉFINITION : Méthode d'ordre p

Une méthode est dite d'ordre p si

$$\exists M > 0 / \forall i, |c_i| \leq M.h^{p+1}$$

Propriétés

Propriété 1

Une méthode est convergente si et seulement si elle est stable et consistante.

Propriété 2

Si $\vec{\Phi}$ est lipschitzienne en \mathbf{Y} alors le schéma numérique est stable.

Démonstration :

Soit le schéma numérique et le schéma perturbé avec $\vec{\Phi}$ k -lipschitzienne.

$$\text{On a : } |\mathbf{Y}_i^* - \mathbf{Y}_i| = |\mathbf{Y}_{i-1}^* - \mathbf{Y}_{i-1} + h \cdot (\vec{\Phi}(\mathbf{Y}_j^*, t_j, h) - \vec{\Phi}(\mathbf{Y}_j, t_j, h)) + \zeta_{i-1}|$$

Or $\vec{\Phi}$ est k -lipschitzienne et par inégalité triangulaire, on obtient

Or $\vec{\Phi}$ est k -lipschitzienne et par inégalité triangulaire, on obtient

$$|Y_i^* - Y_i| \leq (1 + h.k).|Y_{i-1}^* - Y_{i-1}| + |\zeta_{i-1}|$$

Or $\vec{\Phi}$ est k -lipschitzienne et par inégalité triangulaire, on obtient

$$|Y_i^* - Y_i| \leq (1 + h.k) \cdot |Y_{i-1}^* - Y_{i-1}| + |\zeta_{i-1}|$$

Par récurrence simple, $|Y_i^* - Y_i| \leq (1 + h.k)^i \cdot |Y_0^* - Y_0| + \sum_{k=0}^{i-1} |\zeta_k|$

Or $\vec{\Phi}$ est k -lipschitzienne et par inégalité triangulaire, on obtient

$$|Y_i^* - Y_i| \leq (1 + h.k) \cdot |Y_{i-1}^* - Y_{i-1}| + |\zeta_{i-1}|$$

Par récurrence simple, $|Y_i^* - Y_i| \leq (1 + h.k)^i \cdot |Y_0^* - Y_0| + \sum_{k=0}^{i-1} |\zeta_k|$

Comme $1 + h.k > 1$, $1 \leq (1 + h.k)^i \leq (1 + h.k)^N$, ainsi

Or $\vec{\Phi}$ est k -lipschitzienne et par inégalité triangulaire, on obtient

$$|Y_i^* - Y_i| \leq (1 + h.k).|Y_{i-1}^* - Y_{i-1}| + |\zeta_{i-1}|$$

Par récurrence simple, $|Y_i^* - Y_i| \leq (1 + h.k)^i . |Y_0^* - Y_0| + \sum_{k=0}^{i-1} |\zeta_k|$

Comme $1 + h.k > 1$, $1 \leq (1 + h.k)^i \leq (1 + h.k)^N$, ainsi

$$\forall i, \quad |Y_i^* - Y_i| \leq M. (|Y_0^* - Y_0| + \sum_{i=0}^N |\zeta_{i}|)$$

Or $\vec{\Phi}$ est k -lipschitzienne et par inégalité triangulaire, on obtient

$$|Y_i^* - Y_i| \leq (1 + h.k).|Y_{i-1}^* - Y_{i-1}| + |\zeta_{i-1}|$$

Par récurrence simple, $|Y_i^* - Y_i| \leq (1 + h.k)^i . |Y_0^* - Y_0| + \sum_{k=0}^{i-1} |\zeta_k|$

Comme $1 + h.k > 1$, $1 \leq (1 + h.k)^i \leq (1 + h.k)^N$, ainsi

$$\forall i, \quad |Y_i^* - Y_i| \leq M. (|Y_0^* - Y_0| + \sum_{i=0}^N |\zeta_{i}|)$$

Donc la borne supérieure est également majorée par ce terme, on en déduit que le schéma est stable.

Intégration numérique

Le problème est d'estimer $\int_{t_i}^{t_{i+1}} \mathbf{F}(\mathbf{Y}, t) dt$ pour passer de \mathbf{Y}_i à \mathbf{Y}_{i+1}

$$\mathbf{Y}_{i+1} = \mathbf{Y}_i + \int_{t_i}^{t_{i+1}} \mathbf{F}(\mathbf{Y}, t) dt$$

$$\int_{t_i}^{t_{i+1}} \mathbf{F}(\mathbf{Y}, t) dt \approx h \cdot \vec{\Phi}(\mathbf{Y}_j, t_j, h) = h \cdot \sum_{j=0}^m \omega_j \cdot \mathbf{F}(\mathbf{Y}_j, \lambda(i, j))$$

Intégration numérique

Le problème est d'estimer $\int_{t_i}^{t_{i+1}} \mathbf{F}(\mathbf{Y}, t) dt$ pour passer de \mathbf{Y}_i à \mathbf{Y}_{i+1}

$$\mathbf{Y}_{i+1} = \mathbf{Y}_i + \int_{t_i}^{t_{i+1}} \mathbf{F}(\mathbf{Y}, t) dt$$

$$\int_{t_i}^{t_{i+1}} \mathbf{F}(\mathbf{Y}, t) dt \approx h \cdot \vec{\Phi}(\mathbf{Y}_j, t_j, h) = h \cdot \sum_{j=0}^m \omega_j \cdot \mathbf{F}(\mathbf{Y}_j, \lambda(i, j))$$

RAPPEL

$$\int_a^b f(x) \cdot dx = \sum_{i=0}^{n-1} \int_{x_i}^{x_{i+1}} f(x) \cdot dx \approx \frac{b-a}{n} \cdot \sum_{i=0}^{n-1} \sum_{j=0}^m \omega_j \cdot f(\lambda(i, j))$$

Intégration de f sur le segment $[t_i, t_{i+1}]$

Méthode	Expression	Poids	Noeuds

Intégration de f sur le segment $[t_i, t_{i+1}]$

Méthode	Expression	Poids	Noeuds
Rec. gauche	$h.f(t_i)$	$\omega_0 = 1$	$\lambda(i, 0) = t_i$

Intégration de f sur le segment $[t_i, t_{i+1}]$

Méthode	Expression	Poids	Noeuds
Rec. gauche	$h.f(t_i)$	$\omega_0 = 1$	$\lambda(i, 0) = t_i$
Rec. droite	$h.f(t_{i+1})$	$\omega_0 = 1$	$\lambda(i, 0) = t_{i+1}$

Intégration de f sur le segment $[t_i, t_{i+1}]$

Méthode	Expression	Poids	Noeuds
Rec. gauche	$h.f(t_i)$	$\omega_0 = 1$	$\lambda(i, 0) = t_i$
Rec. droite	$h.f(t_{i+1})$	$\omega_0 = 1$	$\lambda(i, 0) = t_{i+1}$
Point milieu	$h.f\left(\frac{t_i + t_{i+1}}{2}\right)$	$\omega_0 = 1$	$\lambda(i, 0) = \frac{t_i + t_{i+1}}{2} = t_{i+\frac{1}{2}}$

Intégration de f sur le segment $[t_i, t_{i+1}]$

Méthode	Expression	Poids	Noeuds
Rec. gauche	$h.f(t_i)$	$\omega_0 = 1$	$\lambda(i, 0) = t_i$
Rec. droite	$h.f(t_{i+1})$	$\omega_0 = 1$	$\lambda(i, 0) = t_{i+1}$
Point milieu	$h.f\left(\frac{t_i + t_{i+1}}{2}\right)$	$\omega_0 = 1$	$\lambda(i, 0) = \frac{t_i + t_{i+1}}{2} = t_{i+\frac{1}{2}}$
Trapèze	$\frac{h}{2} \cdot (f(t_i) + f(t_{i+1}))$	$\omega_0 = \omega_1 = \frac{1}{2}$	$\lambda(i, 0) = t_i; \lambda(i, 1) = t_{i+1}$

Intégration de $F(t, y)$

$$\int_{t_i}^{t_{i+1}} f(t).dt \approx h.\omega_i.f(\lambda(i, j)) \quad \text{et} \quad \int_{t_i}^{t_{i+1}} F(Y, t)dt \approx h.\vec{\Phi}(Y_j, t_j, h)$$

Intégration de $F(t, y)$

$$\int_{t_i}^{t_{i+1}} f(t).dt \approx h.\omega_i.f(\lambda(i, j)) \quad \text{et} \quad \int_{t_i}^{t_{i+1}} F(Y, t)dt \approx h.\vec{\Phi}(Y_j, t_j, h)$$

Méthode	Expression
Euler explicite	$h.F(Y_j, t_j)$
Euler implicite	$h.F(Y_{i+1}, t_{i+1})$
Euler semi-implicite	$h.((1 - \lambda).F(Y_j, t_j) + \lambda.F(Y_{i+1}, t_{i+1}))$
Euler semi-implicite	$\frac{h}{2}.(F(Y_j, t_j) + F(Y_{i+1}, t_{i+1}))$
Heun	$\frac{h}{2}.(F(Y_j, t_j) + F(Y_j + h.F(Y_j, t_j), t_{i+1}))$

Sommaire

- 11 Mise en place du problème
- 12 Annexe mathématique
- 13 Méthodes à un pas**
 - Méthode d'Euler explicite
 - Méthode d'Euler implicite
 - Méthode de Runge-Kutta
 - Conclusion
- 14 Mise en forme des systèmes d'équations différentielles
- 15 Résolution avec Scipy

$$\text{Support 1: } y(t) = -\frac{1}{t}$$

$$\text{Support 1: } y(t) = -\frac{1}{t}$$

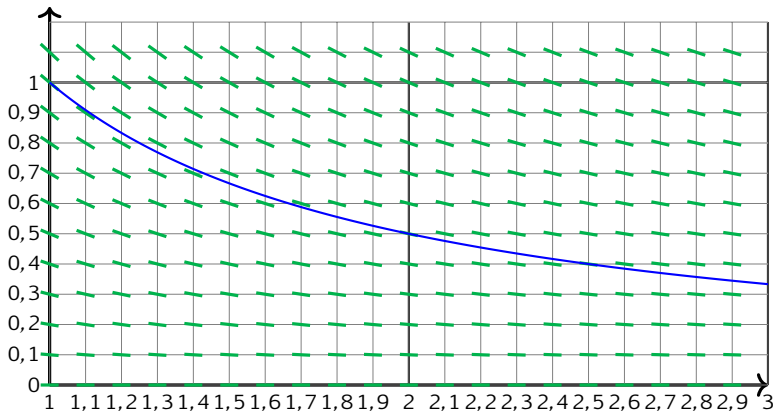
$$\frac{dy(t)}{dt} + \frac{y(t)}{t} = 0 \text{ avec } y(1) = 1$$

$$\text{Support 1: } y(t) = -\frac{1}{t}$$

$$\frac{dy(t)}{dt} + \frac{y(t)}{t} = 0 \text{ avec } y(1) = 1$$

dont la solution exacte est $y(t) = \frac{1}{t}$.

Support 1: $y(t) = -\frac{1}{t}$



Support 2: $y(t) = e^{-t}$

L'exemple élémentaire 2 qui servira de support est la résolution de l'équation différentielle :

Support 2: $y(t) = e^{-t}$

L'exemple élémentaire 2 qui servira de support est la résolution de l'équation différentielle :

$$\dot{y}(t) + y(t) = 0 \text{ avec } y(0) = 1$$

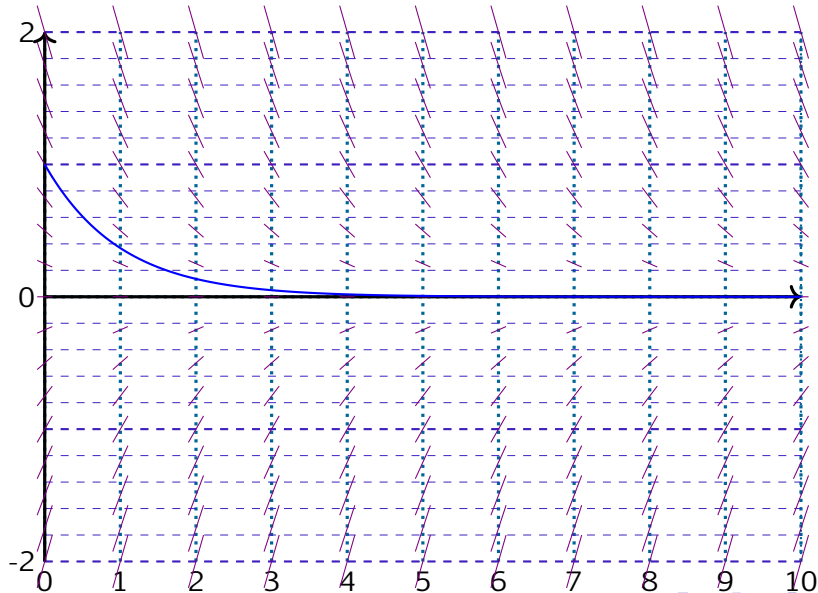
Support 2: $y(t) = e^{-t}$

L'exemple élémentaire 2 qui servira de support est la résolution de l'équation différentielle :

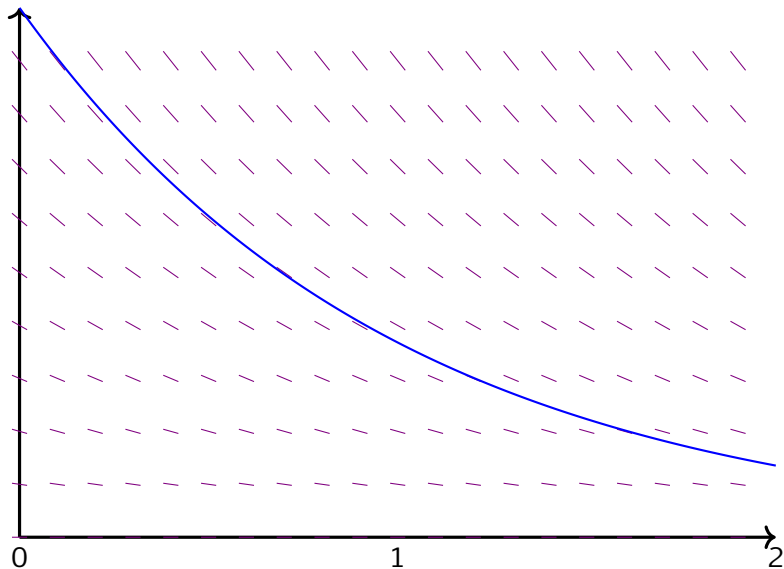
$$\dot{y}(t) + y(t) = 0 \text{ avec } y(0) = 1$$

dont la solution exacte est $y(t) = e^{-t}$.

Gradients



Gradients



Méthode d'Euler explicite

Méthode

La première façon d'approximer l'intégrale est de réaliser une méthode des rectangles à gauche.

$$\text{Ainsi } \int_{t_i}^{t_{i+1}} F(Y, t) dt = h.F(Y_i, t_i).$$

Méthode d'Euler explicite

Méthode

La première façon d'approximer l'intégrale est de réaliser une méthode des rectangles à gauche.

$$\text{Ainsi } \int_{t_i}^{t_{i+1}} F(Y, t) dt = h \cdot F(Y_i, t_i).$$

On en déduit que $\vec{\Phi} = F(Y_i, t_i)$

Méthode d'Euler explicite

Méthode

La première façon d'approximer l'intégrale est de réaliser une méthode des rectangles à gauche.

$$\text{Ainsi } \int_{t_i}^{t_{i+1}} \mathbf{F}(\mathbf{Y}, t) dt = h \cdot \mathbf{F}(\mathbf{Y}_i, t_i).$$

On en déduit que $\vec{\Phi} = \mathbf{F}(\mathbf{Y}_i, t_i)$

Ainsi la relation de récurrence générale est :

Méthode d'Euler explicite

Méthode

La première façon d'approximer l'intégrale est de réaliser une méthode des rectangles à gauche.

$$\text{Ainsi } \int_{t_i}^{t_{i+1}} F(Y, t) dt = h.F(Y_i, t_i).$$

On en déduit que $\vec{\Phi} = F(Y_i, t_i)$

Ainsi la relation de récurrence générale est :

$$Y_{i+1} = Y_i + h.F(Y_i, t_i)$$

Méthode d'Euler explicite

Méthode

La pente permettant de passer de Y_i à Y_{i+1} est la valeur de $F(Y_i, t_i)$ comme l'illustre la figure ci-contre où le champ de la fonction F est tracée.

Méthode d'Euler explicite

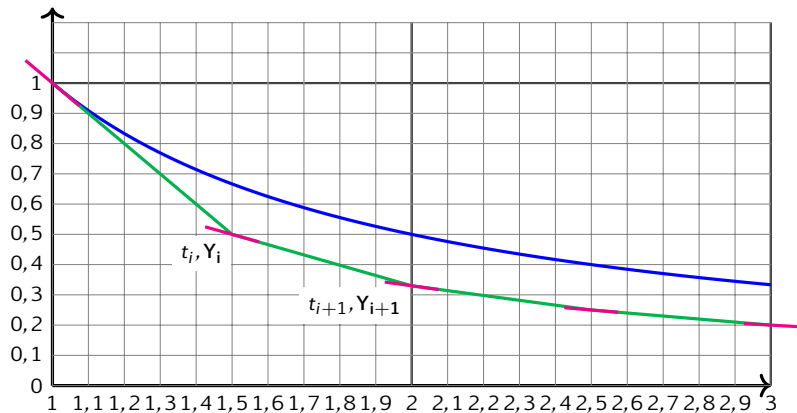
Méthode

La pente permettant de passer de Y_i à Y_{i+1} est la valeur de $F(Y_i, t_i)$ comme l'illustre la figure ci-contre où le champ de la fonction F est tracée.

On peut aussi exprimer la relation sous la forme :

$$F(Y_i, t_i) = \frac{dY}{dt}(t_i) \simeq \frac{Y_{i+1} - Y_i}{h}$$

Méthode d'Euler explicite



Propriétés

Stabilité

Propriétés

Stabilité si $\vec{\Phi}$ est lipschitzienne en Y le schéma est stable. Cela revient à vérifier que F est lipschitzienne en Y .

Consistance

Propriétés

Stabilité si $\vec{\Phi}$ est lipschitzienne en Y le schéma est stable. Cela revient à vérifier que F est lipschitzienne en Y .

Consistance L'erreur de consistance est

$$c_i = Y_{ex}(t_{i+1}) - Y_{ex}(t_i) - h \cdot F(Y_{ex}(t_i), t_i)$$

Or le développement de Taylor-Lagrange donne :

Propriétés

Stabilité si $\vec{\Phi}$ est lipschitzienne en Y le schéma est stable. Cela revient à vérifier que F est lipschitzienne en Y .

Consistance L'erreur de consistance est

$$c_i = Y_{ex}(t_{i+1}) - Y_{ex}(t_i) - h \cdot F(Y_{ex}(t_i), t_i)$$

Or le développement de Taylor-Lagrange donne :

$$Y_{ex}(t_i + h) = Y_{ex}(t_{i+1}) = Y_{ex}(t_i) + h \cdot \dot{Y}_{ex}(t_i) + \frac{h^2}{2} \cdot \ddot{Y}_{ex}(\xi)$$

avec $\xi \in [t_i, t_i + h]$.

Propriétés

Stabilité si $\vec{\Phi}$ est lipschitzienne en Y le schéma est stable. Cela revient à vérifier que F est lipschitzienne en Y .

Consistance L'erreur de consistance est

$$c_i = Y_{ex}(t_{i+1}) - Y_{ex}(t_i) - h \cdot F(Y_{ex}(t_i), t_i)$$

Or le développement de Taylor-Lagrange donne :

$$Y_{ex}(t_i + h) = Y_{ex}(t_{i+1}) = Y_{ex}(t_i) + h \cdot \dot{Y}_{ex}(t_i) + \frac{h^2}{2} \cdot \ddot{Y}_{ex}(\xi)$$

$$\text{avec } \xi \in [t_i, t_i + h].$$

$$\text{Or } \dot{Y}_{ex} = \vec{\Phi}(Y_{ex}(t_i), t_i, h) = F(Y_{ex}(t_i), t_i)$$

Propriétés

Stabilité si $\vec{\Phi}$ est lipschitzienne en Y le schéma est stable. Cela revient à vérifier que F est lipschitzienne en Y .

Consistance L'erreur de consistance est

$$c_i = Y_{ex}(t_{i+1}) - Y_{ex}(t_i) - h \cdot F(Y_{ex}(t_i), t_i)$$

Or le développement de Taylor-Lagrange donne :

$$Y_{ex}(t_i + h) = Y_{ex}(t_{i+1}) = Y_{ex}(t_i) + h \cdot \dot{Y}_{ex}(t_i) + \frac{h^2}{2} \cdot \ddot{Y}_{ex}(\xi)$$

$$\text{avec } \xi \in [t_i, t_i + h].$$

$$\text{Or } \dot{Y}_{ex} = \vec{\Phi}(Y_{ex}(t_i), t_i, h) = F(Y_{ex}(t_i), t_i)$$

$$\text{On en déduit que } c_i = \frac{h^2}{2} \cdot \ddot{Y}_{ex}(\xi) \text{ avec } \xi \in [t_i, t_i + h]$$

Ainsi la méthode est consistante. Comme elle est stable, on en déduit qu'elle est convergente.

L'expression de l'erreur de consistance indique que la méthode est d'ordre 1.

Application

Q - 1 : *Écrire la relation de récurrence de la méthode d'Euler sur l'exemple.*

Application

Q - 1 : *Écrire la relation de récurrence de la méthode d'Euler sur l'exemple.*

$$y'(t) + y(t) = 0 \quad \text{avec} \quad y(0) = 1$$

$$y'(t) = F(t, y(t)) \quad \text{avec} \quad F(t, y(t)) \approx F(t_i, y(t_i)) = -y(t_i)$$

Application

Q - 1 : Écrire la relation de récurrence de la méthode d'Euler sur l'exemple.

$$y'(t) + y(t) = 0 \quad \text{avec} \quad y(0) = 1$$

$$y'(t) = F(t, y(t)) \quad \text{avec} \quad F(t, y(t)) \approx F(t_i, y(t_i)) = -y(t_i)$$

Relation de récurrence :

$$y_{i+1} = y_i + h \cdot (-y_i)$$

Application

Q - 1 : Écrire la relation de récurrence de la méthode d'Euler sur l'exemple.

$$y'(t) + y(t) = 0 \quad \text{avec} \quad y(0) = 1$$

$$y'(t) = F(t, y(t)) \quad \text{avec} \quad F(t, y(t)) \approx F(t_i, y(t_i)) = -y(t_i)$$

Relation de récurrence :

$$y_{i+1} = y_i + h \cdot (-y_i) = y_i \cdot (1 - h)$$

Application

Q - 1 : *Écrire la relation de récurrence de la méthode d'Euler sur l'exemple.*

$$y'(t) + y(t) = 0 \quad \text{avec} \quad y(0) = 1$$

$$y'(t) = F(t, y(t)) \quad \text{avec} \quad F(t, y(t)) \approx F(t_i, y(t_i)) = -y(t_i)$$

Relation de récurrence :

$$y_{i+1} = y_i + h \cdot (-y_i) = y_i \cdot (1 - h) = y_0 \cdot (1 - h)^i$$

Q - 2 : *Montrer que pour certaines valeurs de h , la solution peut présenter des problèmes. Justifier et commenter les courbes de la Fig 15.*

Q - 2 : *Montrer que pour certaines valeurs de h , la solution peut présenter des problèmes. Justifier et commenter les courbes de la FIG 15.*

$$y_{i+1} = y_0 \cdot (1 - h)^i$$

Q - 2 : Montrer que pour certaines valeurs de h , la solution peut présenter des problèmes. Justifier et commenter les courbes de la Fig 15.

$$y_{i+1} = y_0 \cdot (1 - h)^i \quad \forall h \in]0, 2[, \quad \lim_{i \rightarrow \infty} y_n = 0$$

Q - 2 : Montrer que pour certaines valeurs de h , la solution peut présenter des problèmes. Justifier et commenter les courbes de la Fig 15.

$$y_{i+1} = y_0 \cdot (1 - h)^i \quad \forall h \in]0, 2[, \quad \lim_{i \rightarrow \infty} y_n = 0$$

ce qui est bien la limite de $y(t) = e^{-t}$, solution de l'équation différentielle.

Q - 2 : Montrer que pour certaines valeurs de h , la solution peut présenter des problèmes. Justifier et commenter les courbes de la Fig 15.

$$y_{i+1} = y_0 \cdot (1 - h)^i \quad \forall h \in]0, 2[, \quad \lim_{i \rightarrow \infty} y_n = 0$$

ce qui est bien la limite de $y(t) = e^{-t}$, solution de l'équation différentielle.

Aussi, pour $h = 1, \forall n \in \mathbb{N}/n \geq 1, y_i = 0$.

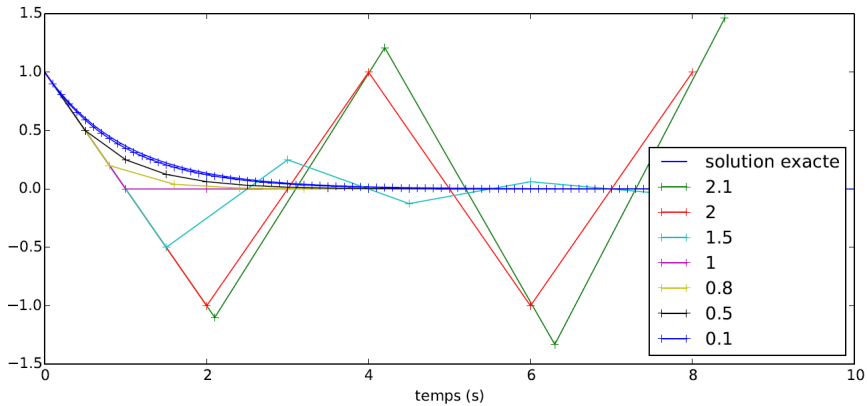
Q - 2 : Montrer que pour certaines valeurs de h , la solution peut présenter des problèmes. Justifier et commenter les courbes de la Fig 15.

$$y_{i+1} = y_0 \cdot (1 - h)^i \quad \forall h \in]0, 2[, \quad \lim_{i \rightarrow \infty} y_n = 0$$

ce qui est bien la limite de $y(t) = e^{-t}$, solution de l'équation différentielle.

Aussi, pour $h = 1, \forall n \in \mathbb{N}/n \geq 1, y_i = 0$.

Enfin si $h \in]-\infty, 0[\cup]2, \infty[$ la série diverge.



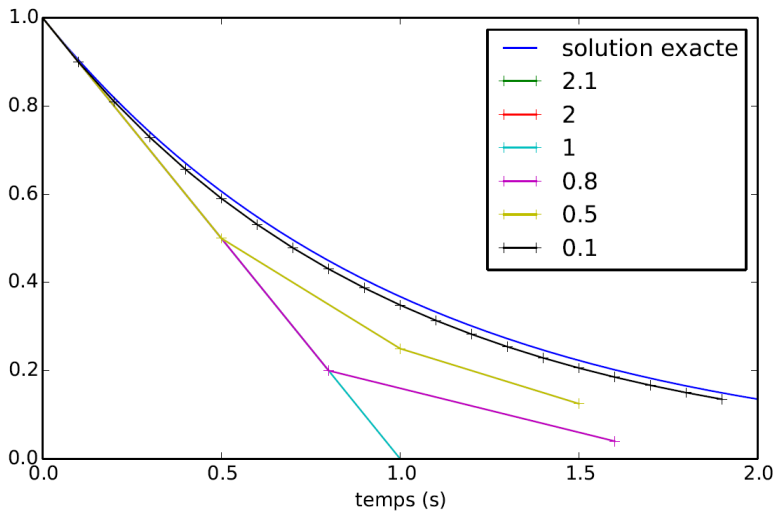


FIGURE 15 – Solution de l'application pour le schéma Euler explicite.

Le tableau suivant donne l'erreur comme étant le maximum de l'écart entre la solution exacte et la solution approchée ainsi que le temps de calcul pour différentes valeurs du pas de temps.

h	10^{-1}	10^{-2}	10^{-3}	10^{-4}	10^{-5}	10^{-6}
N pas de temps	100	1000	10000	100000	1000000	10000000
Erreur	$4,1 \cdot 10^{-1}$	$4,9 \cdot 10^{-2}$	$5,0 \cdot 10^{-3}$	$5,0 \cdot 10^{-04}$	$5,0 \cdot 10^{-5}$	$5,0 \cdot 10^{-06}$
Temps (s)	$5,2 \cdot 10^{-3}$	$1,6 \cdot 10^{-2}$	$7,8 \cdot 10^{-2}$	$5,7 \cdot 10^{-1}$	5,0	50

On constate que l'erreur évolue linéairement en fonction du pas de temps ainsi que le temps de calcul (complexité en $O(N)$ car une seule boucle).

Il faut descendre à un pas de temps très faible avant d'obtenir un niveau d'erreur satisfaisant. On voit bien ici la limitation des méthodes d'ordre 1.

Méthode d'Euler implicite

La seconde façon d'approximer l'intégrale est de réaliser une méthode des rectangles à droite.

$$\text{Ainsi } \int_{t_i}^{t_{i+1}} \mathbf{F}(\mathbf{Y}, t) \cdot dt = h \cdot \mathbf{F}(\mathbf{Y}_{i+1}, t_{i+1})$$

Méthode d'Euler implicite

La seconde façon d'approximer l'intégrale est de réaliser une méthode des rectangles à droite.

$$\text{Ainsi } \int_{t_i}^{t_{i+1}} \mathbf{F}(\mathbf{Y}, t) \cdot dt = h \cdot \mathbf{F}(\mathbf{Y}_{i+1}, t_{i+1})$$

On en déduit que $\vec{\Phi} = \mathbf{F}(\mathbf{Y}_{i+1}, t_{i+1})$

Méthode d'Euler implicite

La seconde façon d'approximer l'intégrale est de réaliser une méthode des rectangles à droite.

$$\text{Ainsi } \int_{t_i}^{t_{i+1}} \mathbf{F}(\mathbf{Y}, t) \cdot dt = h \cdot \mathbf{F}(\mathbf{Y}_{i+1}, t_{i+1})$$

On en déduit que $\vec{\Phi} = \mathbf{F}(\mathbf{Y}_{i+1}, t_{i+1})$

Ainsi la relation de récurrence générale est :

$$\mathbf{Y}_{i+1} = \mathbf{Y}_i + h \cdot \mathbf{F}(\mathbf{Y}_{i+1}, t_{i+1})$$

Méthode d'Euler implicite

La seconde façon d'approximer l'intégrale est de réaliser une méthode des rectangles à droite.

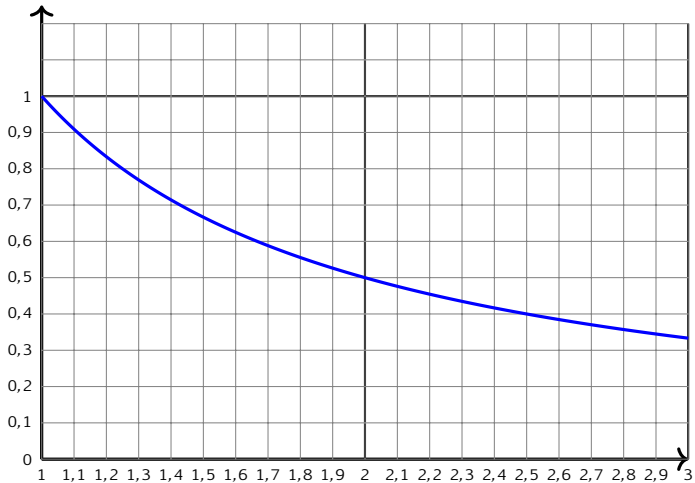
$$\text{Ainsi } \int_{t_i}^{t_{i+1}} \mathbf{F}(\mathbf{Y}, t) \cdot dt = h \cdot \mathbf{F}(\mathbf{Y}_{i+1}, t_{i+1})$$

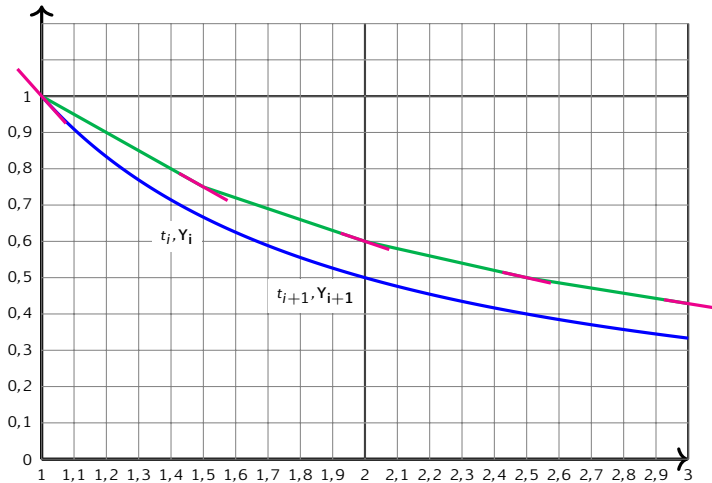
On en déduit que $\vec{\Phi} = \mathbf{F}(\mathbf{Y}_{i+1}, t_{i+1})$

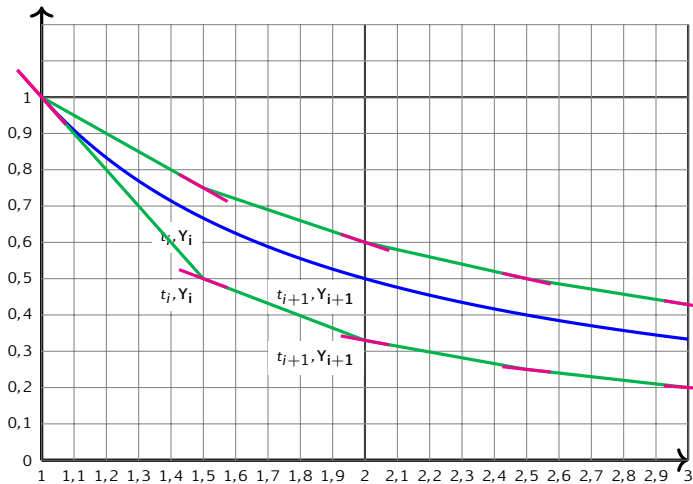
Ainsi la relation de récurrence générale est :

$$\mathbf{Y}_{i+1} = \mathbf{Y}_i + h \cdot \mathbf{F}(\mathbf{Y}_{i+1}, t_{i+1})$$

La pente permettant de passer de \mathbf{Y}_i à \mathbf{Y}_{i+1} est la valeur de $\mathbf{F}(\mathbf{Y}_{i+1}, t_{i+1})$ comme l'illustre la figure ci-contre où le champ de la fonction \mathbf{F} est tracée.







On peut aussi exprimer la relation sous la forme :

$$F(\mathbf{Y}_{i+1}, t_{i+1}) = \frac{d\mathbf{Y}}{dt}(t_{i+1}) \simeq \frac{\mathbf{Y}_{i+1} - \mathbf{Y}_i}{h}$$

On peut aussi exprimer la relation sous la forme :

$$F(Y_{i+1}, t_{i+1}) = \frac{dY}{dt}(t_{i+1}) \simeq \frac{Y_{i+1} - Y_i}{h}$$

On remarque que le second membre dépend aussi de Y_{i+1} . Il faut donc dans les cas non linéaires se ramener à un problème stationnaire à résoudre avec la méthode de Newton sur :

$$Y_{i+1} - Y_i - h.F(Y_{i+1}, t_{i+1}) = 0$$

ce qui alourdit notablement le calcul.

On peut aussi exprimer la relation sous la forme :

$$F(Y_{i+1}, t_{i+1}) = \frac{dY}{dt}(t_{i+1}) \simeq \frac{Y_{i+1} - Y_i}{h}$$

On remarque que le second membre dépend aussi de Y_{i+1} . Il faut donc dans les cas non linéaires se ramener à un problème stationnaire à résoudre avec la méthode de Newton sur :

$$Y_{i+1} - Y_i - h.F(Y_{i+1}, t_{i+1}) = 0$$

ce qui alourdit notablement le calcul.

Dans les cas linéaires, il faut au moins une inversion, ce qui est aussi très lourd pour les problèmes de grande dimension. Il faut donc retenir qu'une méthode implicite est généralement plus coûteuse qu'une méthode explicite à pas de temps égal.

Propriétés

De même que pour la méthode explicite, la méthode est convergente et d'ordre 1.

Application

Q - 3 : Écrire la relation de récurrence de la méthode d'Euler implicite sur l'exemple.

$$y'(t) + y(t) = 0 \quad \text{avec} \quad y(0) = 1$$

$$y'(t) = F(t, y(t)) \quad \text{avec} \quad F(t, y(t)) \approx F(t_{i+1}, y(t_{i+1})) = -y(t_{i+1})$$

Application

Q - 3 : Écrire la relation de récurrence de la méthode d'Euler implicite sur l'exemple.

$$y'(t) + y(t) = 0 \quad \text{avec} \quad y(0) = 1$$

$$y'(t) = F(t, y(t)) \quad \text{avec} \quad F(t, y(t)) \approx F(t_{i+1}, y(t_{i+1})) = -y(t_{i+1})$$

Relation de récurrence :

$$y_{i+1} = y_i + h \cdot (-y_{i+1})$$

Application

Q - 3 : Écrire la relation de récurrence de la méthode d'Euler implicite sur l'exemple.

$$y'(t) + y(t) = 0 \quad \text{avec} \quad y(0) = 1$$

$$y'(t) = F(t, y(t)) \quad \text{avec} \quad F(t, y(t)) \approx F(t_{i+1}, y(t_{i+1})) = -y(t_{i+1})$$

Relation de récurrence :

$$y_{i+1} = y_i + h \cdot (-y_{i+1}) \quad \Rightarrow \quad y_{i+1} = \frac{y_i}{1+h}$$

Application

Q - 3 : Écrire la relation de récurrence de la méthode d'Euler implicite sur l'exemple.

$$y'(t) + y(t) = 0 \quad \text{avec} \quad y(0) = 1$$

$$y'(t) = F(t, y(t)) \quad \text{avec} \quad F(t, y(t)) \approx F(t_{i+1}, y(t_{i+1})) = -y(t_{i+1})$$

Relation de récurrence :

$$y_{i+1} = y_i + h \cdot (-y_{i+1}) \quad \Rightarrow \quad y_{i+1} = \frac{y_i}{1+h} = \frac{y_0}{(1+h)^i}$$

Application

Q - 3 : Écrire la relation de récurrence de la méthode d'Euler implicite sur l'exemple.

$$y'(t) + y(t) = 0 \quad \text{avec} \quad y(0) = 1$$

$$y'(t) = F(t, y(t)) \quad \text{avec} \quad F(t, y(t)) \approx F(t_{i+1}, y(t_{i+1})) = -y(t_{i+1})$$

Relation de récurrence :

$$y_{i+1} = y_i + h \cdot (-y_{i+1}) \quad \Rightarrow \quad y_{i+1} = \frac{y_i}{1+h} = \frac{y_0}{(1+h)^i}$$

Q - 4 : Montrer que quelle que soit la valeur de h , la solution ne peut plus présenter les mêmes problèmes qu'en explicite.

Application

Q - 3 : Écrire la relation de récurrence de la méthode d'Euler implicite sur l'exemple.

$$y'(t) + y(t) = 0 \quad \text{avec} \quad y(0) = 1$$

$$y'(t) = F(t, y(t)) \quad \text{avec} \quad F(t, y(t)) \approx F(t_{i+1}, y(t_{i+1})) = -y(t_{i+1})$$

Relation de récurrence :

$$y_{i+1} = y_i + h \cdot (-y_{i+1}) \quad \Rightarrow \quad y_{i+1} = \frac{y_i}{1+h} = \frac{y_0}{(1+h)^i}$$

Q - 4 : Montrer que quelle que soit la valeur de h , la solution ne peut plus présenter les mêmes problèmes qu'en explicite.

La solution approchée pour différents pas de temps est donnée sur la FIG 16.

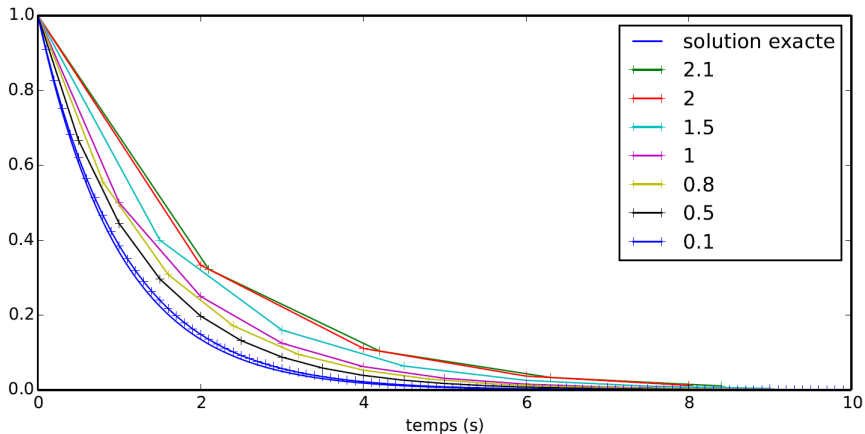


FIGURE 16 – Solution de l'application pour le schéma Euler implicite.

Le tableau suivant donne l'erreur comme étant le maximum de l'écart entre la solution exacte et la solution approchée ainsi que le temps de calcul pour différentes valeurs du pas de temps.

h	10^{-1}	10^{-2}	10^{-3}	10^{-4}	10^{-5}	10^{-6}
N pas de temps	100	1000	10000	100000	1000000	10000000
Erreur	$5,9 \cdot 10^{-1}$	$5,0 \cdot 10^{-2}$	$5,0 \cdot 10^{-3}$	$5,0 \cdot 10^{-4}$	$5,0 \cdot 10^{-5}$	$5,0 \cdot 10^{-6}$
Temps (s)	$5,2 \cdot 10^{-3}$	$3,1 \cdot 10^{-2}$	$3,7 \cdot 10^{-1}$	3,7	37	370

On constate que l'erreur évolue linéairement en fonction du pas de temps ainsi que le temps de calcul (complexité en $O(N)$ car une seule boucle).

On constate que l'erreur évolue linéairement en fonction du pas de temps ainsi que le temps de calcul (complexité en $O(N)$ car une seule boucle).

Le temps de calcul est cependant plus important à cause de la résolution de l'équation avec l'algorithme de Newton (on aurait pu inverser analytiquement la relation dans ce cas d'école).

On constate que l'erreur évolue linéairement en fonction du pas de temps ainsi que le temps de calcul (complexité en $O(N)$ car une seule boucle).

Le temps de calcul est cependant plus important à cause de la résolution de l'équation avec l'algorithme de Newton (on aurait pu inverser analytiquement la relation dans ce cas d'école).

Il faut descendre à un pas de temps très faible avant d'obtenir un niveau d'erreur satisfaisant. On voit bien ici la limitation des méthodes d'ordre 1.

Méthode de Runge-Kutta

Les méthodes d'Euler peuvent être vu comme une approximation du développement de Taylor-Lagrange à l'ordre 1.

Méthode de Runge-Kutta

Les méthodes d'Euler peuvent être vu comme une approximation du développement de Taylor-Lagrange à l'ordre 1.

Pour augmenter la précision, une solution pourrait être d'augmenter l'ordre de développement; c'est cependant impossible sans les expressions des dérivées n-ième de la fonction F .

Méthode de Runge-Kutta

Les méthodes d'Euler peuvent être vu comme une approximation du développement de Taylor-Lagrange à l'ordre 1.

Pour augmenter la précision, une solution pourrait être d'augmenter l'ordre de développement; c'est cependant impossible sans les expressions des dérivées n-ième de la fonction F .

Le schéma de Runge-Kutta permet d'obtenir des méthodes d'ordres plus élevés sans utiliser les fonctions dérivées de F .

Expression générale des méthodes de Runge-Kutta explicites

La méthode générale consiste à déterminer $\vec{\Phi}$ à partir d'approximations successives.

Expression générale des méthodes de Runge-Kutta explicites

La méthode générale consiste à déterminer $\vec{\Phi}$ à partir d'approximations successives.

On choisit

$$\vec{\Phi}(\mathbf{Y}, \tau, h) = \sum_{i=1}^q a_i \cdot k_i(\mathbf{Y}, \tau, h) \quad \text{avec} \quad k_1(\mathbf{Y}, \tau, h) = \mathbf{F}(\mathbf{Y}, \tau)$$

Expression générale des méthodes de Runge-Kutta explicites

La méthode générale consiste à déterminer $\vec{\Phi}$ à partir d'approximations successives.

On choisit

$$\vec{\Phi}(\mathbf{Y}, \tau, h) = \sum_{i=1}^q a_i \cdot k_i(\mathbf{Y}, \tau, h) \quad \text{avec} \quad k_1(\mathbf{Y}, \tau, h) = \mathbf{F}(\mathbf{Y}, \tau)$$

pour $i \geq 2$

$$k_i(\mathbf{Y}, \tau, h) = \mathbf{F} \left(\tau + \alpha_i \cdot h, \mathbf{Y} + \sum_{j=1}^{i-1} \beta_{ij} \cdot k_j(\mathbf{Y}, \tau, h) \right)$$

On choisit les paramètres q , a_i , α_i et β_{ij} pour que la méthode soit d'ordre p .

On remarque que la méthode d'Euler explicite est le cas particulier de Runge-Kutta d'ordre 1.

En pratique, on va rarement au delà de l'ordre 4.

Runge-Kutta d'ordre 2 (RK2) – Méthode de Heun

Le relation de récurrence est

$$Y_{i+1} = Y_i + \frac{h}{2} \cdot F(Y_i, t_i) + \frac{h}{2} \cdot F(Y_i + h \cdot F(t_i, Y_i), t_i + h)$$

Runge-Kutta d'ordre 2 (RK2) – Méthode de Heun

Le relation de récurrence est

$$Y_{i+1} = Y_i + \frac{h}{2} \cdot F(Y_i, t_i) + \frac{h}{2} \cdot F(Y_i + h \cdot F(t_i, Y_i), t_i + h)$$

On a donc pris $\vec{\Phi}(Y, \tau, h) = \frac{k_1 + k_2}{2}$ avec :

Runge-Kutta d'ordre 2 (RK2) – Méthode de Heun

Le relation de récurrence est

$$Y_{i+1} = Y_i + \frac{h}{2} \cdot F(Y_i, t_i) + \frac{h}{2} \cdot F(Y_i + h \cdot F(t_i, Y_i), t_i + h)$$

On a donc pris $\vec{\Phi}(Y, \tau, h) = \frac{k_1 + k_2}{2}$ avec :

$$k_1 = F(Y_i, t_i) \quad \text{et} \quad k_2 = F(Y_i + h \cdot k_1, t_i + h)$$

Runge-Kutta d'ordre 2 (RK2) – Méthode de Heun

Le relation de récurrence est

$$Y_{i+1} = Y_i + \frac{h}{2} \cdot F(Y_i, t_i) + \frac{h}{2} \cdot F(Y_i + h \cdot F(t_i, Y_i), t_i + h)$$

On a donc pris $\vec{\Phi}(Y, \tau, h) = \frac{k_1 + k_2}{2}$ avec :

$$k_1 = F(Y_i, t_i) \quad \text{et} \quad k_2 = F(Y_i + h \cdot k_1, t_i + h)$$

Cette relation est une intégration de type trapèze entre la valeur au pas i et la valeur du pas $i + 1$ estimée à partir d'une méthode d'Euler explicite.

Runge-Kutta d'ordre 2 (RK2) – Méthode de Heun

Le relation de récurrence est

$$Y_{i+1} = Y_i + \frac{h}{2} \cdot F(Y_i, t_i) + \frac{h}{2} \cdot F(Y_i + h \cdot F(t_i, Y_i), t_i + h)$$

On a donc pris $\vec{\Phi}(Y, \tau, h) = \frac{k_1 + k_2}{2}$ avec :

$$k_1 = F(Y_i, t_i) \quad \text{et} \quad k_2 = F(Y_i + h \cdot k_1, t_i + h)$$

Cette relation est une intégration de type trapèze entre la valeur au pas i et la valeur du pas $i + 1$ estimée à partir d'une méthode d'Euler explicite.

La solution approchée de l'application pour différents pas de temps est donnée sur la Fig 17.

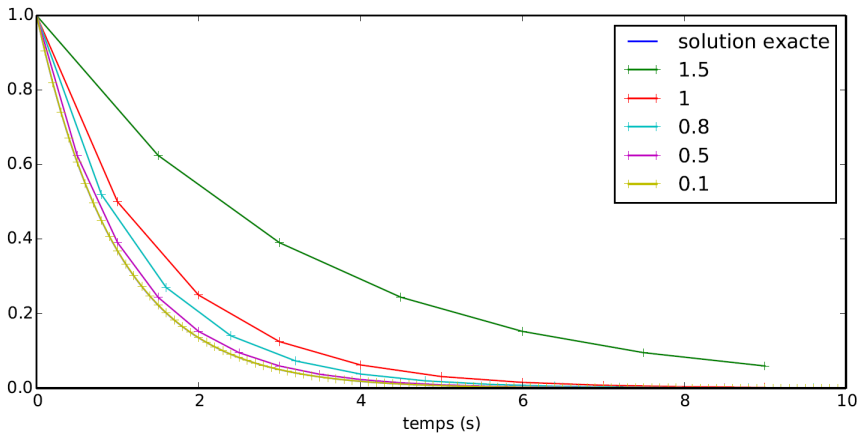


FIGURE 17 – Solution de l'application pour le schéma de Heun.

Le tableau suivant donne l'erreur comme étant le maximum de l'écart entre la solution exacte et la solution approchée ainsi que le temps de calcul pour différentes valeurs du pas de temps.

h	10^{-1}	10^{-2}	10^{-3}	10^{-4}	10^{-5}	10^{-6}
N pas de temps	100	1000	10000	100000	1000000	10000000
Erreur	$1,8 \cdot 10^{-2}$	$1,7 \cdot 10^{-4}$	$1,7 \cdot 10^{-6}$	$1,7 \cdot 10^{-8}$	$1,7 \cdot 10^{-10}$	$1,7 \cdot 10^{-12}$
Temps (s)	$1,5 \cdot 10^{-3}$	$1,4 \cdot 10^{-2}$	$1,3 \cdot 10^{-1}$	1,2	12	120

On constate que l'erreur évolue de façon quadratique en fonction du pas de temps : en divisant le pas de temps par 10, l'erreur diminue d'un facteur 100.

Le temps de calcul est toujours linéaire car la complexité n'a pas augmenté (complexité en $O(N)$). Il est environ double de la méthode d'Euler explicite puisqu'il y a deux évaluations de $F(Y, t)$.

Runge-Kutta d'ordre 4 (RK4)

La fonction $\vec{\Phi}$ est $\vec{\Phi}(\mathbf{Y}, \tau, h) = \frac{1}{6} \cdot (k_1 + 2.k_2 + 2.k_3 + k_4)$ avec quatre évaluations successives de F :

$$k_1(\mathbf{Y}, \tau, h) = F(\mathbf{Y}, \tau)$$

$$k_2(\mathbf{Y}, \tau, h) = F\left(\tau + \frac{h}{2}, \mathbf{Y} + \frac{h}{2} \cdot k_1(\mathbf{Y}, \tau, h)\right)$$

$$k_3(\mathbf{Y}, \tau, h) = F\left(\tau + \frac{h}{2}, \mathbf{Y} + \frac{h}{2} \cdot k_2(\mathbf{Y}, \tau, h)\right)$$

$$k_4(\mathbf{Y}, \tau, h) = F(\tau + h, \mathbf{Y} + h \cdot k_3(\mathbf{Y}, \tau, h))$$

La solution approchée pour différents pas de temps est donnée sur la Fig 18.

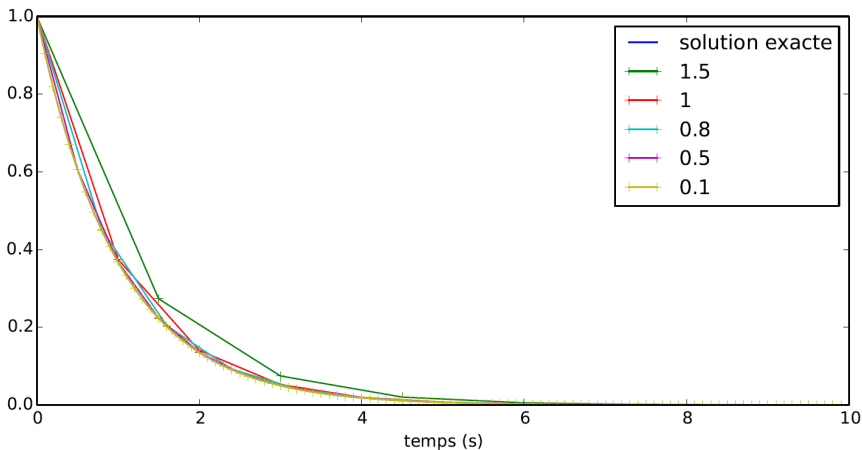


FIGURE 18 – Solution de l'application pour le schéma de Runge-Kutta d'ordre 4.

Le tableau suivant donne l'erreur comme étant le maximum de l'écart entre la solution exacte et la solution approchée ainsi que le temps de calcul pour différentes valeurs du pas de temps.

h	10^{-1}	10^{-2}	10^{-3}
N pas de temps	100	1000	10000
Erreur	$9,0 \cdot 10^{-6}$	$8,4 \cdot 10^{-10}$	$9,2 \cdot 10^{-14}$
Temps (s)	$1,7 \cdot 10^{-3}$	$1,6 \cdot 10^{-2}$	$2,5 \cdot 10^{-1}$

On constate que l'erreur évolue à l'ordre 4 en fonction du pas de temps :

Le tableau suivant donne l'erreur comme étant le maximum de l'écart entre la solution exacte et la solution approchée ainsi que le temps de calcul pour différentes valeurs du pas de temps.

h	10^{-1}	10^{-2}	10^{-3}
N pas de temps	100	1000	10000
Erreur	$9,0 \cdot 10^{-6}$	$8,4 \cdot 10^{-10}$	$9,2 \cdot 10^{-14}$
Temps (s)	$1,7 \cdot 10^{-3}$	$1,6 \cdot 10^{-2}$	$2,5 \cdot 10^{-1}$

On constate que l'erreur évolue à l'ordre 4 en fonction du pas de temps : ainsi, en divisant le pas de temps par 10, l'erreur diminue d'un facteur 10000.

Le tableau suivant donne l'erreur comme étant le maximum de l'écart entre la solution exacte et la solution approchée ainsi que le temps de calcul pour différentes valeurs du pas de temps.

h	10^{-1}	10^{-2}	10^{-3}
N pas de temps	100	1000	10000
Erreur	$9,0 \cdot 10^{-6}$	$8,4 \cdot 10^{-10}$	$9,2 \cdot 10^{-14}$
Temps (s)	$1,7 \cdot 10^{-3}$	$1,6 \cdot 10^{-2}$	$2,5 \cdot 10^{-1}$

On constate que l'erreur évolue à l'ordre 4 en fonction du pas de temps : ainsi, en divisant le pas de temps par 10, l'erreur diminue d'un facteur 10000.

Le temps de calcul est toujours linéaire car la complexité n'a pas augmentée (complexité en $O(N)$).

Conclusion

En terme de schéma explicite, la méthode de Runge-Kutta à l'ordre 4 est souvent utilisée.

Concernant le choix du pas de temps, il obéit à un compromis entre le temps de calcul, la stabilité et le stockage en mémoire.

Conclusion

En terme de schéma explicite, la méthode de Runge-Kutta à l'ordre 4 est souvent utilisée.

Concernant le choix du pas de temps, il obéit à un compromis entre le temps de calcul, la stabilité et le stockage en mémoire.

Le choix de ce pas de temps dépendra essentiellement de la dynamique du système à représenter : si on veut représenter une dynamique à 10 kHz, il faudrait prendre un pas de temps de l'ordre de 10^{-5} (10 points par période).

Le choix du pas de temps est souvent fait de manière empirique par expérience de l'ingénieur.

Sommaire

- 11 Mise en place du problème
- 12 Annexe mathématique
- 13 Méthodes à un pas
- 14 Mise en forme des systèmes d'équations différentielles**
 - Équation harmonique
 - Système masse-ressort-amortisseur entretenu
 - Création du glycol
- 15 Résolution avec Scipy

Équation harmonique

Prenons un système masse-ressort.

Équation harmonique

Prenons un système masse-ressort.

Son équation du mouvement se met sous la forme :

$$\ddot{y}(t) + \omega_0^2 \cdot y(t) = 0$$

Équation harmonique

Prenons un système masse-ressort.

Son équation du mouvement se met sous la forme :

$$\ddot{y}(t) + \omega_0^2 \cdot y(t) = 0$$

Les conditions initiales sont $y(0) = 1$ et $\dot{y}(0) = 0$.

Équation harmonique

Prenons un système masse-ressort.

Son équation du mouvement se met sous la forme :

$$\ddot{y}(t) + \omega_0^2 \cdot y(t) = 0$$

Les conditions initiales sont $y(0) = 1$ et $\dot{y}(0) = 0$.

Q - 3 : Mettre ce problème sous la forme du problème de Cauchy.

Équation harmonique

Prenons un système masse-ressort.

Son équation du mouvement se met sous la forme :

$$\ddot{y}(t) + \omega_0^2 \cdot y(t) = 0$$

Les conditions initiales sont $y(0) = 1$ et $\dot{y}(0) = 0$.

Q - 3 : Mettre ce problème sous la forme du problème de Cauchy.

L'idée est d'introduire des variables supplémentaires. Ici, on introduit la fonction $y_2(t)$ telle qu'elle soit solution de $y_2(t) = \dot{y}_1(t)$ avec $y_1(t) = y(t)$.

L'équation différentielle du système masse-ressort s'écrit ainsi :

L'équation différentielle du système masse-ressort s'écrit ainsi :

$$\dot{y}_2(t) = -\omega_0^2 \cdot y_1(t)$$

L'équation différentielle du système masse-ressort s'écrit ainsi :

$$\dot{y}_2(t) = -\omega_0^2 \cdot y_1(t)$$

En posant

$$\vec{Y}(t) = \begin{pmatrix} y_1(t) \\ y_2(t) \end{pmatrix}, \quad \vec{F}(\vec{Y}, t) = \begin{pmatrix} y_2(t) \\ -\omega_0^2 \cdot y_1(t) \end{pmatrix} \quad \text{et} \quad \vec{Y}_0 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

L'équation différentielle du système masse-ressort s'écrit ainsi :

$$\dot{y}_2(t) = -\omega_0^2 \cdot y_1(t)$$

En posant

$$\vec{Y}(t) = \begin{pmatrix} y_1(t) \\ y_2(t) \end{pmatrix}, \quad \vec{F}(\vec{Y}, t) = \begin{pmatrix} y_2(t) \\ -\omega_0^2 \cdot y_1(t) \end{pmatrix} \quad \text{et} \quad \vec{Y}_0 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

on obtient bien la forme du problème de Cauchy.

La résolution se fait de la même manière que précédemment sauf que cette fois les grandeurs manipulées sont des vecteurs et non plus des scalaires.

Une autre solution classique permettant de résoudre cette équation différentielle avec la méthode d'Euler explicite est d'appliquer deux fois la définition.

On note Y_p la dérivée première. On a $Y_p = \frac{Y_{i+1} - Y_i}{h}$

Une autre solution classique permettant de résoudre cette équation différentielle avec la méthode d'Euler explicite est d'appliquer deux fois la définition.

On note Y_p la dérivée première. On a $Y_{p_i} = \frac{Y_{i+1} - Y_i}{h}$

On note Y_{pp} la dérivée seconde.

$$Y_{pp_i} = \frac{Y_{p_{i+1}} - Y_{p_i}}{h} = \frac{Y_{i+2} - 2.Y_{i+1} + Y_i}{h^2}$$

Une autre solution classique permettant de résoudre cette équation différentielle avec la méthode d'Euler explicite est d'appliquer deux fois la définition.

On note Yp la dérivée première. On a $Yp_i = \frac{Y_{i+1} - Y_i}{h}$

On note Ypp la dérivée seconde.

$$Ypp_i = \frac{Yp_{i+1} - Yp_i}{h} = \frac{Y_{i+2} - 2.Y_{i+1} + Y_i}{h^2}$$

On en déduit que

$$Y_{i+2} = (-h^2 \cdot \omega_0^2 - 1) \cdot Y_i + 2 \cdot Y_{i+1}$$

On obtient une relation de récurrence directe qui peut être programmée.

Une autre solution classique permettant de résoudre cette équation différentielle avec la méthode d'Euler explicite est d'appliquer deux fois la définition.

On note Yp la dérivée première. On a $Yp_i = \frac{Y_{i+1} - Y_i}{h}$

On note Ypp la dérivée seconde.

$$Ypp_i = \frac{Yp_{i+1} - Yp_i}{h} = \frac{Y_{i+2} - 2.Y_{i+1} + Y_i}{h^2}$$

On en déduit que

$$Y_{i+2} = (-h^2 \cdot \omega_0^2 - 1) \cdot Y_i + 2 \cdot Y_{i+1}$$

On obtient une relation de récurrence directe qui peut être programmée.

Système masse-ressort-amortisseur entretenu

Soit un système masse-ressort-amortisseur soumis à une excitation extérieure harmonique.

Système masse-ressort-amortisseur entretenu

Soit un système masse-ressort-amortisseur soumis à une excitation extérieure harmonique.

L'équation différentielle se met sous la forme générale :

Système masse-ressort-amortisseur entretenu

Soit un système masse-ressort-amortisseur soumis à une excitation extérieure harmonique.

L'équation différentielle se met sous la forme générale :

$$\ddot{y}(t) + 2.\xi.\omega_0.\dot{y}(t) + \omega_0^2.y(t) = f_0.\sin(\omega t)$$

Système masse-ressort-amortisseur entretenu

Soit un système masse-ressort-amortisseur soumis à une excitation extérieure harmonique.

L'équation différentielle se met sous la forme générale :

$$\ddot{y}(t) + 2.\xi.\omega_0.\dot{y}(t) + \omega_0^2.y(t) = f_0.\sin(\omega t)$$

Les conditions initiales sont $y(0) = 1$ et $\dot{y}(0) = 0$ (solution Fig 19(a)).

Système masse-ressort-amortisseur entretenu

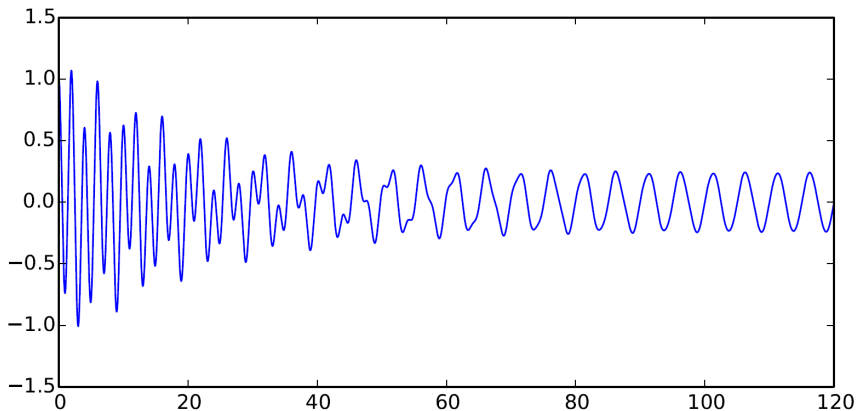
Soit un système masse-ressort-amortisseur soumis à une excitation extérieure harmonique.

L'équation différentielle se met sous la forme générale :

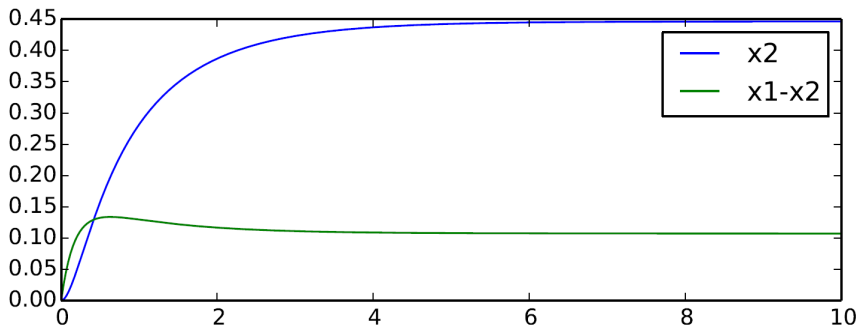
$$\ddot{y}(t) + 2.\xi.\omega_0.\dot{y}(t) + \omega_0^2.y(t) = f_0.\sin(\omega t)$$

Les conditions initiales sont $y(0) = 1$ et $\dot{y}(0) = 0$ (solution Fig 19(a)).

Q - 5 : *Mettre ce problème sous la forme du problème de Cauchy.*



(a) Solution de l'équation différentielle du système masse-ressort-amortisseur pour $\omega_0 = 2$ rad/s, $\xi = 0,015$, $f_0 = 2$ N et $\omega = 5$ rad/s.

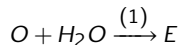


(b) Solution de l'équation différentielle de la réaction de création du glycol.

FIGURE 19 – Solutions des exemples 2 et 3.

Création du glycol

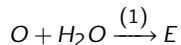
Le glycol (noté E) résulte de l'addition d'eau à l'oxyde d'éthylène (noté O) en phase gazeuse, selon la réaction



de constante de vitesse k_1 .

Création du glycol

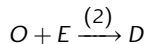
Le glycol (noté E) résulte de l'addition d'eau à l'oxyde d'éthylène (noté O) en phase gazeuse, selon la réaction



de constante de vitesse k_1 .

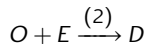
Cette réaction est effectuée à 473 K sous une pression $P = 15,0$ bar.

Industriellement le temps de passage dans le réacteur ne permet pas d'atteindre l'état d'équilibre thermodynamique et on constate l'apparition de diéthylèneglycol (noté D) produit par la réaction



se déroulant également en phase gazeuse, de constante de vitesse k_2 .

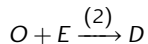
Industriellement le temps de passage dans le réacteur ne permet pas d'atteindre l'état d'équilibre thermodynamique et on constate l'apparition de diéthylèneglycol (noté D) produit par la réaction



se déroulant également en phase gazeuse, de constante de vitesse k_2 .

Les réactions sont supposées d'ordre un par rapport à chacun des réactifs et totales.

Industriellement le temps de passage dans le réacteur ne permet pas d'atteindre l'état d'équilibre thermodynamique et on constate l'apparition de diéthylèneglycol (noté D) produit par la réaction



se déroulant également en phase gazeuse, de constante de vitesse k_2 .

Les réactions sont supposées d'ordre un par rapport à chacun des réactifs et totales.

Pour traduire le fait que l'eau réagit moins vite que le glycol E sur l'oxyde d'éthylène O , les constantes de vitesse k_1 et k_2 sont choisies telles que $k_2 = 5.k_1$.

Le mélange initial est constitué d'oxyde d'éthylène et d'eau à la concentration molaire $c_0 = 1$ mol/L chacun. On notera x_1 l'avancement de la réaction (1) en mol/L et x_2 l'avancement de la réaction (2) en mol/L

Les équations différentielles modélisant le mécanisme réactionnel sont (solution FIG 19(b)):

Q - 34 : *Mettre ce problème sous la forme du problème de Cauchy.*

$$\begin{cases} \frac{dx_1}{dt} = k_1 \cdot (c_0 - x_1 - x_2) \cdot (c_0 - x_1) \\ \frac{dx_2}{dt} = k_2 \cdot (c_0 - x_1 - x_2) \cdot (x_1 - x_2) \end{cases}$$

Sommaire

- 11 Mise en place du problème
- 12 Annexe mathématique
- 13 Méthodes à un pas
- 14 Mise en forme des systèmes d'équations différentielles
- 15 Résolution avec Scipy**

OBJECTIF : Résoudre une équation différentielle du type $\vec{Y}' = \vec{F}(\vec{Y}, t)$.

Le module `scipy.integrate` propose la fonction `odeint` telle que `odeint(F, Y0, T)` permet d'obtenir une estimation de la solution, pour chaque valeur de T , de l'équation différentielle $Y' = F(Y, t)$ en partant de $T[0]$ avec comme valeur de l'état initial $Y0$. L'état du système peut contenir N variables avec $N \in \mathbb{N}^*$.

En ramenant l'équation différentielle scalaire du second ordre (1), à une équation différentielle vectorielle du premier ordre (2), le problème de Cauchy initial (3) peut être traité avec `odeint`.

$$\begin{cases} \frac{1}{\omega_0^2} \cdot y''(t) + \frac{2 \cdot \xi}{\omega_0} \cdot y'(t) + y(t) = K \cdot (u(t - t_0) - u(t - t_1)) \\ y(t_0) = y_0 \text{ et } y'(t_0) = v_0 \end{cases} \quad (1)$$

avec $u(t)$ l'échelon unitaire, nul si t est négatif, égale à 1 sinon.

$$\begin{cases} \frac{dy(t)}{dt} = v(t) \\ \frac{dv(t)}{dt} = (K.(u(t-t_0) - u(t-t_1)) - y(t)).\omega_0^2 - 2.\xi.\omega_0.v(t) \\ y(t_0) = y_0 \\ v(t_0) = v_0 \end{cases} \quad (2)$$

$$\begin{cases} \frac{dY(t)}{dt} = F(Y, t) \\ Y(t_0) = [y_0, v_0] \end{cases} \quad (3)$$

avec $Y(t) = [y(t), v(t)]$

```
def sys_amort(Y, t):
    if t < t0 or t > t1:
        u = 0
    else:
        u = 1
    return [Y[1], (K*u - Y[0]) * w0**2 - 2*xi*w0*Y[1]]
```

```
K, xi, w0, t0, t1, y0, v0 = 3, 0.25, 20, 0, 1, 0, 0
T = np.linspace(t0-1, t1+1, 200)
sim = scint.odeint(sys_amort, [y0, v0], T)
```

- fonction associée au problème de Cauchy
- définition des paramètres et simulations
- présentation des résultats

```
plt.figure(0)
plt.title(r'Déplacement $y(t)$')
plt.plot(T, sim[:,0])
plt.figure(1)
plt.title(r'Vitesse $v(t)$')
plt.plot(T, sim[:,1])
```

Quatrième partie IV

Pivot de Gauss et algèbre linéaire

Sommaire

16 Introduction

- Système d'équations linéaires
- Méthode directe
- Quelques facteurs influents

17 Méthode de Gauss

18 Résoudre un système système de Cramer avec Numpy

Système d'équations linéaires

Dans les domaines du calcul de structure, de la simulation multiphysique ou encore en métrologie, il est fréquent d'avoir à résoudre des systèmes d'équations linéaires.

Système d'équations linéaires

Dans les domaines du calcul de structure, de la simulation multiphysique ou encore en métrologie, il est fréquent d'avoir à résoudre des systèmes d'équations linéaires.

L'algorithme de Gauss permet de résoudre ces systèmes.

Système d'équations linéaires

Dans les domaines du calcul de structure, de la simulation multiphysique ou encore en métrologie, il est fréquent d'avoir à résoudre des systèmes d'équations linéaires.

L'algorithme de Gauss permet de résoudre ces systèmes.

Soit un système de n équations linéaires, on peut l'écrire sous la forme matricielle $A.X = B$:

$$\begin{cases} a_{11}.x_1 + a_{12}.x_2 + \dots + a_{1n}.x_n & = & b_1 \\ a_{21}.x_1 + a_{22}.x_2 + \dots + a_{2n}.x_n & = & b_2 \\ & \vdots & \\ a_{n1}.x_1 + a_{n2}.x_2 + \dots + a_{nn}.x_n & = & b_n \end{cases} \Rightarrow \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$$

Méthode directe

En utilisant la méthode de Gauss, l'objectif est de résoudre ce système en *une fois*. Ce type de méthode est appelé **directe** en opposition aux méthodes **itératives** pour lesquelles on répète des opérations jusqu'à ce que le résultat converge vers une solution.

Méthode directe

En utilisant la méthode de Gauss, l'objectif est de résoudre ce système en *une fois*. Ce type de méthode est appelé **directe** en opposition aux méthodes **itératives** pour lesquelles on répète des opérations jusqu'à ce que le résultat converge vers une solution.

Les méthodes directes sont encore majoritairement employées dans les codes actuels car elles sont robustes et que l'on est capable de déterminer à l'avance leur coût (nombre d'opérations, stockage nécessaire). Cependant, le coût pour la méthode de Gauss augmente fortement pour des systèmes de grande taille.

Quelques facteurs influents

Les systèmes auxquels on va s'intéresser dans la suite sont supposés **linéaires** et **inversibles**.

Si la dimension de la matrice $\dim(A) = n$ est importante le nombre de données à stocker est très élevé (n^2). Afin de réduire les ressources mémoire nécessaires, le stockage des données devra être optimisé.

De plus, lors du stockage numérique de ces données, chaque nombre est codé sur un nombre fini de bits. Selon le principe de la *virgule flottante*, des approximations sont donc commises. Or pour les systèmes de grande taille, ou mal conditionnés, la méthode de Gauss a tendance à amplifier les erreurs d'arrondi (d'où une instabilité numérique).

Sommaire

16 Introduction

17 Méthode de Gauss

- Principe du pivot
- Inversion de matrice

18 Résoudre un système système de Cramer avec Numpy

Transvection

Pour appliquer la méthode de Gauss, on procède par opérations élémentaires sur les lignes L_i d'un système d'équations. Or la solution d'un système linéaire ne change pas si :

- on multiplie tous les termes d'une équation par une constante **non nulle** :

$$L_i \leftarrow \lambda.L_i \quad \text{avec} \quad \lambda \neq 0$$

- on remplace une équation par la somme, membre à membre de **cette équation** et d'une autre équation du système :

$$L_i \leftarrow L_i + \lambda.L_j$$

On construit alors la matrice de transvection $T_{i,j}(\lambda)$, de dimension n , définie par :

$$[T_{i,j}(\lambda)]_{k,l} = \begin{cases} 1 & \text{si } k = l \\ \lambda & \text{si } k = i \text{ et } l = j \\ 0 & \text{sinon} \end{cases}$$

$$\Leftrightarrow T_{i,j}(\lambda) = L_i \begin{matrix} & & & C_j \\ \left[\begin{array}{cccc} 1 & & & 0 \\ & 1 & & \lambda \\ & & \ddots & \\ & 0 & & 1 \\ & & & & 1 \end{array} \right] & = \mathbb{I}_n + \lambda.E_{i,j}$$

avec \mathbb{I}_n , la matrice identité de dimension n et $E_{i,j}$ la matrice constituée de 0 sauf ligne i et colonne j où elle admet 1 comme coefficient.

Triangularisation

L'idée principale de la méthode de Gauss est d'obtenir, par opérations élémentaires, un système triangulaire d'équations, où la ligne L_i s'exprime en fonction de i (resp. $n - i + 1$) inconnues dont une seule n'est pas apparue (resp. n'apparaît pas) dans les lignes précédentes (resp. suivantes) :

$$\begin{cases} a'_{11} \cdot x_1 + a'_{12} \cdot x_2 + \cdots + a'_{1n} \cdot x_n & = & b'_1 \\ & a'_{22} \cdot x_2 + \cdots + a'_{2n} \cdot x_n & = & b'_2 \\ & & \ddots & \vdots & = & \vdots \\ & & & a'_{nn} \cdot x_n & = & b'_n \end{cases} \Rightarrow \begin{bmatrix} a'_{11} & a'_{12} & \cdots & a'_{1n} \\ & a'_{22} & \cdots & a'_{2n} \\ & & \ddots & \vdots \\ & & & a'_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b'_1 \\ b'_2 \\ \vdots \\ b'_n \end{bmatrix}$$

Pour triangulariser, il faut éliminer des coefficients sous la diagonale de la matrice. A l'étape i :

$$\forall k \in \llbracket i + 1, n \rrbracket, L_k \leftarrow L_k - \frac{a_{k,i}}{a_{i,i}} . L_i$$

On utilise donc des **pivots**, les coefficients $a_{i,i}$.

Il apparaît donc un problème mathématiquement si $a_{i,i} = 0$ et numériquement si $a_{i,i} \ll 1$.

Pour éviter ce problème deux solutions existent:

- méthode du pivot partiel : permutation des lignes ou des colonnes \Rightarrow méthode simple
- méthode du pivot total : permutation des lignes et des colonnes \Rightarrow méthode plus robuste

Dans le cas du pivot partiel sur les lignes :

Algorithm 4 Pivot partiel sur les lignes

entrée: A,B deux tableaux

- 1: **pour** i de 1 à $n - 1$ **faire**
 - 2: déterminer j tel que $a_{j,i} = \sup_{i \leq k \leq n} |a_{k,i}|$
 - 3: inverser ligne i et ligne j sur A et sur B
 - 4: $L_i, L_j \leftarrow L_j, L_i$
 - 5: **pour** k de $i + 1$ à n **faire**
 - 6: $L_k \leftarrow L_k - \frac{a_{k,i}}{a_{i,i}} \cdot L_i$
 - 7: **fin pour**
 - 8: **fin pour**renvoi: A,B
-

Remontée

Une fois la matrice triangularisée, si le problème initial est un système de Cramer, on obtient simplement la solution, ligne par ligne :

$$\forall i \in \llbracket 1, n \rrbracket \quad x_i = \frac{1}{a'_{i,i}} \cdot \left(b'_i - \sum_{k=i+1}^n a'_{i,k} \cdot x_k \right)$$

Remontée

Algorithm 5 Remontée

entrée: A', B'

- 1: **pour** i de n à 1 **faire**
 - 2: **pour** j de $i + 1$ à n **faire**
 - 3: $b'_i = b'_i - a'_{i,k} \cdot x_k$
 - 4: **fin pour**
 - 5: $x_i = \frac{b'_i}{a'_{i,i}}$
 - 6: **fin pour**
- renvoi: x
-

Inversion de matrice

Balayage

Lorsque la matrice est triangularisée, on utilise une méthode identique à celle de la triangularisation pour diagonaliser la matrice A .

Algorithm 7 Diagonalisation

entrée: A', B'

- 1: **pour** i de n à 2 **faire**
 - 2: **pour** k de $i - 1$ à 1 **faire**
 - 3: $L_k \leftarrow L_k - \frac{a'_{k,i}}{a'_{i,i}} \cdot L_i$
 - 4: **fin pour**
 - 5: **fin pour**
- renvoi: A', B'
-

Résolution

Le système d'équation ayant été diagonalisé, le problème de Cramer initial peut alors s'écrire :

$$\begin{cases} a''_{11} \cdot x_1 & = & b''_1 \\ a''_{22} \cdot x_2 & = & b''_2 \\ & \vdots & \\ a''_{nn} \cdot x_n & = & b''_n \end{cases}$$

Résolution

Le système d'équation ayant été diagonalisé, le problème de Cramer initial peut alors s'écrire :

$$\begin{cases} a''_{11} \cdot x_1 = b''_1 \\ a''_{22} \cdot x_2 = b''_2 \\ \vdots = \vdots \\ a''_{nn} \cdot x_n = b''_n \end{cases}$$

$$\Rightarrow \begin{bmatrix} a''_{11} & 0 & \cdots & 0 \\ 0 & a''_{22} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & a''_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b''_1 \\ b''_2 \\ \vdots \\ b''_n \end{bmatrix} \Rightarrow \forall x \in \llbracket 1, n \rrbracket x_i = \frac{b''_i}{a''_{i,i}}$$

Résolution

Algorithm 8 Résolution

entrée: A'', B''

1: **pour** i de 1 à n **faire**

2: $x_i = \frac{b_i''}{a_{i,i}''}$

3: **fin pour**

renvoi: x

Matrice inverse

On suppose la matrice A inversible, donnant une solution unique au problème de Cramer $A.X = B$. On appelle alors :

- $(T_i)_1^N$ les N matrices de transvection permettant de triangulariser la matrice A (passer de A à A') :

$$\underbrace{T_N \cdot T_{N-1} \dots T_1}_{T} \cdot A = A'$$

- $(S_i)_1^M$ les M matrices de transvection permettant de diagonaliser A' (passer de A' à A'')

$$\underbrace{S_M \cdot S_{M-1} \dots S_1}_{S} \cdot A' = A''$$

- D , la matrice diagonale ayant pour coefficient diagonaux $d_{i,j} = \frac{1}{a''_{i,j}}$

On obtient donc $D.S.T.A = I_n$. La matrice A étant inversible, $D.S.T$ apparaît donc comme l'**inverse de A** .

Ainsi pour **déterminer l'inverse d'une matrice inversible A** , il suffit d'**appliquer à la matrice identité les mêmes** opérations élémentaires sur les lignes et/ou les colonnes qui permettent de passer de A à I_n .

Algorithm 9 Inversion de matrice

entrée: A

- 1: $M \leftarrow A$ # copier la matrice
 - 2: déterminer la taille n, m de M
 - 3: vérifier que la matrice est carrée
 - 4: construire la matrice identité I_n
 - 5: **pour** i de 1 à n **faire**
 - 6: rechercher le pivot k sous le coefficient $M_{i,i}$
 - 7: **si** $M_{k,i} = 0$ **alors**
 - 8: **stop**, la matrice n'est pas inversible
 - 9: **fin si**
 - 10: **si** $k \neq i$ **alors**
 - 11: inverser ligne i et ligne k
 - 12: **fin si**
 - 13: obtenir des 0 sous le coefficient diagonal $M_{i,i}$
 - 14: **pour** k de $i + 1$ à n **faire**
 - 15: $L_k \rightarrow L_k - \frac{M_{k,i}}{M_{i,i}} \cdot L_i$
 - 16: **fin pour**
 - 17: **fin pour**
 - 18: obtenir $M_{i,i} = 1$
 - 19: **pour** i de 1 à n **faire**
 - 20: $L_i = L_i / M_{i,i}$
 - 21: **fin pour**
 - 22: éliminer tous les coefficients au dessus de $M_{i,i}$
 - 23: **pour** i de n à 2 **faire**
 - 24: **pour** k de 1 à $i - 1$ **faire**
 - 25: $L_k \rightarrow L_k - M_{k,i} \cdot L_i$
 - 26: **fin pour**
 - 27: **fin pour**
- renvoi: I_n

Sommaire

- 16 Introduction
- 17 Méthode de Gauss
- 18 Résoudre un système système de Cramer avec Numpy**

Résoudre un système système de Cramer avec Numpy

OBJECTIF : Résoudre un système linéaire du type $A.X = B$ où A est une matrice inversible.

Pour faire de l'algèbre linéaire avec Python, le plus simple peut être d'utiliser le module `numpy.linalg` contenant les méthodes `det(A)` pour le calcul du déterminant de A , `inv(A)` pour le calcul de la matrice inverse de A et `solve(A, B)` pour obtenir la solution du problème linéaire $A.X = B$ (où X est l'inconnue).

La classe `array` (tableaux numpy) possède beaucoup de méthodes notamment `T` pour obtenir la transposée d'un tableau et `dot(B)` pour multiplier un tableau à droite par B .

Résoudre un système de Cramer avec Numpy

```
A = np.array([[2, 1, -1], [6, -1, -2], [1, 2, 3]])
## Vérification du déterminant
print(nalg.det(A))
X = np.array([2, 3, 7]).T
## Produit matriciel de A par B
B = A.dot(X)
```

```
## [Résolution] par calcul
## de la matrice inverse
inv_A = nalg.inv(A)
X_sol_1 = inv_A.dot(B)
## Résolution avec solveur
X_sol_2 = nalg.solve(A, B)
```