



Méthodes numériques

Delphine BRESCH-PIETRI et Pauline Bernard

Stage Liesse
Mai 2022

Programme

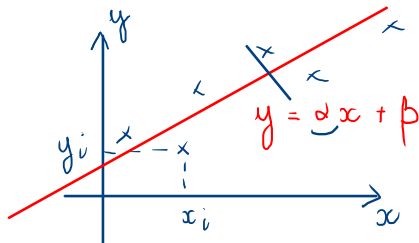
Résolution d'équations

Equations linéaires

Equations non linéaires

Résolution d'équations différentielles

- Régression linéaire
 $(x_i, y_i), 1 \leq i \leq n$



$$\alpha x_i + \beta = y_i \rightarrow \begin{cases} \alpha x_1 + \beta = y_1 \\ \alpha x_2 + \beta = y_2 \\ \vdots \end{cases}$$

- Recherche de points d'équilibre

$$\ddot{\theta} = -\frac{g}{l} \sin \theta$$

$$\begin{aligned} & \dot{x} = f(x) \\ & \hookrightarrow 0 = f(x) \\ & \text{point d'équilibre} \end{aligned}$$

$$n = 2: x_1 = \alpha, x_2 = \beta$$

$$\begin{cases} a_{j1} = x_j \\ a_{j2} = 1 \\ b_j = y_j \end{cases}$$

inconnues

$$\begin{cases} \underbrace{a_{1,1}x_1 + \dots + a_{1,n}x_n}_{= x_1} = b_1 & y_1 \\ \vdots \\ \underbrace{a_{p,1}x_1 + \dots + a_{p,n}x_n}_{= 1} = b_p & y_p \end{cases}$$

que l'on écrit sous forme matricielle

$$\underbrace{A}_{\in \mathcal{M}_{p \times n}(\mathbb{R})} \underbrace{x}_{\in \mathbb{R}^n} = \underbrace{b}_{\in \mathbb{R}^p} = \begin{pmatrix} b_1 \\ \vdots \\ b_p \end{pmatrix} \in \mathbb{R}^p$$

$$p \gg n$$

$$f: \mathbb{R}^m \rightarrow \mathbb{R}, \quad \nabla f: \mathbb{R}^m \rightarrow \mathbb{R}^m$$

Cas $p > n$ (sur-déterminé) : on cherche à minimiser l'écart entre les mesures et la régression linéaire

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|^2 = \sum_{i=1}^p \left(\sum_{j=1}^n a_{ij} x_j - b_j \right)^2$$

où $\|\cdot\|$ est la norme euclidienne dans \mathbb{R}^p . On peut montrer que la solution de ce problème est unique et telle que

$$\nabla f(x^*) = 0 = 2A^T Ax^* - 2A^T b \Leftrightarrow \boxed{x^* = (A^T A)^{-1} A^T b}$$

$$\begin{aligned}
 f(x) &= \|Ax - b\|^2 = (Ax - b)^T (Ax - b) \\
 &= \underline{x^T A^T A x} - \underline{2x^T A^T b} + b^T b
 \end{aligned}$$

A de rang plein (de rang n) $\Rightarrow A^T A$ inversible
 (on peut toujours le faire)

$n \times n$, A de rang $n \neq 0 \Rightarrow A$ inversible
(et A^T aussi)

$$x^* = (A^T A)^{-1} A^T b$$

$$= A^{-1} \cancel{(A^T)^{-1}} \cancel{A^T} b$$

$$= A^{-1} b$$

$$\min \|x\|_1 = \sum_{i=1}^n |x_i|$$

Cas $n > p$ (sous-déterminé) :

$$\min_{Ax=b} \|x\|^2$$

Il existe une unique solution qui est $x^* = A^T(AA^T)^{-1}b$.

`lstsq(A, b, rcond)[0]`

de la librairie `numpy.linalg`, où

- ▶ `A` : array (p, n) des coefficients $\|$
- ▶ `b` : array ($p, 1$)
- ▶ `rcond` : optionnel, valeur de conditionnement de la matrice A . Mettre "rcond = None" pour éviter un warning.

La fonction renvoie :

- ▶ `x` : array à une dimension n contenant la solution $\rightarrow x^*$
- ▶ residuals : scalaire, somme des résidus au carré (fonction coût) $\rightarrow f(x^*)$
- ▶ `rank` : rang de la matrice A
- ▶ `s` : valeurs singulières de la matrice A

$$(\sum a_{ij} x_j - b_i)^2$$

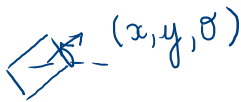
$$f(x) \in \mathbb{R}^m = \begin{pmatrix} f_1(x) \\ f_2(x) \\ f_3(x) \\ \vdots \\ f_p(x) \end{pmatrix} = 0 \quad f: \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$(x_1, x_2) \mapsto \begin{pmatrix} x_1 x_2 - 3 \\ \cos x_1 \\ \sin x_2 \end{pmatrix}$$

Dans le cas non-linéaire, il est souvent impossible de déterminer de façon explicite une solution à l'équation

$$f(x) = 0$$

où $x \in \mathbb{R}^n$ et $f: \mathbb{R}^n \rightarrow \mathbb{R}^p$. On dispose néanmoins de méthodes numériques permettant de réaliser cela.



$$\begin{cases} \dot{x} = v \cos \theta \\ \dot{y} = v \sin \theta \\ \dot{\theta} = \omega \end{cases} \rightarrow \begin{cases} v \cos \theta = 0 \\ v \sin \theta = 0 \\ \omega = 0 \end{cases}$$

f continue et tq $f(a_0)f(b_0) \leq 0$

Cas $n = p = 1$:



jusqu'à ce que $l([a_k, b_k]) < \underline{Tol}$

Algorithme

A partir de $[a_0, b_0]$, itérer

- ▶ $c_k = \frac{a_k + b_k}{2}$
- ▶ si $f(a_k)f(c_k) < 0$, choisir $\underline{a_{k+1}} = a_k$ et $\underline{b_{k+1}} = c_k$
- ▶ sinon, choisir $\underline{a_{k+1}} = c_k$ et $\underline{b_{k+1}} = b_k$

→ permet de trouver une solution x^*

à chaque itération, on améliore l'estimation "linéairement"

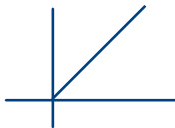
$$\|x_{k+1} - x^*\| \leq \frac{1}{2} \|x_k - x^*\| \leq \frac{1}{2^k} l([a_0, b_0])$$

\Rightarrow convergence linéaire de rapport $\frac{1}{2}$

$$n = p :$$

$$f(w) = 0 \Leftrightarrow x = g(w)$$

$$\Leftrightarrow x - f(w) = x$$



Si la fonction $g = Id - f$ est contractante, cad s'il existe une constante $\kappa \in [0, 1[$ telle que

$$\forall (x_1, x_2) \in \mathbb{R}^n \times \mathbb{R}^n \quad \|g(x_1) - g(x_2)\| \leq \underline{\underline{\kappa}} \|x_1 - x_2\|$$

on peut alors résoudre l'équation en itérant la relation

$$\| \underline{x_{k+1} = g(x_k)} = x_k - f(x_k) \|$$

Théorème du point fixe de Banach :

si g est contractante, $\exists ! x^* \text{ tq } x^* = g(x^*)$
 et la suite $\begin{cases} x_{n+1} = g(x_n) \\ x_0 \in \mathbb{R}^n \end{cases}$ est convergente de limite x^*

$$\dot{T} = -\underline{a} T + \underline{b} \underline{u}(t)$$

$$\rightarrow T(t) = e^{-at} T_0 + \int_0^t e^{-a(t-s)} b u(s) ds$$

or seek however k_1 s.t. $T(k_1) = T_1$
with s.t. $f(k_1) = T(k_1) - T_1 = 0$ ||

$$\Leftrightarrow k_1 = -\frac{1}{a} \ln \left(\frac{T_1 - \int_0^{k_1} e^{-a(k_1-s)} b u(s) ds}{T_0} \right)$$

$$\boxed{k_1 \stackrel{\Delta}{=} g(k_1)}$$

Cas $n = p$:

$$f: \mathbb{R}^m \rightarrow \mathbb{R}^m$$

$$\begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_m} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_m} \end{pmatrix}$$

Si f est différentiable, Taylor à l'ordre 1

$$f(x^*) = 0 = f(x) + J_f(x^*)(x^* - x) + o(\|x - x^*\|)$$

où $J_f(x^*)$ est la matrice jacobienne de f au point x^* .

$$\left[x^* = x - J_f(x)^{-1} f(x) \right]$$

La méthode de Newton-Raphson consiste à utiliser cette approximation et à choisir

$$x_{k+1} = x_k - J_f(x_k)^{-1} f(x_k)$$

$\Rightarrow \|x_{k+1} - x^*\| \leq M \|x_k - x^*\|^2 \leq (M \|x_0 - x^*\|)^{2^k}$
convergence quadratique

on voudrait que cette méthode converge vers une (ou la) solution x^* du problème

$f \in C^2$ et $f: \mathbb{R} \rightarrow \mathbb{R}$

$$x_{k+1} = x_k - \frac{1}{f'(x_k)} f(x_k) \quad (**)$$

et Taylor-Young avec reste intégral donne

$$\forall x \in \mathbb{R} \quad f(x) = f(x_k) + f'(x_k)(x - x_k) + \int_{x_k}^x f''(t)(x-t) dt \quad (**)$$

$$\text{donc } x_{k+1} = x_k - \frac{1}{f'(x_k)} \left[\underbrace{f(x^*)}_{=0} - f'(x_k)(x^* - x_k) - \int_{x_k}^{x^*} f''(t)(x^* - t) dt \right]$$

par (*) dans (**) pour $x = x^*$

$$x_{k+1} = x_k - \frac{1}{f'(x_k)} \left[\overbrace{f(x^*)}^{x_k} - f'(x_k)(x^* - x_k) - \int_{x_k}^{x^*} f''(t)(x^* - t) dt \right]$$

donc

$$x_{k+1} - x^* = \cancel{x_k - x^*} - \frac{1}{f'(x_k)} \left[\cancel{-f'(x_k)(x^* - x_k)} \right] + \frac{1}{f'(x_k)} \int_{x_k}^{x^*} f''(t)(x^* - t) dt$$

$\left\{ \begin{array}{l} f' \\ f'' \end{array} \right\}$ bornées au $v(x^*)$
" " " " " "

$$\begin{aligned} \Rightarrow \|x_{k+1} - x^*\| &\leq M \left| \int_{x_k}^{x^*} (x^* - t) dt \right| \\ &\leq M_0 \|x_k - x^*\|^2 \end{aligned}$$

Cas $p \geq n$:

On s'intéresse au problème

$$\begin{aligned} f(x) = 0 &\Leftrightarrow \|f(x)\|^2 = 0 \\ &\Leftrightarrow \min_x \|f(x)\|^2 \end{aligned}$$

$$\min_{x \in \mathbb{R}^n} \|f(x)\|^2 = F(x) \quad (1)$$

$$\min_{x \in \mathbb{R}^n} F$$

$$\begin{aligned} F(x) &= \|f(x)\|^2 \text{ avec } F: \mathbb{R}^n \rightarrow \mathbb{R} \\ &= \sum_{i=1}^p f_i(x)^2 \end{aligned}$$

Si F est différentiable, on a

$$F(x) = F(x_k) + \nabla F(x_k)^T (x - x_k) + o(\|x - x_k\|) \quad (2)$$

$$\simeq F(x_k) + \nabla F(x_k)^T (x - x_k) = F(x_k + \lambda)$$

Méthode de gradient à pas fixe

$$x_{k+1} = x_k - \lambda \nabla F(x_k)$$

on prend $x_{k+1} = x_k - \lambda \nabla F(x_k)$

$$-\lambda \|\nabla F(x_k)\|^2 \quad (3) \leq 0$$

≥ 0

$\overbrace{\text{bisect}(f, a, b, \text{args})}^{\text{obligatoire}}$
 $\overbrace{\text{args}}^{\text{optionnel}}$
→ dichotomique

de la librairie scipy.optimize, où

- ▶ ~~f~~ : fonction $\mathbb{R} \rightarrow \mathbb{R}$,
- ▶ a : scalaire de début d'intervalle
- ▶ b : scalaire de fin d'intervalle
- ▶ options : nombre maximal d'itérations, tolérance, ...

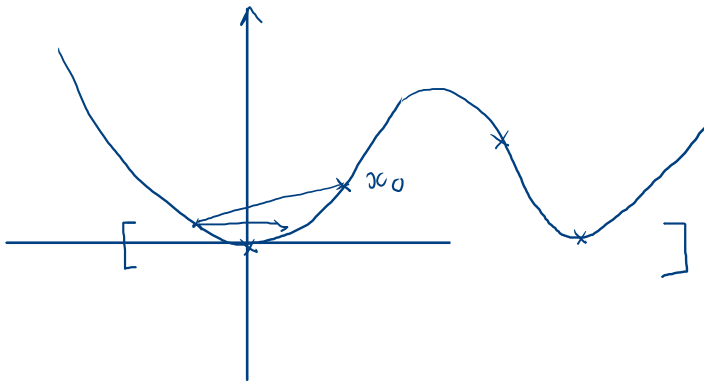
$\overbrace{\text{fsolve}(f, x_0, \text{args})}^{\text{obligatoire}}$
 $\overbrace{\text{args}}^{\text{optionnel}}$

→ Newton - Raphson /
descente de gradient
(Gauss-Newton)

de la librairie scipy.optimize, où

- ▶ ~~f~~ : fonction $\mathbb{R}^n \rightarrow \mathbb{R}^n$,
- ▶ x0 : condition initiale de taille n |
- ▶ options : jacobien, nombre maximal d'itérations, tolérance, ...

scipy.optimize.minimize ←



Programme

Résolution d'équations

Equations linéaires

Equations non linéaires

$$f(w) = 0$$

\nwarrow
 $w \in \mathbb{R}^n$

Résolution d'équations différentielles

$$\dot{x} = f(t, x)$$

\searrow
 $w(t) \in \mathbb{R}^n$

Exemples

- concentrations des espèces chimiques dans une réaction $A + B \xrightleftharpoons[k']{k} C$

$$\frac{dc_A}{dt} = \dot{c}_A = -k c_A c_B + k' c_C$$

$$\dot{c}_B = -k c_A c_B + k' c_C$$

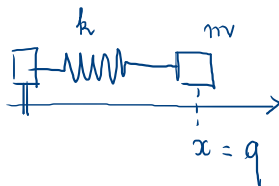
$$\dot{c}_C = k c_A c_B - k' c_C$$

→ 3 EDO d'ordre 1

- mécanique Newtonienne ou Lagrangienne

EDO ordre 2

$$\underbrace{M}_{\text{masse}} \ddot{q} = \sum_k F_k(t, q, \dot{q})$$



y de classe C^p vérifie

$$\| y^{(p)}(t) = \psi(t) y(t), \dot{y}(t), \dots, y^{(p-1)}(t) \|$$

si et seulement si $x = (y, \dot{y}, \dots, y^{(p-1)})$ de classe C^1 vérifie

$$\rightarrow \| \dot{x} = f(t, x),$$

où f est définie par

$$\dot{x} = \begin{pmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \vdots \\ \dot{x}_p \end{pmatrix} = \begin{pmatrix} \psi(t) y(t) \\ \dot{y}(t) \\ \vdots \\ y^{(p-1)}(t) \end{pmatrix}$$

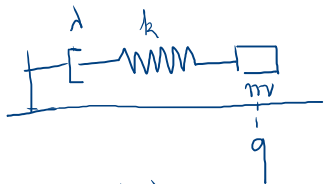
$$f(t, x_1, x_2, \dots, x_p) = \begin{pmatrix} x_2 \\ x_3 \\ \vdots \\ x_p \\ \psi(t, x_1, \dots, x_p) \end{pmatrix}$$

\Rightarrow toujours se mettre sous la forme $\dot{x} = f(t, x)$: crucial pour étude mathématique et numérique !

équation scalaire
d'ordre p



équation vectorielle (dans \mathbb{R}^p)
d'ordre 1



$$m\ddot{q} = -kq - \lambda\dot{q}$$

$$x = \begin{pmatrix} q \\ \dot{q} \end{pmatrix}$$

$$\dot{x} = \begin{pmatrix} \dot{q} = x_2 \\ \ddot{q} = -\frac{1}{m}(kq + \lambda\dot{q}) = -\frac{1}{m}(kx_1 + \lambda x_2) \end{pmatrix} = f(x)$$

$$f: x \in \mathbb{R}^2 \mapsto \begin{pmatrix} x_2 \\ -\frac{1}{m}(kx_1 + \lambda x_2) \end{pmatrix}$$

$$\begin{array}{l}
 \dot{\bullet} \rightarrow \\
 T(q_T)
 \end{array}
 \leftarrow \dot{\bullet} s(q_S)
 \left\{
 \begin{array}{l}
 m_T \ddot{q}_T = \dots \\
 m_S \ddot{q}_S = \dots
 \end{array}
 \right.$$

$$x = (q_T, \dot{q}_T, q_S, \dot{q}_S)$$

$$\dot{u} = \begin{pmatrix} \frac{\dot{q}_T}{m_T} () \\ q_S \\ \frac{1}{m_S} () \end{pmatrix}$$



$$\dot{x} = f(t, x) \quad , \quad x(t_0) = x_0$$

$$x(t) = x_0 + \int_{t_0}^t f(s, x(s)) ds = x_0 + \sum_{j=0}^{J-1} \int_{t_j}^{t_{j+1}} f(s, x(s)) ds$$

$$\underline{t_0} < \underline{t_1} < \dots < \underline{t_J} = t$$

Idée : approximer l'intégrale sur des intervalles $[t_j, t_{j+1}]$ suffisamment petits.

x^j : approximation de $x(t_j)$

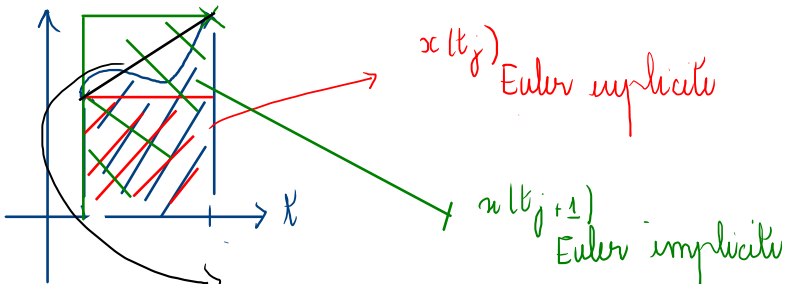
$$\Delta t_j = t_{j+1} - t_j$$

On implémente

$$\underline{x^{j+1}} = \underline{x^j} + \underbrace{\Delta t_j \Phi(t_j, x^j, \Delta t_j)}$$

tel que

$$\Phi(t_j, x^j, \Delta t_j) \approx \frac{1}{\Delta t_j} \int_{t_j}^{t_{j+1}} f(s, x(s)) ds$$



Trapezio $\frac{1}{2} (x(t_j) + x(t_{j+1}))$

$$\begin{cases} x_{j+1} = x_j + \Delta t_j f(t_j, x_j) \\ x_0 = x_0 \end{cases} \quad \begin{cases} \dot{x} = f(t, x) \\ x(0) = x_0 \end{cases}$$

Taylor ordre 1 :

$$\begin{aligned} x(t_{j+1}) &= x(t_j) + \Delta t_j \dot{x}(t_j) + o(\Delta t_j) \\ &= x(t_j) + \Delta t_j f(t_j, x(t_j)) + o(\Delta t_j) \end{aligned}$$

donc $\max_j \|x_j - x(t_j)\| \leq \epsilon \Delta t$ avec $\Delta t = \max_j \Delta t_j$
 \Rightarrow schéma d'ordre 1

Limites $\begin{cases} \rightarrow$ ordre 1 "||"
 \rightarrow systèmes "raides" (constantes de temps $\neq t_0$)
 \rightarrow systèmes conservatifs / hamiltoniens

$$\begin{cases} \underline{x_{j+1}} = x_j + \Delta t_j f(t_{j+1}, \underline{x_{j+1}}) \\ x_0 = x_0 \end{cases} \quad \begin{cases} \dot{x} = f(t, x) \\ x(t_0) = x_0 \end{cases}$$

on cherche à résoudre $x_{j+1} = F(x_{j+1})$ //

où $F: x \in \mathbb{R}^n \mapsto x_j + \Delta t_j f(t_{j+1}, x)$

si $f: (t, u) \mapsto f(t, u)$ est K -lipschitzienne par rapport à u :

$$\exists K \geq 0 \quad \forall (u_1, u_2) \in \mathbb{R}^n \times \mathbb{R}^n$$

$$\forall t \in \mathbb{R} \quad \|f(t, u_1) - f(t, u_2)\| \leq K \|u_1 - u_2\|$$

$$\begin{aligned} \|F(a) - F(b)\| &= \|x_j + \Delta t_j f(t_{j+1}, a) - [x_j + \Delta t_j f(t_{j+1}, b)]\| \\ &= \Delta t_j \|f(t_{j+1}, a) - f(t_{j+1}, b)\| \leq \Delta t_j K \|a - b\| \\ \Rightarrow \text{contractante si } \underline{\Delta t_j K} < 1 : & \underline{\text{on résout par point fixe}} \end{aligned}$$

$$\begin{cases} x_{j+1}^{k+1} = F(x_{j+1}^k) \\ x_{j+1}^0 = x_j + \Delta t_j f(t_j, x_j) \end{cases} \text{ Euler explicit}$$

$$x(t) = \exp(-\lambda t)$$

$$\leftarrow \dot{x} = -\lambda x \quad x(0) = 1 \text{ avec } \lambda > 0$$



Euler explicite : $x_{j+1} = x_j - \Delta t \lambda x_j$
 $= (1 - \Delta t \lambda) x_j$

|| si $\Delta t \lambda < 2$, alors $x_j \xrightarrow{j \rightarrow +\infty} 0$

Euler implicite : $x_{j+1} = x_j - \Delta t \lambda x_{j+1}$

donc $(1 + \lambda \Delta t) x_{j+1} = x_j$

stabilité inconditionnelle, $x_{j+1} = \frac{1}{1 + \lambda \Delta t} x_j \xrightarrow{j \rightarrow +\infty} 0$

\Rightarrow Explicite ou implicite pour $\dot{x} = \begin{pmatrix} -1 & 0 \\ 0 & -\mu \end{pmatrix} x$ avec $\mu \gg 1$?

Schéma convergent à l'ordre p si $\exists c_v > 0$ (indépendent de Δt) tel que

$$\max_{1 \leq j \leq J} \|x^j - x(t_j)\| \leq c_v (\Delta t)^p$$

Détermination numérique de l'ordre :

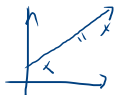
plusieurs (n_j) correspondant à plusieurs valeurs de Δt
 $\log(\text{ordre } \|x_j - x(t_j)\|) \propto \log(\Delta t)$

Détermination analytique de l'ordre :

$$x(t_{j+1}) = x(t_j) + \Delta t_j \dot{x}(t_j) + \frac{\Delta t_j^2}{2} \ddot{x}(t_j) + o(\Delta t_j^3)$$

$$= f(t_j, x(t_j))$$

$$= \partial_t f(t_j, x(t_j)) + \partial_x f(t_j, x(t_j)) \dot{x}(t_j) = f(t_j, x(t_j))$$



- méthode de Heun (ordre 2) :




$$x^{j+1} = x^j + \frac{\Delta t_j}{2} \left(\underline{f(t_j, x^j)} + \overbrace{f(t_{j+1}, x^j + \Delta t_j f(t_j, x^j))} \right),$$

ou sa version implicite méthode des trapèzes (ou Crank–Nicolson) :

$$x^{j+1} = x^j + \frac{\Delta t_j}{2} \left(f(t_j, x^j) + f(t_{j+1}, x^{j+1}) \right). \quad ||$$

- schéma de Runge–Kutta d'ordre 4 (RK4) largement utilisée:

$$x^{j+1} = x^j + \underline{\underline{\Delta t_j}} \frac{F_1 + 2F_2 + 2F_3 + F_4}{6} \quad \text{avec} \quad \left\{ \begin{array}{l} F_1 = f(t_j, x^j) \quad || \\ F_2 = f\left(t_j + \frac{\Delta t_j}{2}, x^j + \frac{\Delta t_j}{2} F_1\right) \quad | \\ F_3 = f\left(t_j + \frac{\Delta t_j}{2}, x^j + \frac{\Delta t_j}{2} F_2\right) \quad | \\ F_4 = f(t_j + \Delta t_j, x^j + \Delta t_j F_3), \quad | \end{array} \right.$$

- adaptation du pas 

→ solve - inq

→ odeint(fun, x0, t, args)

de la librairie scipy.integrate, où

- ▶ fun : fonction prenant en entrée (y, t) et renvoyant un array $f(y, t)$,
- ▶ y0 : la condition initiale ||
- ▶ t : suite des temps auxquels on souhaite résoudre y
- ▶ options : jacobien de la fonction, tolérance, points critiques, pas de temps min ou max, ...

La fonction renvoie :

- || ▶ y : array à deux dimensions de taille $(len(t), len(y_0))$


```
solve_ivp(fun,t_span,x0,method,t_eval)
```

de la librairie `scipy.integrate`, où

- ▶ `fun` : fonction prenant en entrée (t, x) et renvoyant un array $f(t, x)$,
- || ▶ `t_span` : intervalle d'intégration $[t_0, t_f]$
- ▶ `x0` : la condition initiale
- ▶ `method` : méthode d'intégration : par défaut, en mettant `[]` ou en nommant les entrées suivantes, c'est "RK45", le Runge-Kutta explicite d'ordre 5 (avec contrôle du pas à l'ordre 4). Lorsque cela échoue ou demande trop d'itérations, le système est peut-être raide : utiliser alors les méthodes implicites "Radau" ou "BDF".
- ||| ▶ `t_eval` : array donnant les temps auxquels vous voulez qu'il vous renvoie la solution (doit être dans l'intervalle `t_span`)
- ▶ `options` :

La fonction renvoie un objet `sol` de type `bunch` :

- || ▶ `sol.t` : array à une dimension contenant les temps
- ||| ▶ `sol.y` : array à deux dimensions (`dim_x, dim_t`), de façon à ce que `sol.y[k]` donne l'array du k -ième état au cours du temps.