

Lecture 5

Zed

March 26, 2017

1 Newton's Method for Solving Nonlinear Equations

We want to find roots for nonlinear equation $g(x) = 0$.

Algo. Newton's Method

- Initialize a starting point x_0 .
- Update x by $x_{n+1} \leftarrow x_n - \frac{g(x_n)}{g'(x_n)}$, suppose $g'(x_n) \neq 0$.
- Stop when $\|g(x_n)\| \leq \text{thres}$ (small), $\|x_{n+1} - x_n\| \leq \text{thres}$ (small) or $n \geq K$ (fail to converge).

The Newton's method only converges in a local sense. Suppose the real zero is α , s.t. $g(\alpha) = 0$, then using Taylor expansion at x_n :

$$\begin{aligned} 0 = g(\alpha) &= g(x_n) + g'(x_n)(\alpha - x_n) + \frac{1}{2}g''(\xi_n)(\alpha - x_n)^2 \\ \Rightarrow 0 &= \left(\frac{g(x_n)}{g'(x_n)} - x_n \right) + \alpha + \frac{g''(\xi_n)}{2g'(x_n)}(\alpha - x_n)^2 \end{aligned} \quad (1)$$

Hence exist bound M such that

$$|\alpha - x_{n+1}| = \frac{g''(\xi_n)}{2g'(x_n)}|\alpha - x_n|^2 \leq M|\alpha - x_n|^2$$

That is, if the error at n -th iteration $|\alpha - x_n|$ is already small, the next error at $(n+1)$ -th iteration will be the square of it. Which implies a (locally) quadratic convergence rate.

2 Implicit Methods

2.1 Implicit Runge-Kutta Method

We consider the Runge-Kutta method with table:

$$\begin{array}{c|cccccc} c_1 & a_{11} & a_{12} & \dots & a_{1,s-1} & a_{1s} \\ c_2 & a_{21} & a_{22} & \dots & a_{2,s-1} & a_{2s} \\ c_3 & a_{31} & a_{32} & \dots & a_{3,s-1} & a_{3s} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ c_s & a_{s1} & a_{s2} & \dots & a_{s,s-1} & a_{ss} \\ \hline & b_1 & b_2 & \dots & b_{s-1} & b_s \end{array} = \begin{array}{c|c} 1/2 & 1/2 \\ \hline & 1 \end{array}$$

A table of this kind with a full matrix of entries \mathbf{A} (instead of entries only in the lowertriangular area) suggests that this is an *implicit* RK method.

For the listed example, we have

$$\begin{cases} y_{n+1} = y_n + hb_1k_1 \\ k_1 = f(t_n + c_1h, y_n + ha_{11}k_1) \end{cases} \quad (2)$$

And insert the values of coefficients:

$$y_{n+1} = y_n + h \left(f(t_n + \frac{1}{2}h, \frac{y_{n+1}}{2} + \frac{y_n}{2}) \right)$$

Which is usually referred to as the *Implicit Midpoint Method*. This method is the simplest implicit RK method, the simplest *Gauss-Legendre* method. And a *Symplectic method*, which is energy-preserving.

The Midpoint method is of order 2 (see the first question in HW1).

Ex. (Energy Preserving) Consider the Hamilton system: $p' = -q, q' = p$; with $p(0) = q(0) = 1$. We can find a Hamilton function H such that $p' = -\frac{\partial H}{\partial q}, q' = \frac{\partial H}{\partial p}$. And hence $\frac{\partial H}{\partial t} = 0$. However, when we use the explicit RK method to solve the system (for example, the `ode45` in `Matlab`), we will find that the method does not preserve energy, i.e. $\frac{\partial H}{\partial t} \neq 0$.

3 Symplectic Methods

Our motivation to derive such a family of method is that it is natural to look into these discrete systems which preserves as much as possible the intrinsic properties of the continuous system. (Feng Kang, 1985; Ruth, 1983). We use the term *Symplectic* to refer to such systems, objects and methods.

Ex. We consider a simple system: $p' = -q, q' = p$, i.e. $H(p, q) = \frac{1}{2}p^2 + \frac{1}{2}q^2 = \text{const}$. And there is a numerical method to solve this ode system. At n -th step we are at (p_n, q_n) , and the numerical method will map this value to (p_{n+1}, q_{n+1}) in the next step.

We consider all the (p_n, q_n) 's within an infinitesimal area dV_n , and in next step it is mapped to dV_{n+1} . Furthermore, for all initial values (p_0, q_0) in the infinitesimal area dV_0 , after n steps of numerical iteration, it will go to dV_n . We want a kind of *stability* such that initial values that are close to each other will yield numerical solutions that are also close to each other. In the setting of this problem we want the area of $dV_n \approx dV_0$. Intuitively, this suggests the same thing in the sense of stability. We investigate some methods.

1. Explicit Euler

$$\begin{cases} p_m = p_{m-1} - hq_{m-1} \\ q_m = q_{m-1} + hp_{m-1} \end{cases} \quad (3)$$

The Jacobian

$$\mathbf{J}_{ex} = \frac{\partial(p_m, q_m)}{\partial(p_{m-1}, q_{m-1})} = \begin{vmatrix} \frac{\partial p_m}{\partial p_{m-1}} & \frac{\partial p_m}{\partial q_{m-1}} \\ \frac{\partial q_m}{\partial p_{m-1}} & \frac{\partial q_m}{\partial q_{m-1}} \end{vmatrix} = \begin{vmatrix} 1 & -h \\ h & 1 \end{vmatrix} = 1 + h^2$$

So we will see $dV_m = (1 + h^2)^m dV_0$. Any two initial values will evolve to be infinitely farwary.

2. Implicit Euler

$$\begin{cases} p_m = p_{m-1} - hq_m \\ q_m = q_{m-1} + hp_m \end{cases} \Rightarrow \begin{cases} (1 + h^2)p_m = p_{m-1} - hq_{m-1} \\ (1 + h^2)q_m = q_{m-1} + hp_{m-1} \end{cases} \quad (4)$$

The Jacobian

$$\mathbf{J}_{im} = \frac{\partial(p_m, q_m)}{\partial(p_{m-1}, q_{m-1})} = \begin{vmatrix} \frac{1}{1+h^2} & \frac{-h}{1+h^2} \\ \frac{h}{1+h^2} & \frac{1}{1+h^2} \end{vmatrix} = \frac{1}{1+h^2}$$

Therefore we will see $dV_m = (1+h^2)^{-m}dV_0$, the solution will eventually degenerate to one point.

3. Symplectic Euler

$$\begin{cases} p_m = p_{m-1} - hq_{m-1} \\ q_m = q_{m-1} + hp_m \end{cases} \Rightarrow \begin{cases} p_m = p_{m-1} - hq_{m-1} \\ q_m = q_{m-1} + h(p_{m-1} - hq_{m-1}) \end{cases} \quad (5)$$

$$\mathbf{J}_{symp} = \frac{\partial(p_m, q_m)}{\partial(p_{m-1}, q_{m-1})} = \begin{vmatrix} 1 & -h \\ h & 1-h^2 \end{vmatrix} = 1$$

4. Implicit Midpoint

$$\begin{cases} p_m = p_{m-1} - \frac{h}{2}q_{m-1} - \frac{h}{2}q_m \\ q_m = q_{m-1} + \frac{h}{2}p_{m-1} + \frac{h}{2}p_m \end{cases} \Rightarrow \begin{cases} (1 + \frac{h^2}{4})p_m = (1 - \frac{h^2}{4})p_{m-1} - hq_{m-1} \\ (1 + \frac{h^2}{4})q_m = (1 - \frac{h^2}{4})q_{m-1} + hp_{m-1} \end{cases} \quad (6)$$

$$\mathbf{J}_{mid} = \frac{\partial(p_m, q_m)}{\partial(p_{m-1}, q_{m-1})} = \begin{vmatrix} \frac{4-h^2}{4+h^2} & \frac{-4h}{4+h^2} \\ \frac{4h}{4+h^2} & \frac{4-h^2}{4+h^2} \end{vmatrix} = \frac{(4-h^2)^2 + 16h^2}{(4+h^2)^2} = 1$$

Def. Hamiltonian System We investigate Hamiltonian system in higher space dimensions: generalized coordinate $\mathbf{q} = (q_1, \dots, q_d)^\top$ and generalized momentum: $\mathbf{p} = (p_1, \dots, p_d)^\top$,

$$\begin{cases} \dot{\mathbf{p}} = -\frac{\partial H}{\partial \mathbf{q}} \\ \dot{\mathbf{q}} = \frac{\partial H}{\partial \mathbf{p}} \end{cases} \quad (7)$$

It is easy to check that the *Hamiltonian Energy* $H = \text{const}$ over time.

$$\frac{d}{dt}H(\mathbf{p}, \mathbf{q}) = \left(\frac{\partial H}{\partial \mathbf{p}}\right)^\top \dot{\mathbf{p}} + \left(\frac{\partial H}{\partial \mathbf{q}}\right)^\top \dot{\mathbf{q}} = 0$$

We can write in another form, let $\mathbf{y} = (\mathbf{p}, \mathbf{q})^\top = (p_1, \dots, p_d; q_1, \dots, q_d)^\top$, $\mathbf{J} := \begin{pmatrix} \mathbf{O} & \mathbf{I}_d \\ -\mathbf{I}_d & \mathbf{O} \end{pmatrix}$, one can easily check that $\mathbf{J}^{-1} = -\mathbf{J}$. Then we can write the system as $\dot{\mathbf{y}} = \mathbf{J}^{-1}\nabla H(\mathbf{y})$ (*).

Proof. $\mathbf{J}^{-1} = -\mathbf{J}$ because:

$$\mathbf{J}\mathbf{J} = \begin{pmatrix} \mathbf{O} & \mathbf{I}_d \\ -\mathbf{I}_d & \mathbf{O} \end{pmatrix} \begin{pmatrix} \mathbf{O} & \mathbf{I}_d \\ -\mathbf{I}_d & \mathbf{O} \end{pmatrix} = \begin{pmatrix} -\mathbf{I}_d & \mathbf{O} \\ \mathbf{O} & -\mathbf{I}_d \end{pmatrix} = -\mathbf{I}_{2d}$$

And

$$RHS(*) = \begin{pmatrix} \mathbf{O} & -\mathbf{I}_d \\ \mathbf{I}_d & \mathbf{O} \end{pmatrix} \begin{pmatrix} \nabla_{\mathbf{p}} H \\ \nabla_{\mathbf{q}} H \end{pmatrix} = \begin{pmatrix} -\nabla_{\mathbf{q}} H \\ \nabla_{\mathbf{p}} H \end{pmatrix} = \begin{pmatrix} -\frac{\partial H}{\partial \mathbf{q}} \\ \frac{\partial H}{\partial \mathbf{p}} \end{pmatrix} = (\dot{\mathbf{p}}, \dot{\mathbf{q}})^\top = LHS$$

Def. Flow Map: Given an autonomous differential equation $\mathbf{y}' = \mathbf{f}(\mathbf{y})$, the *flow map* $\varphi_t(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$, such that $\varphi_t(\mathbf{y}_0) = \mathbf{y}(t)$ maps the initial value to the solution at t . The definition can be extended to the flow map of *measurable sets* $\Omega \subseteq \mathbb{R}^d$:

$$\begin{aligned} \varphi_t : \mathcal{B}(\mathbb{R}^d) &\rightarrow \mathcal{B}(\mathbb{R}^d) \\ \varphi_t(\Omega) &= \{\mathbf{y}(t) : \mathbf{y}_0 \in \Omega\} \end{aligned} \quad (8)$$

where $\mathcal{B}(\mathbb{R}^d)$ is the Borel sigma-algebra on \mathbb{R}^d (the collection of all measurable sets).

Def. Symplecticity: For the Hamiltonian system $(\mathbf{p}, \mathbf{q}) \in \mathbb{R}^d \times \mathbb{R}^d$:

1. If p, q has only one dimension, i.e. $d = 1$, the system is *symplectic* if the area is preserved in phase space $\mathbb{R} \times \mathbb{R}$ under the flow map φ_t .
2. $d \geq 2$, $(\mathbf{p}, \mathbf{q}) = (p_1, \dots, p_d; q_1, \dots, q_d)$ the system is symplectic if the quantity

$$w^2 = \sum_{i=1}^d dp_i \wedge dq_i$$

is preserved under the flow map φ_t . Where \wedge is the wedge product.

Prop. The equivalent definition of symplecticity: a dynamic system is symplectic \iff

$$\mathbf{J} = \Phi_t^\top \mathbf{J} \Phi_t$$

where $\mathbf{J} = \begin{pmatrix} \mathbf{O} & \mathbf{I}_d \\ -\mathbf{I}_d & \mathbf{O} \end{pmatrix}$, $\Phi_t = \Phi_t(\mathbf{y}) = \partial \varphi_t(\mathbf{y}) / \partial \mathbf{y}$ is the Jacobian of the flow map.

*Proof.**

$$\begin{aligned} w^2 &= \sum_{i=1}^d dp_i \wedge dq_i = \frac{1}{2} \left(\sum_{i=1}^d dp_i \wedge dq_i - \sum_{i=1}^d dq_i \wedge dp_i \right) \\ &= \frac{1}{2} \mathbf{J}^{-1} d\varphi_t \wedge d\varphi_t = \frac{1}{2} \mathbf{J}^{-1} \Phi_t d\mathbf{y} \wedge \Phi_t d\mathbf{y} \\ &= \frac{1}{2} \Phi_t^\top \mathbf{J}^{-1} \Phi_t d\mathbf{y} \wedge d\mathbf{y} \end{aligned} \tag{9}$$

And we require $w^2(t) = w^2(0)$, knowing $\Phi_t(0) = \mathbf{I} \Rightarrow \frac{1}{2} \Phi_t^\top \mathbf{J}^{-1} \Phi_t d\mathbf{y} \wedge d\mathbf{y} = \frac{1}{2} \mathbf{J}^{-1} d\mathbf{y} \wedge d\mathbf{y}$. So we have $\Phi_t^\top \mathbf{J} \Phi_t = \mathbf{J}$. \square

Thm. (Poincaré): Hamiltonian system \iff Symplectic system.

Proof. We only show (\Rightarrow) , i.e. showing $\Phi_t^\top \mathbf{J} \Phi_t = \mathbf{J}$ from Hamiltonian.

We use $\dot{\mathbf{y}} = \mathbf{J}^{-1} \nabla H(\mathbf{y})$, $\varphi_t(\mathbf{y})$ is flow map. Since $\varphi_t(\mathbf{y}) = \mathbf{y}(t)$, it must satisfy the ode when we regard t as its variable (in the time dimension). That is

$$\frac{d\varphi_t(\mathbf{y})}{dt} = \mathbf{J}^{-1} \nabla H(\varphi_t(\mathbf{y}))$$

$$\frac{d\Phi_t(\mathbf{y})}{dt} = \frac{\partial}{\partial \mathbf{y}} \frac{d\varphi_t(\mathbf{y})}{dt} = \mathbf{J}^{-1} \nabla^2 H(\varphi_t(\mathbf{y})) \frac{\partial \varphi_t(\mathbf{y})}{\partial \mathbf{y}} = \mathbf{J}^{-1} \nabla^2 H(\varphi_t(\mathbf{y})) \Phi_t(\mathbf{y})$$

where $\nabla^2 H$ is the hessian. Therefore

$$\begin{aligned} \frac{d}{dt} (\Phi_t^\top \mathbf{J} \Phi_t) &= \left(\frac{d\Phi_t}{dt} \right)^\top \mathbf{J} \Phi_t + \Phi_t^\top \mathbf{J} \frac{d\Phi_t}{dt} \\ &= (\mathbf{J}^{-1} \nabla^2 H(\varphi_t(\mathbf{y})) \Phi_t)^\top \mathbf{J} \Phi_t + \Phi_t^\top \mathbf{J} \mathbf{J}^{-1} \nabla^2 H(\varphi_t(\mathbf{y})) \Phi_t \\ &= \Phi_t^\top \nabla^2 H(\varphi_t(\mathbf{y})) \mathbf{J}^{-\top} \mathbf{J} \Phi_t + \Phi_t^\top \mathbf{J} \mathbf{J}^{-1} \nabla^2 H(\varphi_t(\mathbf{y})) \Phi_t \\ &= -\Phi_t^\top \nabla^2 H(\varphi_t(\mathbf{y})) \Phi_t + \Phi_t^\top \nabla^2 H(\varphi_t(\mathbf{y})) \Phi_t \\ &= 0 \end{aligned} \tag{10}$$

Using the fact that $\mathbf{J}^{-\top} = (-\mathbf{J})^\top = \begin{pmatrix} \mathbf{O} & -\mathbf{I}_d \\ \mathbf{I}_d & \mathbf{O} \end{pmatrix}^\top = \mathbf{J}$, and $\mathbf{J}\mathbf{J} = -\mathbf{I}$. Therefore $\Phi_t^\top \mathbf{J} \Phi_t = \Phi_t^\top \mathbf{J} \Phi_t|_{t=0}$. Note that $\Phi_t|_{t=0} = \det(\mathbf{I})$. Hence we have $\Phi_t^\top \mathbf{J} \Phi_t = \mathbf{J}$. \square