

DYNAMICAL SYSTEMS
TUTORIAL 1: NUMERICAL ANALYSIS OF ODE

ROY MALKA, 2008; MICHAL PNUELI, 2019

1. INTRODUCTION

Following the notation of [1].

Our aim is to obtain a numerical solution of the ordinary differential equation with initial condition (Cauchy problem):

$$(1) \quad \dot{y} = f(t, y), \quad y(t_0) = y_0$$

where $\dot{y} = \frac{dy}{dt}$.

How do computers calculate such a solution of (1) which is continuous in time? Exactly as numbers are all rational on the computer, the computer must make time discrete. So we replace differential equations with a difference equation.

2. THE EULER ONE-STEP METHOD

Dividing the time interval $[t_0, t_M]$ to $N + 1$ mesh points, with a constant distance (i.e. time-step, step size or mesh size): $h = (t_M - t_0)/N$.

Given $y(t_0) = y_0$ (this is the initial point so the numerical approximation is *equal* to the actual value), let us suppose that we have calculated y_n up to some n , $0 \leq n \leq N - 1$, $N \geq 1$, thus we define

$$(2) \quad y_{n+1} = y_n + hf(t_n, y_n)$$

Taking $n = 0, 1, \dots, N - 1$, iterating step by step, we obtain $\{y_n\}$ the approximation of $y(t)$ at the mesh points t_n .

What is the motivation for this method: expanding $y(t_{n+1}) = y(t_n + h)$ into a Taylor series about t_n we get:

$$y(t_n + h) = y(t_n) + hf(t_n, y(t_n)) + \mathcal{O}(h^2)$$

replacing $y(t_n)$ with its numerical approximation y_n we get Eq. (2). By which we get $\frac{y_{n+1} - y_n}{h} \approx f(t_n, y_n)$.

This is an example of a general one-step method:

$$(3) \quad y_{n+1} = y_n + h\Phi(t_n, y_n; h)$$

Date: April 1, 2019.

where Φ is a continuous function of its arguments.

The meaning of "one step" is that y_{n+1} is expressed in terms of y_n , where y_n is the numerical approximation of $y(t_n)$.

Similarly, a k -step method would be expressing y_{n+1} in terms of y_{n-k+1}, \dots, y_n . The simplest example of a one-step method is Euler's method.

In order to analyze the accuracy of numerical methods we need the aid of the following definitions.

Definition 1. *The global error is defined as*

$$(4) \quad e_n = y(t_n) - y_n$$

Measuring the distance between the analytic solution ($y(t_n)$) and the numerical approximation (y_n) at time t_n .

Definition 2. *The truncation error is defined as*

$$(5) \quad T_n = \frac{y(t_{n+1}) - y(t_n)}{h} - \Phi(t_n, y(t_n); h)$$

Notice that for $n = 0$: $T_0 = \frac{y(t_1) - y(t_0)}{h} - \Phi(t_0, y(t_0); h) = \frac{y(t_1) - y_0}{h} - \Phi(t_0, y_0; h)$ since $y(t_0) = y_0$, using the global error, we get $T_0 = \frac{y_1 \pm e_1 - y_0}{h} - \Phi(t_0, y_0; h)$ assign the numerical approximation by a one step method we get $T_0 = \frac{y_0 + h\Phi(t_0, y_0; h) \pm e_1 - y_0}{h} - \Phi(t_0, y_0; h) = \pm \frac{e_1}{h}$. This shows that the size of the truncation error is a measure of local error.

Now that we see that the definitions are connected we are ready for the following theorem that give an upper bound on the global error using the truncation error.

Theorem 1. *Given the one step method with update equation (3), where $\Phi(t, y; h)$ is continuous in all its arguments and satisfy Lipschitz condition with respect to the second argument, that is, there exists a positive constant L_Φ s.t. for $0 \leq h \leq h_0$ and for all (t, u) and (t, v) in $D = \{(t, y) : t_0 \leq t \leq t_M, |y - y_0| \leq C\}$, we have that*

$$|\Phi(t, u; h) - \Phi(t, v; h)| \leq L_\Phi |u - v|.$$

Then, assuming that $|y_n - y_0| \leq C$, for $n = 1, 2, \dots, N$ it follows that

$$(6) \quad |e_n| \leq \frac{T}{L_\Phi} (e^{L_\Phi(t_n - t_0)} - 1), \quad n \in \{1, \dots, N\},$$

where $T = \max_{0 \leq n \leq N-1} |T_n|$.

Proof. Rewriting (5)

$$y(t_{n+1}) = y(t_n) + h\Phi(t_n, y(t_n); h) + hT_n$$

and subtracting (3) from this we obtain

$$y(t_{n+1}) - y_{n+1} = y(t_n) + h\Phi(t_n, y(t_n); h) + hT_n - y_n - h\Phi(t_n, y_n; h)$$

$$e_{n+1} = e_n + h[\Phi(t_n, y(t_n); h) - \Phi(t_n, y_n; h)] + hT_n$$

Since $(t_n, y(t_n)), (t_n, y_n)$ are in D the Lipschitz condition implies that (first triangle inequality)

$$|e_{n+1}| \leq |e_n| + h|[\Phi(t_n, y(t_n); h) - \Phi(t_n, y_n; h)]| + h|T_n| \leq |e_n| + hL_\Phi|e_n| + h|T_n|$$

rearrange to get

$$|e_{n+1}| \leq (1 + hL_\Phi)|e_n| + h|T_n|$$

solving the recursion we get (reminder $e_0 = 0$)

$$|e_n| \leq \frac{T}{L_\Phi} |(1 + hL_\Phi)^n - 1|.$$

Observing that $(1 + hL_\Phi) \leq e^{(hL_\Phi)}$, and $Nh = (t_M - t_0)$ we are done. \square

Implication let us apply this general result to analyze Euler's method, in the case that $y \in \mathcal{C}^2[t_0, t_M]$ i.e., twice differentiable with respect to time. Expanding $y(t_{n+1})$ in a Taylor series with reminder:

$$y(t_{n+1}) = y(t_n) + h\dot{y}(t_n) + \frac{h^2}{2!}\ddot{y}(\xi_n), \quad t_n \leq \xi_n \leq t_{n+1}$$

Substituting this into the (5) where $\Phi(t, y; h) \equiv f(t, y)$, we get,

$$T_n = \frac{1}{2}h\ddot{y}(\xi_n).$$

Let $H_2 = \max_{\xi \in [t_0, t_M]} |\ddot{y}(\xi)|$. Then, $|T_n| \leq \frac{1}{2}hH_2, n = 0, 1, \dots, N-1$. Substituting into (6), noting that $L_\Phi = L$ the Lipschitz constant of f , we have

$$|e_n| \leq \frac{1}{2}H_2 \left[\frac{e^{L(t_M - t_0)} - 1}{L} \right] h, \quad n = 0, 1, \dots, N.$$

In order to demonstrate the relevance of the above error analysis in practice, we turn to the following example.

Example 1. *Let*

$$\dot{y} = \tan^{-1}(y), \quad y(0) = y_0.$$

In order to obtain an upper bound on the global error in Euler's approximation we need to determine the constants H_2 and L .

Here $f(t, y) = \tan^{-1}(y)$; so by Mean Value theorem,

$$|f(t, u) - f(t, v)| = \left| \frac{\partial f}{\partial y}(t, \eta)(u - v) \right| = \left| \frac{\partial f}{\partial y}(t, \eta) \right| |u - v|$$

where η lies between u and v . In our case,

$$\left| \frac{\partial f}{\partial y}(t, y) \right| = |(1 + y^2)^{-1}| \leq 1$$

Thus, $L = 1$. In order to obtain H_2 we need to bound $|\ddot{y}(t)|$ with out solving the initial value! We do that by differentiation of the ODE:

$$\ddot{y} = \frac{d}{dt}(\tan^{-1}y) = (1 + y^2)^{-1}\dot{y} = (1 + y^2)^{-1}\tan^{-1}y.$$

Therefore, $|\ddot{y}(t)| \leq H_2 = \frac{1}{2}\pi$. Assigning L and H_2 to (6) we get

$$|e_n| \leq \frac{1}{4}\pi(e^{t_n})h, \quad n = 0, 1, \dots, N.$$

Now, if we are givgin a tolerance TOL , we can ensure that the global error does not exceed this tolerance by choosing h accordingly

$$h \leq \frac{4}{\pi(e^{t_M} - 1)}TOL.$$

So for each n , $|y(t_n) - y_n| \leq |e_n| \leq TOL$.

3. CONSISTENCY AND CONVERGENCE

There are three main issues that should be considered in the analysis of a numerical method:

- (1) **Stability** (of method): the difference equation is stable if its solution is continuous and Lipschitz in the initial conditions (how close is $\{z_n\}$ to $\{y_n\}$ if $z_0 = y_0 + \varepsilon$).
- (2) **Consistency** of difference equation with an ode: connection between step size and local error.
- (3) **Convergence** (of sol. an ode sol.): when h goes to zero the solution of the difference equation goes to the solution of the ode

The error analysis in theorem 1 suggests that for a one-step method, if the truncation error 'approaches zero' as $h \rightarrow 0$ then the global error 'approaches zero' as well. This is the motivation for the following definition.

Definition 3 (Consistency). *The numerical method (3) is consistent with the initial value ODE problem (1) if the truncation error (5) is such that for any $\varepsilon > 0$ there exists a positive $h(\varepsilon)$ for which $|T_n| < \varepsilon$, for $0 < h < h(\varepsilon)$ and any pair of points $(t_n, y(t_n)), (t_{n+1}, y(t_{n+1}))$ on any solution curve in D .*

In the limit of: $h \rightarrow 0$ and $\lim_{n \rightarrow \infty} t_n = t \in [t_0, t_M]$, we have $T_n \rightarrow \frac{dy}{dt} - \Phi(t, y(t); h = 0)$, this implies that one step method is consistent if and only if

$$(7) \quad \Phi(t, y; h)|_{h=0} \equiv f(t, y).$$

One can show that consistency is a necessary condition for convergence.

Definition 4 (Order of accuracy). *Numerical method (3) is said to have order of accuracy p , if p is the largest positive integer such that, for any sufficiently smooth solution curve $(t, y(t))$ in D of the initial value problem (1), there exists constants K and h_0 such that*

$$|T_n| \leq Kh^p \text{ for } 0 < h \leq h_0,$$

for any points $(t_n, y(t_n)), (t_{n+1}, y(t_{n+1}))$ on the solution curve.

4. IMPLICIT ONE STEP METHOD

Euler Implicit method

$$(8) \quad y_{n+1} = y_n + hf(t_{n+1}, y_{n+1})$$

The main difference between Euler's method and (8) is in the appearance of y_{n+1} on both sides of (8). This complication requires additional computational steps.

5. RUNGE-KUTTA METHODS

The Euler method give first order accuracy, but it is very simple and it requires a single evaluation of f at (t_n, y_n) in order to get the approximation of y_{n+1} . The Runge-Kutta methods are coming to improve accuracy by increasing the number of evaluations of f at intermidate points between (t_n, y_n) and (t_{n+1}, y_{n+1}) .

6. STIFF SYSTEMS

A stiff equation is a differential equation for which certain numerical methods for solving the equation are numerically unstable, unless the step size is taken to be extremely small. There is no unique definition of stiffness in the literature. However, essential properties of stiff systems are as follows:

- There exist, for certain initial conditions, solutions that change slowly.
- Solutions in a neighborhood of these smooth solutions converge quickly to them.

Usually, stiffness appears in a system of differential equations, but let's start with a simple 1-dimensional example:

$$\dot{y} = \lambda y, \quad y(0) = y_0$$

where λ is a constant. The solution is $y(t) = y_0 e^{\lambda t}$. When $\lambda < 0$ the solution decreases exponentially to zero. We expect the numerical approximation will behave accordingly. Expressing the numerical approximation in Euler's method:

$$y_n^E = (1 + h\lambda)^n y_0$$

When $\lambda < 0$, $|1 + h\lambda| < 1$ if and only if $0 < h|\lambda| < 2$. This gives the problem-dependent restriction on the size of h : $h < \frac{2}{|\lambda|}$. For higher values of h the solution of Euler's method will oscillate with increasing magnitude with increasing n , instead of converging to zero as $n \rightarrow \infty$.

While the above example indicates the problem, in a single equation it is not too hard, but in a system of equations it becomes more intricate.

REFERENCES

- [1] E. Suli and D. Mayers. An Introduction to Numerical Analysis. Cambridge University Press, 2003.
- [2] https://www.wias-berlin.de/people/john/LEHRE/NUMERIK_II/ode_2.pdf