

# Euler's Method

## Math490: Week1

Ayberk Nardan

February 2023

### Abstract

The notes are about the first chapter *The Euler's Method* of the **Iserles's** book. It starts with introducing the Lipschitz Condition and then passes into the Euler's method. While giving the introduction about the method, the chapter also tells what the convergence of a method means and why it is necessary to prove it.

Also, the chapter introduces error analysis. It tries to give information about how error behaves in Euler's method and other methods given in the chapter.

### Contents

<b>1</b>	<b>Lipschitz Condition</b>	<b>2</b>
<b>2</b>	<b>Euler's Method</b>	<b>2</b>
2.1	Convergence . . . . .	3
<b>3</b>	<b>Trapezoidal Rule</b>	<b>4</b>
3.1	Implicit Midpoint Method . . . . .	5
<b>4</b>	<b>The Theta Method</b>	<b>6</b>
4.1	Order of the Theta Method . . . . .	6
4.2	Error Analysis of the Theta Method . . . . .	6

## 1 Lipschitz Condition

We want to approximate solutions to problems in the following form when we mean solving an ODE numerically.

$$y' = f(t, y) \quad , t \geq t_0 \quad y(t_0) = y_0 \quad (1)$$

Here, we need  $f(t, y)$  be a sufficiently well-behaved function maps the space  $[t_0, \inf](R)^D \rightarrow (R)^D$ . Also, by the mean of well-behaved, we want to ensure that at least, the  $f$  obeys the given vector norm and the Lipschitz Condition, which is

$$\exists \lambda > 0 \text{ s.t. } \|f(t, y) - f(t, x)\| < \lambda \|x - y\| \quad \forall x, y \in \mathbb{R}^D \quad t \geq t_0 \quad (2)$$

**Note:** As a stronger condition, we can also assume that the  $f$  is analytic. In that case, there is a theory that says the solution is also analytic.

## 2 Euler's Method

It is a basic method to solve ODEs in the following form.

$$y' = f(t, y) \quad \text{where } y(t_0) = y_0 \quad t > t_0 \quad (3)$$

To estimate the solution, We write the  $y(t)$  in the form,

$$y(t) = y_0 + \int_{t_0}^t f(t, y(t)) dt$$

Then do the approximation

$$f(t, y(t)) \approx f(t_0, y_0) \quad t \in [t_0, t_0 + \epsilon]$$

Consequently, our estimation becomes

$$y(t) = y_0 + (t - t_0)f(t_0, y_0) \quad (4)$$

This estimation is the basis of Euler's method. Then we create a sequence of time as  $\{t_0, t_0 + h = t_1, t_1 + h = t_2, \dots\}$ , and use the time sequence to calculate a sequence of  $y$  values  $\{y_n\}$  which is an estimation of the exact solution.

$$\begin{aligned} y_1 &= y_0 + hf(t_0, y_0) \\ y_2 &= y_1 + hf(t_1, y_1) \\ &\vdots \\ y_{n+1} &= y_n + hf(t_n, y_n) \end{aligned}$$

Consequently, the generalized form of Euler's method is

$$y_{n+1} = y_n + hf(t_n, y_n) \quad (5)$$

## 2.1 Convergence

Euler's method is quite simple and easy to use. However, to ensure that our estimation of the solution is reliable, we should check if the method converges or does not. This step must be done for all the numerical methods actually.

Fortunately, Euler's method always converges. Before proving it, we need to define the convergence of a method first.

**Definition 2.1.** *A method is said to be convergent if, for every ODE in the form of (1) with a lipshcitz function  $f$  and for every  $t^* > 0$ , it is true that*

$$\lim_{h \rightarrow 0} \max_{n=0,1,\dots,\lfloor t^*/h \rfloor} \|y_{n,h} - y(t_n)\| = 0$$

**Theorem 2.1.** *Euler's Method is convergent.*

*Proof.* Firstly, for the given  $h > 0$ , lets define  $e_{n,h} = y_{n,h} - y(t_n)$  as the error. Then, let's assume the function  $f$  is analytic, and so does  $y(t)$ . By the Taylor Exp., we have

$$y(t_{n+1}) = y(t_n) + hf(t_n, y(t_n)) + \mathcal{O}(h^2) \quad (6)$$

Then, by subtracting the (4) from Euler's Method solution, we get

$$y_{n+1} - y(t_{n+1}) = y_n - y(t_n) + h[f(t_n, y_n) - f(t_n, y(t_n))] + \mathcal{O}(h^2) \quad (7)$$

$$y_n = y(t_n) + e_{n,h}$$

$$e_{n+1,h} = e_{n,h} + h[f(t_n, y_n) - f(t_n, y(t_n))] + \mathcal{O}(h^2) \quad (8)$$

Now consider the  $\mathcal{O}(h^2)$ . Because  $y$  analytic,  $\mathcal{O}(h^2)$  can be bounded for all  $h > 0$  and  $n < \lfloor t^*/h \rfloor$  by a term  $ch^2$  where  $c > 0$  as  $\mathcal{O}(h^2) \leq ch^2$ . Consequently, we get

$$\|e_{n+1,h}\| \leq \|e_{n,h}\| + h\|f(t_n, y_n) - f(t_n, y(t_n))\| + ch^2 \quad (9)$$

By the Lipschitz Condition,  $\exists \lambda > 0$  such that

$$\|e_{n+1,h}\| \leq \|e_{n,h}\| + h\|y_n - y(t_n)\| + ch^2 \quad (10)$$

Now we claim, which can be proved by induction, that

$$\|e_{n,h}\| \leq \frac{ch}{\lambda} [(1 + \lambda h)^{n+1} - 1] + ch^2 \quad (11)$$

By putting the claim into the (8), we obtain

$$\|e_{n+1,h}\| \leq \frac{ch}{\lambda} [(1 + \lambda h)^{n+1} - 1] \quad (12)$$

To reduce the dependencies on the  $h$ , we use  $e^{nh\lambda} \leq e^{\lfloor t^*/h \rfloor \lambda} < e^{t^*\lambda}$  (For the required algebra to get the inequality and (10), check the handwritten notes.)

$$\Rightarrow \|e_{n+1,h}\| < \frac{ch}{\lambda}(e^{t^*\lambda} - 1)$$

$$\lim_{h \rightarrow 0} \frac{ch}{\lambda}(e^{t^*\lambda} - 1) = 0$$

,so

$$\lim_{h \rightarrow 0} \|e_{n+1,h}\| = 0$$

□

One of the most important results of the proof is that the error is bounded as

$$\|e_{n+1,h}\| < \frac{ch}{\lambda}(e^{t^*\lambda} - 1) \quad \forall n \in 0, 1, \dots, \lfloor t^*/h \rfloor \quad (13)$$

However, this boundary is a lot bigger than the actual error. If we write the exact solution into Euler's Method (3) and Taylor expands the  $y(t_{n+1})$  term, we get

$$[y(t_n) + hf(t_n, y(t_n)) + \mathcal{O}(h^2)] - [y(t_n) + hf(t_n, y(t_n))] = \mathcal{O}(h^2) \quad (14)$$

This is to say that *The Euler's Method is an order of one method*. In general, for the given method

$$y_{n+1} = Y(f, h, y_0, \dots, y_n)$$

The method is called the order of  $p$  if

$$y(t_{n+1}) - Y(f, h, y(t_0), \dots, y(t_n)) = \mathcal{O}(h^{p+1})$$

For an order of  $p$  method, the local error (error between two time steps) decays as  $\mathcal{O}(h^{p+1})$ . Also, the number of steps increases as  $\mathcal{O}(h^{-1})$ . However, our main interest is how the global error behaves for a given method. Assuming that the global error decreases as  $\mathcal{O}(h^p)$  may be reasonable for many cases. Yet, this is not true for all methods, but it is still valid for Euler's method. The error decreases as  $\mathcal{O}(h)$ .

### 3 Trapezoidal Rule

A better approximation for the  $y(t) = y(t_n) + \int_{t_n}^t f(t, y(t))dt$  is to take the average of the function values at the integral limit.

$$y(t) \approx y(t_n) + \frac{1}{2}(t - t_n)[f(t, y(t)) + f(t_n, y(t_n))] \quad (15)$$

Using this approximation method to estimate the solution of a given ODE is called the Trapezoidal Rule.

$$y_{n+1} = y_n + \frac{1}{2}h[f(t_n, y_n) + f(t_{n+1}, y_{n+1})] \quad (16)$$

The order of the trapezoidal rule is two, which can be calculated using the same method as (12). As a difference, one should also Taylor expand the  $y'(t_{n+1}) = f(t_{n+1}, y_{n+1})$  around the  $t_n$ .

$$y(t_{n+1}) - y(t_n) - \frac{1}{2}h[f(t_n, y(t_n)) + f(t_{n+1}, y(t_{n+1}))] = \mathcal{O} \quad (17)$$

$$y(t_{n+1}) = y(t_n) + hf(t_n, y(t_n)) + \frac{1}{2}hy''(t_n) + \mathcal{O}(h^3) \quad (18)$$

$$f(t_{n+1}, y_{n+1}) = y'(t_n) + hy''(t_n) + \mathcal{O}(h^2) \quad (19)$$

By putting (16) and (17) into the (15), we get the difference is  $\mathcal{O}(h^3)$ .

**Theorem 3.1.** *Trapezoidal rule is convergent.*

*Proof.* !!!!! □

As a result, we can say that the global error is decreasing as  $\mathcal{O}(h^2)$  because the method is convergent.

Another important difference between Euler Method and Trapezoidal Rule is that Euler Method is an explicit method, whereas Trapezoidal Rule is implicit. This means that, to calculate  $y_{n+1}$  in each step, we should solve an algebraic equation

$$y_{n+1} - \frac{1}{2}hf(t_{n+1}, y_{n+1}) = y_n + \frac{1}{2}hf(t_n, y_n) \quad (20)$$

### 3.1 Implicit Midpoint Method

The Implicit Midpoint method, a second order convergent method, is a special form of the Runge-Kutta Method. It is pretty simple to Euler's Method. The difference between the two methods is that the Implicit Midpoint Method uses the approximation

$$\int_{t_n}^{t_{n+1}} f(t, y(t))dt = hf\left(t + \frac{h}{2}, \frac{1}{2}(y(t_n) + y(t_{n+1}))\right), \quad h = t_{n+1} - t_n$$

Hence the actual method is

$$y_{n+1} = y_n + hf\left(t_n + \frac{h}{2}, \frac{1}{2}(y_n + y_{n+1})\right) \quad (21)$$

## 4 The Theta Method

Theta method is the generalized method of Euler's Method and the Trapezoidal Rule. It is in the form

$$y_{n+1} = y_n + h(\theta f(t_n, y_n) + (1 - \theta)f(t_{n+1}, y_{n+1})) , \quad \text{where } \theta \in [0, 1] \quad (22)$$

The  $\theta = 1$  case is Euler's Method. The  $\theta = \frac{1}{2}$  case is the Trapezoidal Rule

### 4.1 Order of the Theta Method

The order of the Theta Method depends on the choice of the theta value. We can again calculate it by Tylor expanding  $y(t_{n+1})$  and the method then subtracting from each other.

$$y(t_{n+1}) = y(t_n) + hy'(t_n) + \frac{h^2}{2}y''(t_n) + \frac{h^3}{6}y^{(3)}(t_n) + \mathcal{O}(h^4) \quad (23)$$

By expanding the Theta method

$$y(t_{n+1}) \approx y(t_n) + h[\theta y'(t_n) + (1 - \theta)(y'(t_n) + hy''(t_n) + \frac{h^2}{2}y^{(3)}(t_n) + \mathcal{O}(h^3))] \quad (24)$$

Then, by subtracting (22) from (21), we get

$$h^2(\frac{1}{2} - \theta)y''(t_n) + h^3(\frac{\theta}{2} - \frac{1}{2})y^{(3)}(t_n) + \mathcal{O}(h^4) \quad (25)$$

As the result of the calculation, the order of the Theta Method is 2 if  $\theta = \frac{1}{2}$ . Otherwise, it is always the order of one.

### 4.2 Error Analysis of the Theta Method

If we put the exact solution  $y(t_{n+1})$  into the method (20), then subtract from the numeric method as  $y(t_{n+1}) - \mathcal{Y}(t_n, y_n, f)$  we get

$$e_{n+1} - e_n - h\theta[f(t_n, y_n) - f(t_n, y(t_n))] - h(1 - \theta)[f(t_{n+1}, y_{n+1}) - f(t_{n+1}, y(t_{n+1}))], \quad \text{where } e_n = y_n - y(t_n)$$

After Taylor expanding the  $[f(t_{n+1}, y_{n+1}) - f(t_{n+1}, y(t_{n+1}))]$  around  $t_n$ , we get the following result according to the above order calculation

$$e_{n+1} - e_n = \left\{ \begin{array}{ll} -\frac{h^3}{12}y^{(3)}(t_n) + \mathcal{O}(h^4) & , \theta = \frac{1}{2} \\ (\theta - \frac{1}{2})h^2y^{(2)}(t_n) + \mathcal{O}(h^2) & , \theta \neq \frac{1}{2} \end{array} \right\} \quad (26)$$