Our starting point is the following ODE:

$$\begin{array}{lll} y'(t) &=& f(t,y(t)), \ a \leq t \leq b \\ y(a) &=& \eta. \end{array} \right\}$$
 (1)

Why do we need numerical methods for solving such an "innocuous" differential equation? Because problems of this form are often difficult, if not impossible, to solve analitically. This is specially the case when the equation is nonlinear. Although we will concentrate on first-order scalar equations, the numerical methods are easily extended to cover first-order systems of equations. Furthermore, it is easy to convert higher order scalar initial value problems to first-order systems. Therefore, we are not restricted to first order equations. In fact, the methods can be used as part of numerical algorithms to find solutions of evolutionary **partial** differential equations. Additionally, we will consider also methods for solving **Boundary Value Problems**.

Solving Ordinary Differential equations

Outline of Study of ODE's:

Single-Step Methods for I.V. Problems

- a. Euler.
- b. Trapezoidal.
- c. Taylor methods.
- d. General Runge-Kutta
- e. Adaptive step-size control.

2 Stiff ODEs

- Boundary Value Problems
 - 1. Shooting Method
 - 2. Finite Difference Method

Discretisation

The central idea behind numerical methods is that of discretisation. That is we partition the continuous interval [a, b] by a discrete set of N + 1 points:

$$a = t_0 < t_1 < t_2 < \ldots < t_{N-1} < t_N = b.$$

The parameters

$$h_n = t_{n+1} - t_n, \quad n = 0, 1, ..., N - 1$$
 (2)

are called the step-sizes. We will be often be interested in using an equally spaced partition where

$$h_n = h = \frac{(b-a)}{N}, \ n = 0, 1, ..., N-1.$$

・ロト・母・・由・・ 日・ うらの

We will let y_n denote the numerical approximation to the exact solution $y(t_n)$. A numerical solution of (1) consists of a set of discrete approximations $\{y_n\}_{n=0}^N$. A numerical method is a difference equation involving a number of consecutive approximations

$$y_j, j = 0, 1, ..., k$$

from which we sequentially compute the sequence

$$y_{k+n}, n = 1, 2, ..., N$$

▲ □ ▶ ▲ □ ▶ ▲ □ ▶

The derivation of a number of numerical methods begins by integrating (1) between t_n and t_{n+1} . This gives

$$\int_{t_n}^{t_{n+1}} \frac{dy}{dt} dt = \int_{t_n}^{t_{n+1}} f(t, y) dt$$
$$\Rightarrow y(t_{n+1}) - y(t_n) = \int_{t_n}^{t_{n+1}} f(t, y) dt.$$

Now, if we make the approximation

$$f(t, \mathbf{y}) \approx f(t_n, \mathbf{y}(t_n)), \ t \in (t_n, t_{n+1})$$

then

$$y(t_{n+1}) - y(t_n) \approx \int_{t_n}^{t_{n+1}} f(t_n, y(t_n)) dt = (t_{n+1} - t_n) f(t_n, y(t_n)).$$

Therefore

$$y(t_{n+1}) \approx y(t_n) + (t_{n+1} - t_n)f(t_n, y(t_n)).$$

This suggest the numerical method

$$y_{n+1} = y_n + (t_{n+1} - t_n)f(t_n, y(t_n)), \quad n = 0, 1, ..., N-1$$
 (3)

which is called the forward or explicit Euler method. Note that from the initial condition $y_0 = \eta$ we can explicitly calculate y_1 by applying (3). This in turn allows us to calculate y_2 and then y_3 and so on. A geometrical interpretation of the forward Euler method is that instead of following the possibly curved solution trajectory passing through (t_n, y_n) , the method actually follows a straight line trajectory which has slope $f(t_n, y_n)$.

The forward Euler method is, of course, an approximate method which will only be exact in the trivial case of a linear solution in t. But we hope that the method will be closer to the exact solution as the step-size h is taken smaller. This is a neccesary condition for any reasonable numerical method.

Definition

We will say that a numerical method is convergent when for all IVP (1) with solution sufficiently differentiable, the following condition applies

$$\lim_{h\to 0}\left(\max_{1\leq n\leq N}||y_n-y(t_n)||\right)=0.$$

being $y_0 = y(t_0)$. We will say that the order of conve

We will say that the order of convergence of the method is p, if as h is taken smaller (ie, with N large enough), then

$$\left(\max_{1\leq n\leq N}||y_n-y(t_n)||\right)=\mathcal{O}(h^p), \ Nh= ext{constante}$$

ie, if there $\exists C \ge 0$ such that $\max_{1 \le n \le N} ||y_n - y(t_n)|| \le C|h|^p$ for N large enough.

Theorem

For all IVP (1) with f continuous and satisfying a Lipschitz condition on D, the forward Euler method is convergent and its order of convergence is 1. The error of the Euler method can be bound as follows

$$||y(t_n) - y_n|| \le \frac{C}{2L} (e^{(t_n - a)L} - 1)h, \ 0 \le n \le N$$

being $y_0 = y(t_0)$, $C = \max_{x \in [a,b]} ||y''(x)||$ and L a Lipschitz constant.

▲□▶▲□▶▲□▶▲□▶ □ シペ?

An obvious question is whether we can easily improve upon the forward Euler method. Remembering that Euler method replaces f(t, y(t)) by the slope $f(t_n, y_n)$, it seems likely that an improved approximation would be the average of the slopes at t_n and t_{n+1} . That is

$$y(t_{n+1}) - y(t_n) \approx (t_{n+1} - t_n) \frac{1}{2} (f(t_n, y(t_n)) + f(t_{n+1}, y(t_{n+1})))$$

This suggest the following numerical method:

$$y_{n+1} = y_n + \frac{h_n}{2} \left(f(t_n, y_n) + f(t_{n+1}, y_{n+1}) \right) . \tag{4}$$

This method is called the trapezoidal method and differs from the forward Euler method in an important way:

At the (n + 1)st step we have to solve the (generally no linear) equation

$$g(y_{n+1}) \equiv y_{n+1} - y_n - \frac{h_n}{2} \left(f(t_n, y_n) + f(t_{n+1}, y_{n+1}) \right) = 0$$
 (5)

to determine y_{n+1} . Therefore y_{n+1} is defined implicitly and for that reason the trapezoidal method is an example of an implicit method. The forward Euler method, on the other hand, is an example of an explicit method.

For the trapezoidal method, the following convergence theorem can be stablished

Theorem

For all IVP (1) satisfying a Lipschitz condition, the trapezoidal method is convergent and for hL < 2 (being L a Lipschitz constant) the error can be bound by

$$|e_n| \leq rac{Ch^2}{L} \exp\left(rac{L(t_n-a)}{1-rac{hL}{2}}
ight),$$

where |y'''| < C. Therefore, the trapezoidal method is convergent with order 2.

Solving Ordinary Differential equations

Taylor methods can be used to build explicit methods with higher order of convergence than Euler's method. The main difficult of these methods is the computation of the derivatives of f(t, y). The idea behind these methods is simple: considering the Taylor series of $y(t_{n+1})$ up to a certain order. For example, by truncating up to order h^2 we will have the Taylor method of order 2. Let us build this method. Consider:

$$y(t_{n+1}) = y(t_n) + hy'(t_n) + \frac{h^2}{2}y''(t_n) + \frac{h^3}{6}y'''(\zeta_n)$$

We will take up to order h^2 . For computing $y'(t_n)$ and $y''(t_n)$ we will use the differential equation:

y'(t) = f(t, y(t)) $y''(t) = f_t(t, y(t)) + f_y(t, y(t))y'(t) = f_t(t, y(t)) + f_y(t, y(t))f(t, y(t))$

Therefore we have the method

$$y_{n+1} = y_n + hf(t_n, y_n) + \frac{h^2}{2}(f_t(t_n, y_n) + f_y(t, y_n)f(t_n, y_n))$$

◆□▶ ◆□▶ ◆ □ ▼ ◆ □ ▼ ◆ □ ▼ ◆ □ ▼

Solving Ordinary Differential equations

By defining as the local truncation error the term which is neglected for building the method, we see that in our method this error is $T_n = \frac{h^3}{6} y'''(\zeta_n)$, which satisfies $||T_n|| \le Ch^3$. This means that the method is order 2.

Error Analysis

- Truncation Error (Truncating Taylor Series) [Large step size \implies large errors]
- O Rounding Error (Machine precision) [very small step size → roundoff errors]

Two kinds of Truncation of Error:

- Local error within one step due to application of method
- **Propagation** error due to previous local errors.

Global Truncation Error = Local + Propagation Generally if local truncation error is $\mathcal{O}(h^{(n+1)})$ then, the Global truncation error is $\mathcal{O}(h^{(n)})$.

General Single-step methods:

Before introducing the Runge-Kutta methods, we will briefly describe the general set-up for Single-step methods. In general, these methods can be written as:

$$y_{n+1} = y_n + h\Phi(x_n, y_n, y_{n+1}, h)$$
 (1)

where Φ is related to *f* and it is called the Step Function. When Φ depends on y_{n+1} the method is implicit and in other case the method is explicit.

The examples of single-step methods considered so far are:

- **1** Forward Euler, $\Phi(x_n, y_n, y_{n+1}, h) = f(x_n, y_n)$ (explicit).
- 2 Trapezoidal, $\Phi(x_n, y_n, y_{n+1}, h) = (f(x_n, y_n) + f(x_{n+1}, y_{n+1}))/2$ (implicit).
- Taylor of order 2,

 $\Phi(x_n, y_n, y_{n+1}, h) = f(t_n, y_n) + \frac{h}{2}(f_t(t_n, y_n) + f_y(t, y_n)f(t_n, y_n) + f_y(t_n, y_n))$ (explicit).

For Single-step methods

Definition

We define the local truncation error, T_n , as:

$$T_n(h) = y(t_{n+1}) - y(t_n) - h\Phi(t_n, y(t_n), y(t_{n+1}), h)$$
(2)

・ロト ・聞 ト ・ 国 ト ・ 国 トー

크

Runge-Kutta methods are very popular methods which allow us to obtain high-order methods avoiding the computation of the derivatives of the Taylor methods.

Let's see how to build an explicit Runge-Kutta method of order 2:

$$y_{n+1} = y_n + h\Phi(t_n, y_n, h) \Phi(t_n, y_n, h) = ak_1 + bf(t_n + \alpha h, y_n + \beta hk_1), \ k_1 = f(t_n, y_n)$$
(3)

The constants of the Step function will be determine by imposing that the method behaves like the Taylor method up to order 2. The Taylor method of order 2 was:

$$y_{n+1} = y_n + hf(t_n, y_n) + \frac{h^2}{2}(f_t(t_n, y_n) + f_y(t, y_n)f(t_n, y_n))$$

Then, in order to compare with the new method, we just need to expand $f(t_n + \alpha h, y_n + \beta h k_1)$ up to order 1:

$$\begin{aligned} f(t_n + \alpha h, y_n + \beta hk_1) &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta hk_1 + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) f(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_y(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_t(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_t(t_n, y_n) \beta h + f_t(t_n, y_n) \beta h + \mathcal{O}(h^2) \\ &= f(t_n, y_n) + f_t(t_n, y_n) \alpha h + f_t(t_n, y_n) \beta h + f_$$

Going to Eq. (3), we have that up to order 2 the new method is:

$$y_{n+1} = y_n + h(a+b)f(t_n, y_n) + bh^2(\alpha f_t(t_n, y_n) + \beta f(t_n, y_n)f_y(t_n y_n))$$

Then, if the method must behave as the Taylor method up to order 2, we have the following equations for the parameters of the method:

$$a+b=1$$
, $b\alpha=b\beta=1/2$

The methods verifying these equations will be convergents with order of convergence 2.

Consider, for instance, the methods with $\alpha = \beta$. Then, $b = 1/2\alpha$ and $a = 1 - 1/2\alpha$. The methods with $\alpha = 1/2$ and $\alpha = 1$ are called the Modified Euler Method and the Heun Method, respectively.

The Modified Euler Method ($\alpha = 1/2$) can be written as

$$y_{n+1} = y_n + hf(t_n + \frac{h}{2}, y_n + \frac{h}{2}f_n).$$
 (4)

and using the notation

$$k_1 = f(t_n, y_n), \quad k_2 = f(t_n + \frac{h}{2}, y_n + \frac{h}{2}k_1),$$

we have that

$$y_{n+1}=y_n+hk_2.$$

▲口 ▶ ▲ □ ▶ ▲ 三 ▶ ▲ □ ▶ ▲ □ ▶ ▲ □ ▶

The Heun Method ($\alpha = 1$), can be written as

$$y_{n+1} = y_n + \frac{h}{2} \left(f(t_n, y(t_n)) + f(t_{n+1}, y_n + hf(t_n, y_n)) \right) .$$
 (5)

lf

$$k_1 = f(t_n, y_n), \quad k_2 = f(t_n + h, y_n + hk_1),$$

then, we have:

$$y_{n+1} = y_n + \frac{h}{2}(k_1 + k_2).$$

▲口▶▲圖▶▲圖▶▲圖▶ ▲国▼

The previous notation can be generalized

$$k_1 = f(t_n, y_n), \ k_2 = f(t_n + c_2 h, y_n + a_{21} h k_1),$$

and

$$y_{n+1} = y_n + h(b_1k_1 + b_2k_2),$$

where c_2 , a_{21} , b_1 and b_2 are constants.

All these methods are members of a big family of methods which are called the Runge-Kutta Methods.

イロト イヨト イヨト イヨト

The general Runge-Kutta s-Stages Method is

$$y_{n+1}=y_n+h\sum_{i=1}^s b_ik_i\,,$$

where

$$k_i = f\left(t_n + c_i h, y_n + h \sum_{j=1}^{s} a_{ij} k_j\right), \ i = 1, 2, ..., s.$$

A convenient way of displaying the coefficients of the Runge-Kutta Methods is the use of the Butcher tableaux:

Example

Find the Butcher tableaux of the following Runge-Kutta Method:

$$y_{n+1} = y_n + h\left(\frac{1}{4}k_1 + \frac{3}{4}k_2\right)$$
$$k_1 = f\left(t_n, y_n + h\left(\frac{1}{4}k_1 - \frac{1}{4}k_2\right)\right),$$
$$k_2 = f\left(t_n + \frac{2}{3}h, y_n + h\left(\frac{1}{4}k_1 + \frac{5}{12}k_2\right)\right).$$

Sol.:



If $p^*(s)$ is the maximum order of convergence which can be obtained by using an explicit Runge-Kutta *s*-stages method, then we have

$$\begin{array}{rcl} p^*(s) &=& s, \ s=1,2,3,4\\ p^*(5) &=& 4\\ p^*(6) &=& 5\\ p^*(7) &=& 6\\ p^*(8) &=& 6\\ p^*(9) &=& 7\\ p^*(s) &\leq& s-2, \ s=10,11,\ldots \end{array}$$

This behaviour explains the popularity of the 4-stages Runge-Kutta Methods of order 4.

Adaptive Step-size Control

Goal: with little additional effort estimate (bound) the magnitude of local truncation error at each step so that step size can be reduced/increased if local error increases/decreases. Basic idea: 2. Use a matched pair of Runge-Kutta formulas of order p and p + 1 that use common values of ki, and yield estimate or bound

local truncation error.

R-K Method of Order *p*:

$$y_{n+1} = y_n + h \sum_{i=1}^s b_i k_i$$

R-K method of Order p + 1:

$$y_{n+1} = y_n + h \sum_{i=1}^s \tilde{b}_i k_i$$

These methods are called embedded Runge-Kutta Methods and can be expressed by using the following modified Butcher tableaux:

$$\frac{\begin{array}{c|c} \boldsymbol{c} & \boldsymbol{A} \\ \hline \boldsymbol{b}^{T} \\ \hline \boldsymbol{\tilde{b}}^{T} \\ \hline \end{array}}{\boldsymbol{b}^{T} - \boldsymbol{\tilde{b}}^{T}}$$

イロト イヨト イヨト イヨト

The local truncation error for the order-*p* method can be estimated as follows:

$$ilde{T}_{n+1} = h \sum_{i=1}^s (b_i - ilde{b}_i) k_i \, .$$

Then, the step-size can be adapted by imposing the following condition:

$$\tilde{T}_{n+1} < Tol,$$

being *Tol* the tolerance by step unit. The "new" step size is then selected as follows:

$$h_{new} = \left(\frac{q\text{Tol}}{|\tilde{T}_{n+1}|}\right)^{\frac{1}{p+1}}h.$$
 (6)

where *q* is a factor with a value of $q \approx 0.8$.

Example: The Runge-Kutta-Fehlberg (4,5) scheme has the following Butcher tableaux



Higher-Order ODEs and Systems of Equations

An n-th order ODE can be converged into a system of n coupled 1st-order ODEs. Systems of first order ODEs are solved just as one solves a single ODE.

Example:

Consider the 4th-order ODE

$$y'''' + a(x)y''' + b(x)y'' + c(x)y' + d(x)y = f(x)$$

By letting

$$y''' = v_3; y'' = v_2; and y' = v_1$$

this 4th-order ODE can be written as a system of four coupled 1st-order ODEs.

Stiff Differential Equations

A stiff system of ODE's is one involving rapidly changing components together with slowly changing ones. In many cases, the rapidly varying components die away quickly, after which the solution is dominated by the slow ones. Consider the following ODE:

$$y' = -50(y - \cos(t)), \ 0 \le t \le 1.25, \ y(0) = 0,$$

The figure shows the approximations obtained by using the Forward Euler method for h = 1.25/31 and h = 1.25/32.

・ロト ・聞 ト ・ 国 ト ・ 国 ト



In order to ensure the stability of the numerical solution, we have to choose an step size h < 2/50. This means that, for stiff problems, we would need a very small step to capture the behavior of the rapid transient and to preserve a stable and accurate solution, and this would then make it very laborious to compute the slowly evolving solution.

Euler's method is known as an explicit method because the derivative is taken at the known point *i*. An alternative is to use an implicit approach, in which the derivative is evaluated at a future step i + 1. The simplest method of this type is the backward Euler method that yields

$$y_{n+1} = y_n + hf(y_{n+1}, t_{n+1}).$$

In general, for this kind of problems, suitable methods are the BDF *backward differentiation formulae* methods. These methods have the general form:

$$\frac{1}{h}\sum_{i=1}^{k}\frac{1}{j}\nabla^{j}y_{n+1}=f_{n+1}.$$

STANDARD 4-STAGES EXPLICIT RUNGE-KUTTA METHOD

The method is given by

$$k_{1} = f(t_{n}, y_{n})$$

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

$$k_{3} = f\left(t_{n} + \frac{h}{2}, y_{n} + \frac{h}{2}k_{2}\right),$$

$$k_{4} = f\left(t_{n+1}, y_{n} + hk_{3}\right),$$

$$y_{n+1} = y_{n} + \frac{h}{6}\left(k_{1} + 2k_{2} + 2k_{3} + k_{4}\right).$$