

III. Ordinary Differential Equations

- Goals:
- solve IVPs numerically
 - basic methods (explicit/implicit) and their characteristics
 - stability and stiff equations

Why? E.g. time evolution of CSTR

" " " " tubular reactor
MATLAB: ode45, ode23s, ...

III.1 Problem statement and examples

Def.: A scalar first order Initial Value

Problem (IVP) is given by

$$\frac{dy}{dt} = \dot{y}(t) = f(t, y(t))$$

scalar first order

Ordinary Differential Equation (ODE)

because only first derivative

$$y(t_0) = y_0$$

Initial Value (IV)

and an interval $I = [t_0, T]$ on which the solution $y(t)$ is sought.

initial final time

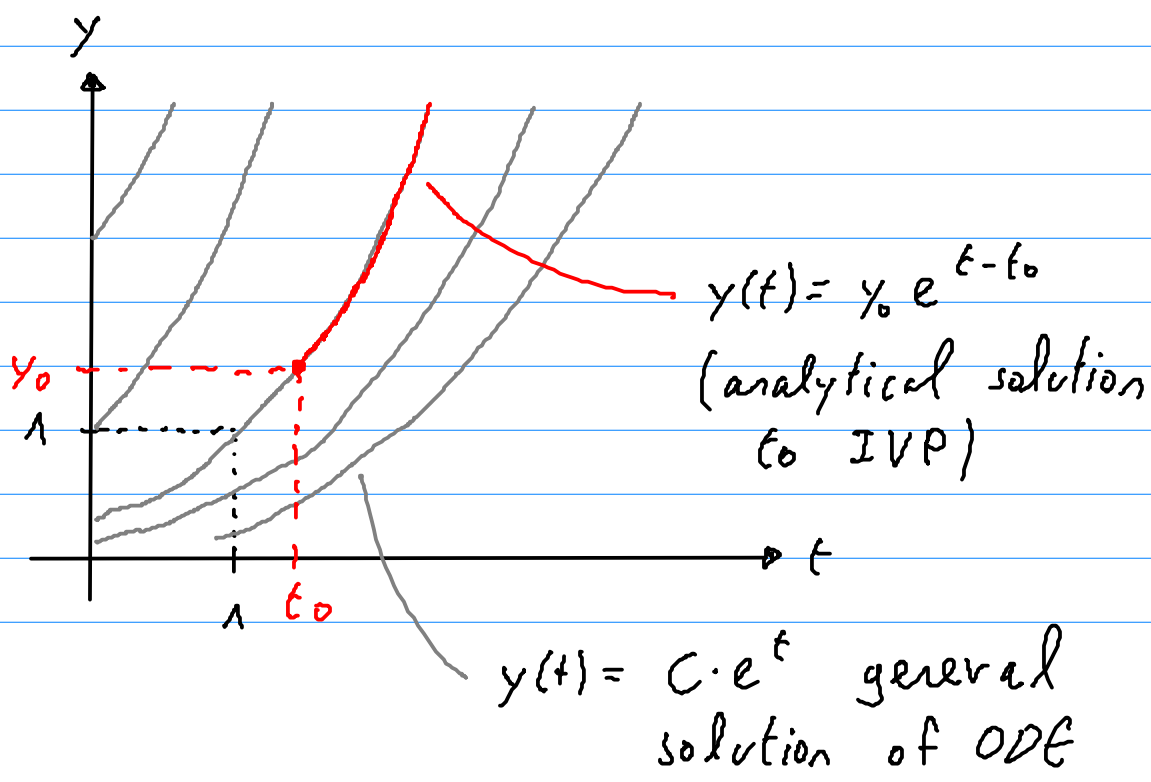
Rem.: (i) Chosen time t as the independent variable for convenience... But it could be anything else x ...

(ii) We seek an entire function

Ex.: (1)
$$\left. \begin{aligned} \dot{y}(t) &= y(t) \\ y(t_0) &= y_0 \end{aligned} \right\} \text{ IVP}$$

For each value of t and y the ODE prescribes the derivative \dot{y} of the solution

\rightsquigarrow slope fields (also direction field or vector field)



3

Often not only one ODE, i.e. scalar, but a whole system

Def.: A (general) first order IVP is given by

$$\dot{\vec{y}}(t) = \vec{f}(t, \vec{y}(t)) \quad \text{system of ODEs}$$

$$\vec{y}(t_0) = \vec{y}_0 \quad \text{IVs}$$

and an interval $I = [t_0, T]$ on which the solution $\vec{y}(t)$ is sought.

Ex.: (2) CSTR

(3) Linear system of ODEs

$$\dot{\vec{y}}(t) = A \vec{y}(t)$$

$$\vec{y}(t_0) = \vec{y}_0$$

————— A matrix

Analytical solution $\vec{y}(t) = e^{A(t-t_0)} \vec{y}_0$

$$= \left(\sum_{k=0}^{\infty} \frac{1}{k!} A^k (t-t_0)^k \right) \vec{y}_0$$

matrix exponential

What about n -th order IVPs?

Def.: A scalar n -th order IVP is given by

$$y^{(n)}(t) = f(t, y(t), \dot{y}(t), \dots, y^{(n-1)}(t))$$

n -th derivative
scalar n -th order ODE

$$y(t_0) = y_0$$

$$\dot{y}(t_0) = \dot{y}_0$$

⋮

$$y^{(n-1)}(t_0) = y_0^{(n-1)}$$

IVPs

and an interval $I = [t_0, T]$ on which the solution $y(t)$ is sought.

Ex.: (4) Equations of motion (Newton's 2nd Law)

$$m \ddot{x}(t) = F$$

mass
acceleration
force

$$x(t_0) = x_0 \quad \text{initial position}$$

$$\dot{x}(t_0) = v_0 \quad \text{initial velocity}$$

In the following we will only see numerical methods for first order IPV.

Reduction to first order system

Given an n -th order ODE

$$y^{(n)}(t) = F(t, y(t), \dot{y}(t), \ddot{y}(t), \dots, y^{(n-1)}(t))$$

We define

$$\begin{aligned} z_0(t) &= y(t) \\ z_1(t) &= \dot{y}(t) = \dot{z}_0(t) \\ z_2(t) &= \ddot{y}(t) = \dot{z}_1(t) \\ &\vdots \\ z_{n-1}(t) &= y^{(n-1)}(t) = \dot{z}_{n-2}(t) \end{aligned} \quad \left. \vphantom{\begin{aligned} z_0(t) \\ z_1(t) \\ z_2(t) \\ \vdots \\ z_{n-1}(t) \end{aligned}} \right\} \begin{array}{l} n-1 \text{ first} \\ \text{order ODEs} \end{array}$$

Now

$$\begin{aligned} z_n(t) &= y^{(n)}(t) = \underline{\dot{z}_{n-1}(t)} \\ &= F(t, y(t), \dot{y}(t), \ddot{y}(t), \dots, y^{(n-1)}(t)) \\ &= \underline{F(t, z_0(t), z_1(t), z_2(t), \dots, z_{n-1}(t))} \end{aligned}$$

\leadsto n first order ODEs!

By defining

$$\vec{z}(t) = \begin{pmatrix} z_0(t) \\ z_1(t) \\ \vdots \\ z_{n-2}(t) \\ z_{n-1}(t) \end{pmatrix} \quad \text{and} \quad \vec{g}(t, \vec{z}) = \begin{pmatrix} z_1(t) \\ z_2(t) \\ \vdots \\ z_{n-1}(t) \\ f(t, z_0(t), \dots, z_{n-1}(t)) \end{pmatrix}$$

we get a first order system of ODEs

$$\dot{\vec{z}}(t) = \vec{g}(t, \vec{z}(t))$$

Together with the IVs we get a first order IPV.

Ex.: (S) $\ddot{y} = \dot{y} + y^2 - e^t$, $t \in [t_0, T]$

$$y(t_0) = 1, \quad \dot{y}(t_0) = 0$$

Reduction to first order:

$$z_0(t) = y(t)$$

$$z_1(t) = \dot{y}(t) = \dot{z}_0(t)$$

$$z_2(t) = \ddot{y}(t) = \dot{z}_1(t)$$

$$= \dot{y}(t) + y^2(t) - e^t$$

$$= z_1(t) + z_0^2(t) - e^t$$

With

$$\vec{z}(t) = \begin{pmatrix} z_0(t) \\ z_1(t) \end{pmatrix}, \quad \vec{g}(t, \vec{z}(t)) = \begin{pmatrix} z_1(t) \\ z_1(t) + z_0(t) - e^t \end{pmatrix}$$

and IV

$$\vec{z}(t_0) = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

we have a first order IVP ✓.

Solving IVP analytically is in general difficult or even impossible:

$$\frac{dy(t)}{dt} = f(t, y(t)) \quad \left| \cdot dt \right., \quad \int_{t_0}^t$$

...

$$\underline{y(t)} = y(t_0) + \int_{t_0}^t f(\tau, y(\tau)) d\tau$$

"solution depends on solution..."

~ Numerical methods

$$\approx y(t_0) + \sum_{\substack{\text{steps} \\ \text{until} \\ t}} f(\tau, y_\tau) \cdot \Delta\tau$$

(time) / step size
approx. sol. at τ

III.2 The Euler methods

Consider the scalar IVP

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

$$y(t_0) = y_0$$

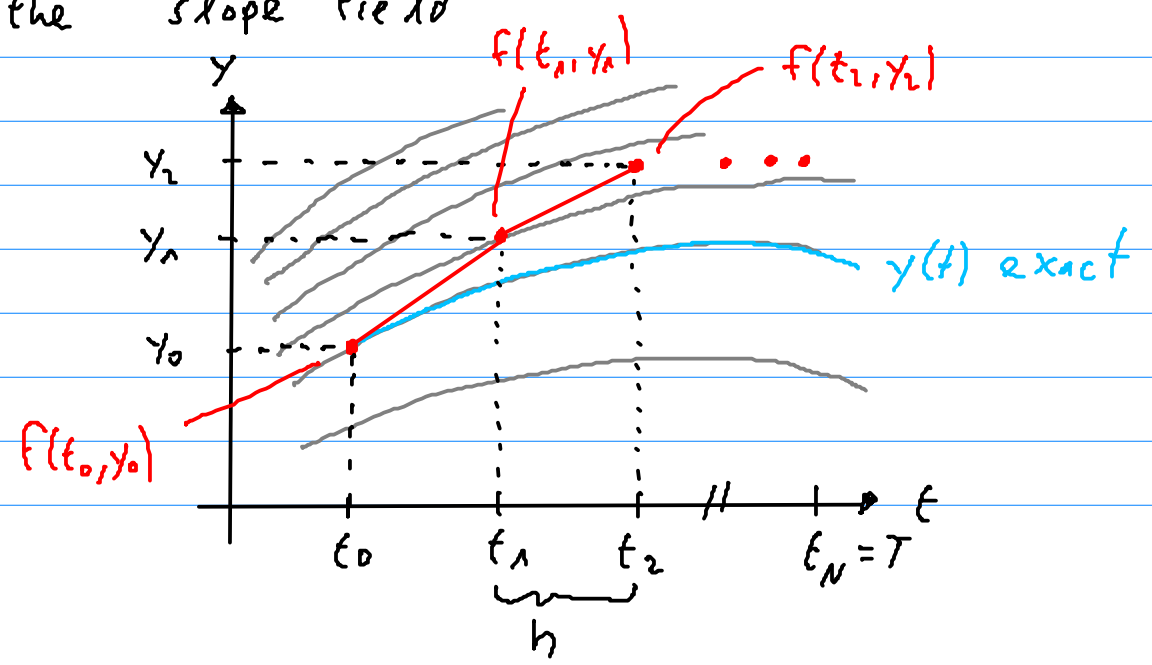
Goal: Find approximation of the solution $y(t)$ for $t_0 \leq t \leq T$.

Idea: Partition the interval $[t_0, T]$ regularly with

$$t_j = t_0 + j \cdot h \quad \text{step size}$$

$$= t_0 + j \cdot \frac{T - t_0}{N} \quad j = 0, 1, \dots, N-1$$

and approx. the solution by "following" the slope field



9

$$\rightsquigarrow y_1 = y_0 + h \cdot f(t_0, y_0)$$

$$y_2 = y_1 + h \cdot f(t_1, y_1)$$

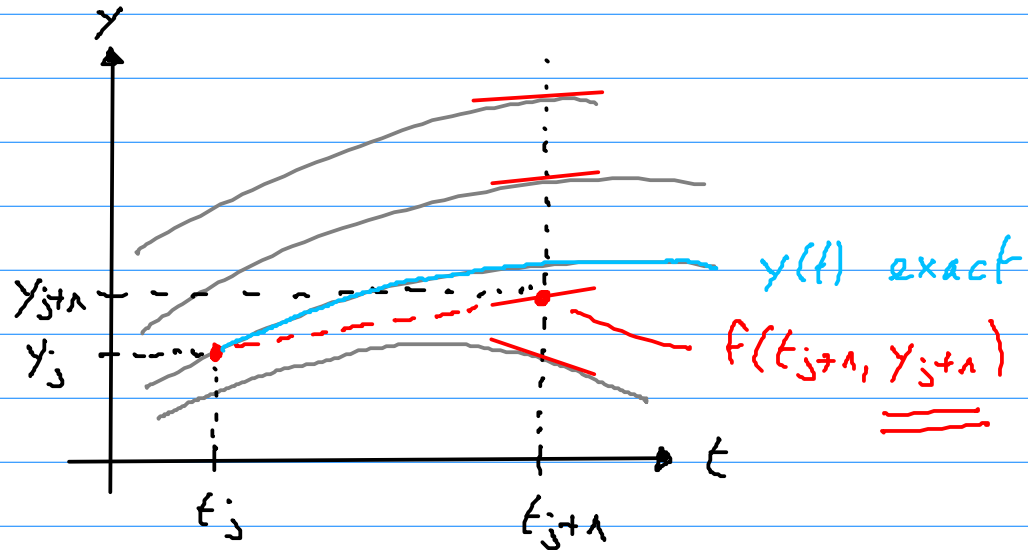
$$\vdots$$

We obtain the (explicit) Euler method (EE)

$$y_{j+1} = y_j + h \cdot f(t_j, y_j), \quad j=0, 1, \dots, N-1$$

A.k.a. Forward Euler

A different possibility



We obtain the implicit Euler method (IE)

$$\underline{y_{j+1}} = y_j + h \cdot \underline{f(t_{j+1}, y_{j+1})}$$

03.10.24

A.k.a. backward Euler

Need to solve
for y_{j+1} !
no implicit

Ex.: (6) $\dot{y}(t) = -y(t) + 2 \cdot \cos(t)$
 $y(0) = 1$

~ Slides (Euler Methods)

Rem.: For systems of ODEs, just replace
 y, f by \vec{y}, \vec{f}

III.3 Error estimation and convergence

To analyze the accuracy of so-called one-step methods we consider the general expression

$$y_{j+1} = y_j + h \cdot \phi(t_j, y_j, y_{j+1}, h)$$

and call ϕ the increment function.

Rem.: (i) $\phi(t_j, y_j, y_{j+1}, h) = f(t_j, y_j) \quad \text{EE}$

(ii) $\phi(t_j, y_j, y_{j+1}, h) = f(t_{j+1}, y_{j+1}) \quad \text{IE}$

Def.: the local truncation error (LTE) at t_j is defined by

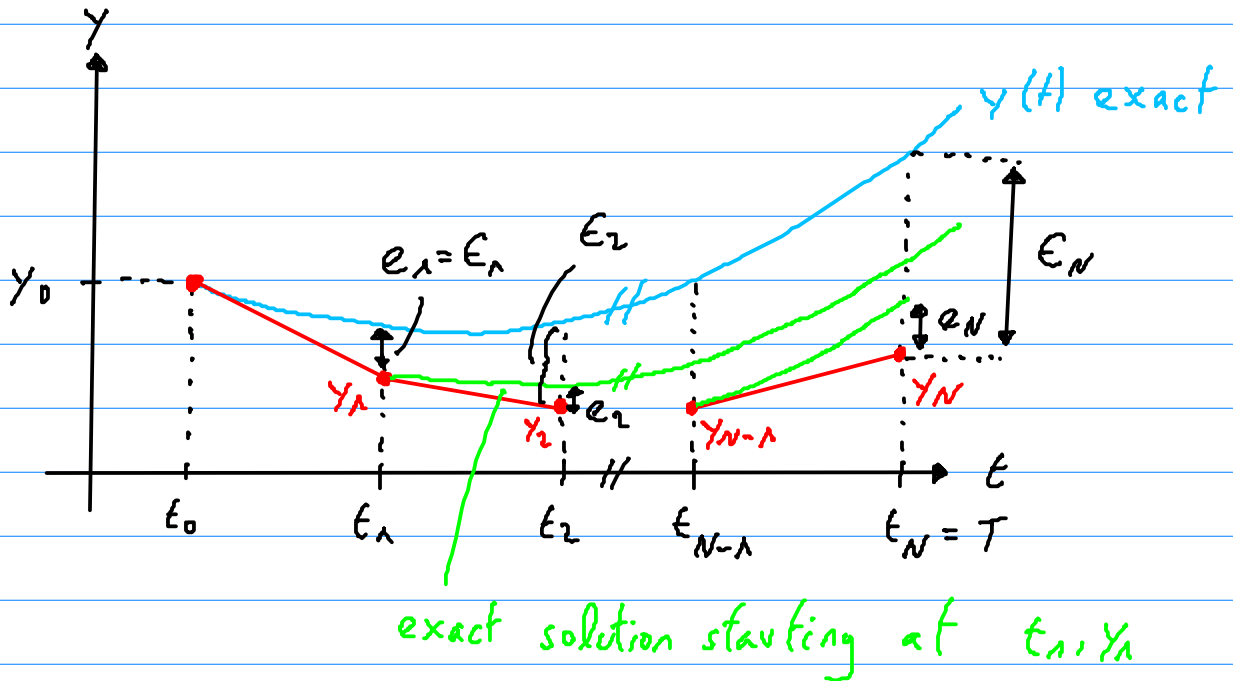
$$e_j = \underset{\substack{\uparrow \\ \text{exact sol. at } t_j}}{y(t_j)} - \left(\underset{\substack{\uparrow \\ t_{j-1}}}{y(t_{j-1})} + h \cdot \phi(t_{j-1}, y(t_{j-1}), y(t_j), h) \right)$$

The LTE is the error after one step of a method executed with the exact solution.

Def.: the global truncation error (GTE) at t_j is defined by

$$E_j = \underset{\substack{| \\ \text{exact}}}{y(t_j)} - \underset{\substack{| \\ \text{approx. solution at } t_j}}{y_j}$$

Graphically:



The LTE can "easily" be obtained by a Taylor expansion:

Ex.: (7) LTE of $\text{EE } \phi(t_j, y_j, y_{j+n}, h) = f(t_j, y_j)$

$$e_j \stackrel{\text{Def.}}{=} y(t_j) - \left(y(t_{j-1}) + h \cdot \phi(t_{j-1}, y(t_{j-1}), y(t_j)) \right)$$

$$= y(t_j) - y(t_{j-1}) - h \cdot f(t_{j-1}, y(t_{j-1}))$$

$$= y(t_{j-1} + h) - y(t_{j-1}) - h \cdot f(t_{j-1}, y(t_{j-1}))$$

$$\text{Taylor} \quad = \cancel{y(t_{j-1})} + \cancel{h \cdot \dot{y}(t_{j-1})} + \frac{h^2}{2} \ddot{y}(t_{j-1}) + \frac{h^3}{6} \dddot{y}(t_{j-1}) + \dots$$

ODE $f(t_{j-1}, y(t_{j-1}))$

$$- \cancel{y(t_{j-1})} - h \cdot \cancel{f(t_{j-1}, y(t_{j-1}))}$$

$$= \frac{h^2}{2} \ddot{y}(t_{j-1}) + h^3(\dots) + h^4(\dots) + \dots = \mathcal{O}(h^2)$$

Ex.: (8) LTE of IE $\phi(t_j, y_j, y_{j+1}, h) = f(t_{j+1}, y_{j+1})$

$$e_j \stackrel{\text{def.}}{=} y(t_j) - y(t_{j-1}) - h \cdot f(t_j, y(t_j))$$

$$= \dots$$

$$= \mathcal{O}(h^2)$$

$y(t_{j-h}) = \dots$ Taylor...

Rem.: The LTE depends on the smoothness of the solution (e.g. y above!) and therefore on the smoothness of the right-hand side function $f(t, y(t))$:

$$\ddot{y}(t) = \frac{d}{dt}(\dot{y}(t))$$

$$\stackrel{\text{ODE}}{=} \frac{d}{dt}(f(t, y(t)))$$

$$= \frac{\partial f}{\partial t} + \frac{\partial f}{\partial y} \frac{dy}{dt} = \frac{\partial f}{\partial t} + \frac{\partial f}{\partial y} f$$

$$\underbrace{\hspace{10em}}_{\text{ODE } \dot{y} = f(t, y)}$$

So we have seen that the LTE error is easy to compute with Taylor series...

How is the LTE related to the practically more meaningful CTE?

Ex.: (3) EE & IE GTE

→ sliders (Euler Methods)

We observe that although the LTE behaves as $O(h^2)$, the GTE is $O(h)$

Actually, one can show (under requirements that are anyway necessary for the existence and uniqueness of solutions to an IVP) that the LTEs accumulate in each step:

$$|E_N| \leq N \cdot \max_{1 \leq j \leq N} |e_j| = O(h)$$

\uparrow $\sim \frac{1}{h}$ \uparrow $O(h^2)$

And we have convergence: $y_N \xrightarrow{h \rightarrow 0} y(T)$

This motivates the following definition

Def.: We say that a method has order of accuracy p if the LTE is

$$|e_j| = O(h^{p+1})$$

or equivalently

$$|E_j| = O(h^p)$$

Rem.: EE and IE have order of accuracy $p=1$.

III.4 Runge-Kutta methods

So far: Euler methods with error $\mathcal{O}(h)$

→ want 100 x smaller error

↳ ? x smaller step size h

100

100

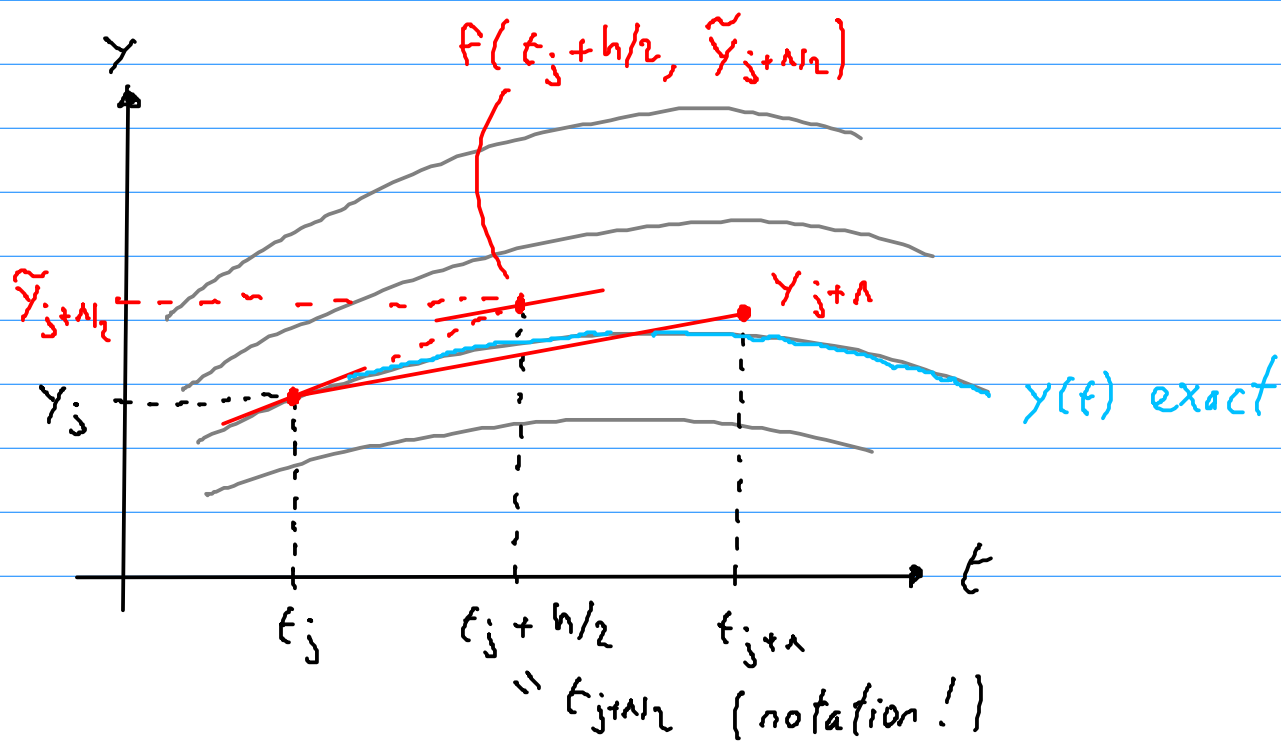
? x more (computational) work

Can we do better? YES! ▽

Try: $y_{j+1} = y_j + h \cdot f\left(t_j + \frac{h}{2}, \tilde{y}_{j+1/2}\right)$

↑
Approx. of sol. at $t_j + h/2$?

Let's go back to the slope field:



Idea: approx. $\tilde{y}_{j+1/2}$ with half EE step

Is this better?

Ex.: (10) Approx. IVP of Ex. (6) with above scheme

~> slides (Runge's Method)

We observe in the above numerical experiment that the LTE of this scheme is $\mathcal{O}(h^2)$. Hence its order of accuracy is $p=2$ (and it is indeed more accurate than the Euler methods).

We could compute the LTE.

$$e_j = y(t_j) - \left[y(t_{j-1}) + h \cdot f\left(t_{j-1} + \frac{h}{2}, y(t_{j-1}) + \frac{h}{2} f(t_{j-1}, y(t_{j-1}))\right) \right]$$

= ... easy, but tedious ...

$$= \mathcal{O}(h^3)$$

The scheme we just derived is known as Runge's (~ 1895) or modified Euler's or explicit midpoint method (EM).

We write it as

$$k_1 = f(t_j, y_j)$$

$$k_2 = f(t_j + h/2, y_j + h/2 \cdot k_1)$$

$$y_{j+1} = y_j + h \cdot k_2$$

← slope approx.!

$\tilde{y}_{j+1/2}$

16.10.24

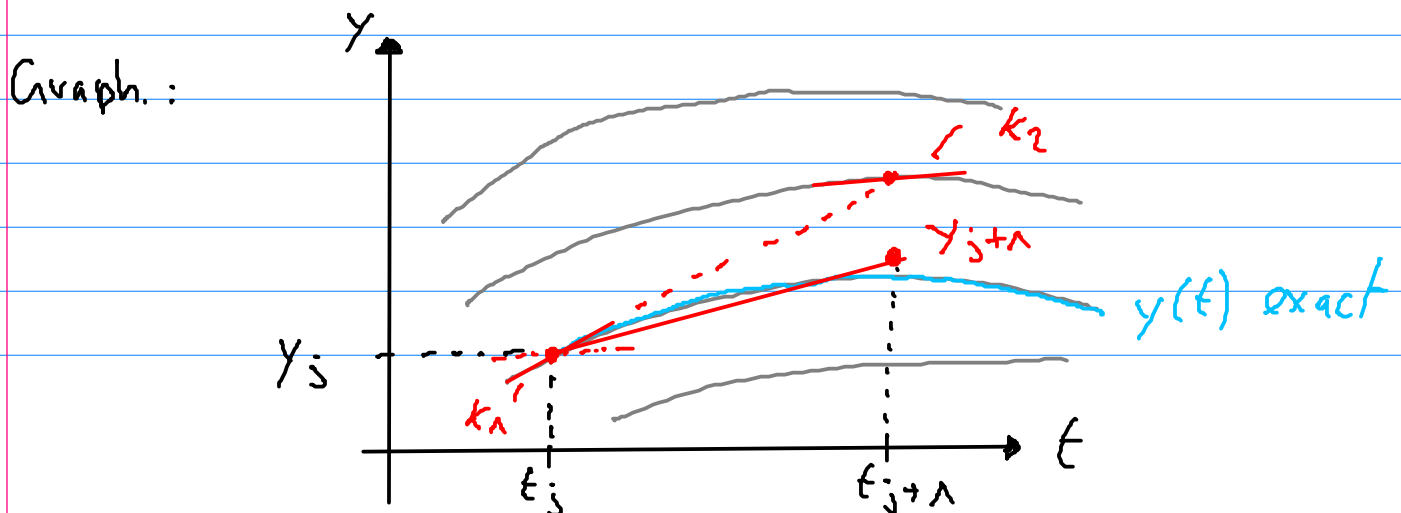
Another possibility is Heun's method

$$k_1 = f(t_j, y_j)$$

$$k_2 = f(t_j + h, y_j + h \cdot k_1)$$

$$y_{j+1} = y_j + \frac{h}{2} (k_1 + k_2)$$

or the explicit trapezoidal method (ET)



One can (easily) show, that Heun's method has order of accuracy $p=2$.

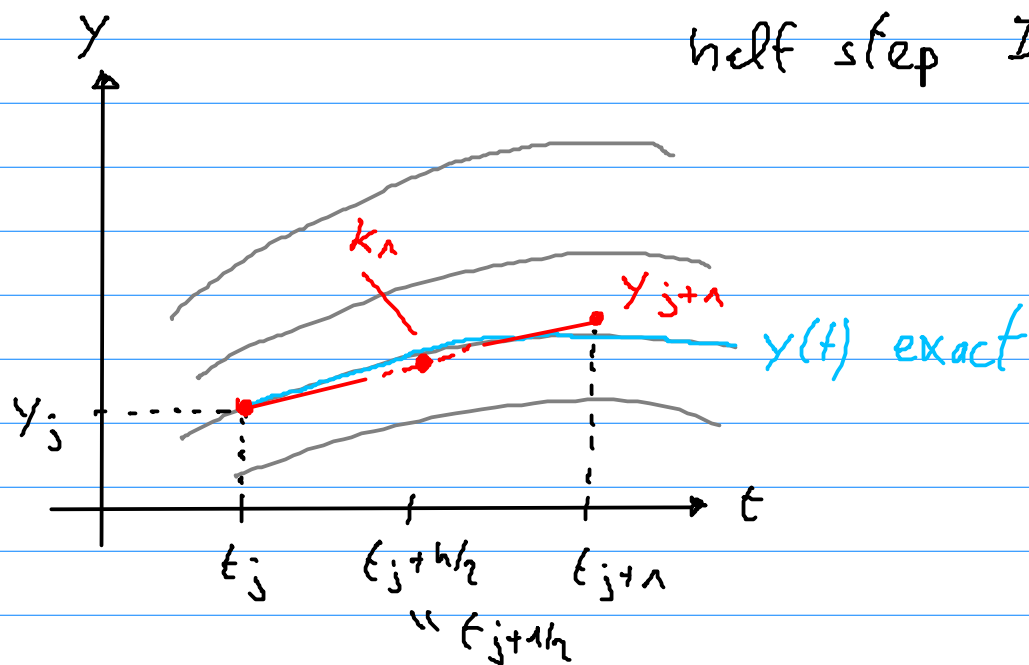
So far we have used the EE to estimate the slopes...

Let's try with IE:

$$y_{j+1} = y_j + h \cdot f(t_j + h/2, \tilde{y}_{j+1/2})$$

with $\tilde{y}_{j+1/2} = y_j + \frac{h}{2} \cdot f(t_j + h/2, \tilde{y}_{j+1/2})$

half step IE!



This method is known as the implicit midpoint method (IMM):

$$\underline{k_n} = f(t_j + h/2, y_j + \frac{h}{2} \underline{k_n})$$

$$y_{j+1} = y_j + h \cdot k_n$$

Another popular scheme is the implicit trapezoidal method (IT):

$$k_1 = f(t_j, y_j)$$

$$\underline{k_2} = f(t_{j+1}, y_j + \frac{h}{2}(k_1 + \underline{k_2}))$$

$$y_{j+1} = y_j + \frac{h}{2}(k_1 + k_2)$$

Sometimes written as

$$y_{j+1} = y_j + \frac{h}{2} \left(f(t_j, y_j) + f(t_{j+1}, y_{j+1}) \right)$$

Clearly looks like the trapezoidal rule and hence the name!

Both IM and IT have order of accuracy $p=2$.

The methods we have seen so far are all representatives of a large class of one-step methods known as Runge-Kutta (RK) methods:

$$k_i = f(t_j + c_i \cdot h, y_j + h \cdot \sum_{l=1}^s a_{il} \cdot k_l)$$

$$y_{j+1} = y_j + h \cdot \sum_{i=1}^s b_i \cdot k_i$$

where:

- s ... number of stages
- c_i ... nodes
- b_i ... weights
- a_{il} ... RK matrix
- k_i ... slopes

It's convenient to write RK methods with a so-called Butcher tableau (BT):

$$\begin{array}{c|cccc}
 c_1 & a_{11} & a_{12} & \dots & a_{1s} \\
 c_2 & a_{21} & a_{22} & \dots & a_{2s} \\
 \vdots & \vdots & \vdots & & \vdots \\
 c_s & a_{s1} & a_{s2} & \dots & a_{ss} \\
 \hline
 & b_1 & b_2 & \dots & b_s
 \end{array}
 = \begin{array}{c|c}
 \vec{c} & A \\
 \hline
 & \vec{b}^T
 \end{array}$$

Ex.: (11) BTs for methods seen so far
 \rightsquigarrow Slides (RK Methods)

(12) Comparison of RK methods on IVP
 of Ex. (6) \rightsquigarrow Slides (RK Methods)

Rem.: (i) For explicit methods the RK matrix
 A is lower triangular

$$A = \begin{pmatrix} 0 & & & & 0 \\ a_{21} & & & & \\ \vdots & & & & \\ a_{s1} & \dots & & & 0 \end{pmatrix}$$

(ii) For system: replace y, f by \vec{y}, \vec{f}

There is indeed a tight connection between RK methods and quadrature. To see this, let's have a look at the LTE:

$$\begin{aligned} e_{j+1} &\stackrel{\text{Def.}}{=} y(t_{j+1}) - \left(y(t_j) + h \cdot \sum_{i=1}^s b_i \cdot k_i \right) \\ &= \underbrace{y(t_{j+1}) - y(t_j)}_{\substack{\downarrow \\ \int_{t_j}^{t_{j+1}} \dot{y}(t) dt}} - h \cdot \sum_{i=1}^s b_i \cdot k_i \\ &= \int_{t_j}^{t_{j+1}} \dot{y}(t) dt - h \cdot \sum_{i=1}^s b_i \cdot k_i \\ &= \int_{t_j}^{t_{j+1}} \underbrace{f(t, y(t))}_{\substack{\downarrow \text{ODE} \\ \uparrow \text{QW} \quad \uparrow f(\dots)}} dt - h \cdot \sum_{i=1}^s b_i \cdot k_i \end{aligned}$$

III.5 Absolute stability

Let's consider the very simple IVP

$$\begin{aligned} \dot{y}(t) &= \lambda y(t) & (\lambda \in \mathbb{R} \text{ or } \mathbb{C}) \\ y(0) &= y_0 \end{aligned}$$

which is known as the Dahlquist test equation (DTE) (in the present context).

The analytical solution is simply

$$y(t) = y_0 e^{\lambda t}$$

Let's apply EE to this IVP

$$y_{j+1} = y_j + h \cdot f(t_j, y_j)$$

$$= y_j + h\lambda y_j$$

$$= (\lambda + h\lambda) y_j$$

$$= (\lambda + h\lambda)^2 y_{j-1}$$

$$\vdots$$

$$= (\lambda + h\lambda)^{j+1} y_0$$

$$y_j = (\lambda + h\lambda) y_{j-1}$$

$$y_{j-1} = (\lambda + h\lambda) y_{j-2}$$

Now let's qualitatively compare the exact to the approx. solution:

(i) $\lambda > 0$: $y(t)$ increases

y_j ? increases \checkmark .

(ii) $\lambda < 0$: $y(t)$ decreases

y_j ? depends on the step size h : $\alpha h < \frac{2}{|\lambda|}$

Let's apply IE to the same IVP

$$y_{j+1} = y_j + h \cdot f(t_{j+1}, y_{j+1})$$

$$= y_j + h\lambda \cdot y_{j+1}$$

(solve for y_{j+1})
IMPLICIT!

$$y_{j+1} = \frac{\lambda}{\lambda - h\lambda} y_j$$

$$y_j = \frac{\lambda}{\lambda - h\lambda} y_{j-1}$$

$$= \left(\frac{\lambda}{\lambda - h\lambda} \right)^2 y_{j-2}$$

\vdots

$$= \left(\frac{\lambda}{\lambda - h\lambda} \right)^{j+1} y_0$$

How does it look here for $\lambda < 0$?

$y(t)$ decreases

y_j ? decreases (does not depend on the step size h !)

30.10.24

Def.: A one-step method applied to the DTE can be written as

$$y_{j+1} = g(z) y_j$$

where $z = h\lambda$ and $g(z)$ is called the stability function (SF) of the method

Ex.: (13) EE : $g(z) = 1 + z$

IE : $g(z) = 1 / (1 - z)$

(14) EM : $g(z) = 1 + z + \frac{z^2}{2}$

We want at least that the approx. solution follows the trend of the exact one!

decay!

For the practically meaningful case $\lambda < 0$ we demand the approx. solution decreases in magnitude:

$$|y_{j+1}| < |y_j| \quad (\text{Absolute stability})$$

or simply

$$\underline{A\text{-stability (AS)}}$$

So with y_{j+1} (SF)

$$|y_{j+1}| = |g(z) y_j| = |g(z)| |y_j| < |y_j|$$

we get a requirement on the SF

$$|g(z)| < 1$$

This motivates the following definition

Def.: Given a one-step method with SF $g(z)$
 For $\lambda \in \mathbb{R}$ we call

$$SI = \{ x = h\lambda \in \mathbb{R} \mid |g(x)| < 1 \}$$

the method's stability interval (SI)

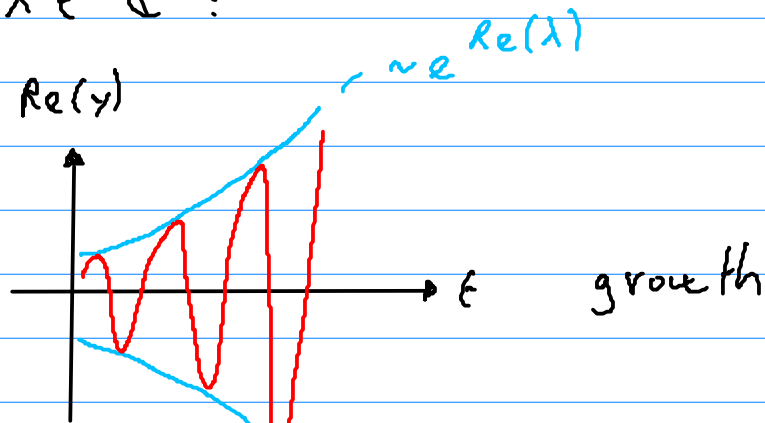
For $\lambda \in \mathbb{C}$ we call

$$SR = \{ z = h\lambda \in \mathbb{C} \mid |g(z)| < 1 \}$$

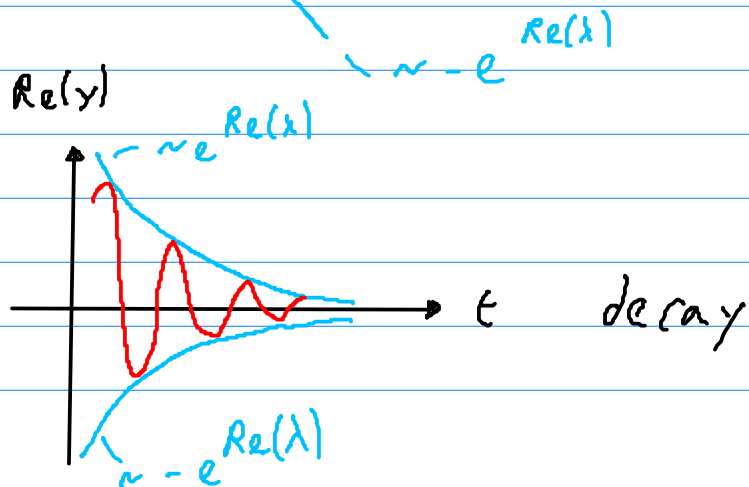
the method's stability region (SR)
 (or "domain")

In general $\lambda \in \mathbb{C}$:

(i) $\text{Re}(\lambda) > 0$



(ii) $\text{Re} < 0$



Def.: A method is called A-stable if the whole left complex half plane is contained in the SR:

$$\{z \in \mathbb{C} \mid \operatorname{Re}(z) < 0\} \subset SR$$

Hence, for A-stable methods there is no time step restriction for $\lambda < 0$!

Ex.: (15) A-stability of EE, IE, EM, ET, IM and RK4 methods

→ slides (Stability regions)

Rem.: In general, explicit RK methods are not A-stable

III.6 Stiff equations

Basically: stiff problems are problems for which "explicit methods don't work"

The step size has to be chosen much smaller for stability rather than from accuracy considerations

Ex.: (16) Stiff linear system

→ Slides (Stiff linear IVP)

A linear inhomogeneous system of ODEs

$$\dot{\vec{y}}(t) = A\vec{y}(t) + \vec{b}(t)$$

↙
n×n Matrix

is called stiff if the eigenvalues of A , λ_i ($i=1,2,\dots,n$), have very different negative real parts:

$$S = \frac{\max_j |\operatorname{Re}(\lambda_j)|}{\min_j |\operatorname{Re}(\lambda_j)|}, \quad \operatorname{Re}(\lambda_j) < 0$$

That is: the stiffness ratio S is "large".

In ex. (16) $S = \frac{1000}{1/2} = 2000$

Stiffness is also very common in nonlinear systems of ODEs

$$\dot{\vec{y}}(t) = \vec{F}(t, \vec{y}(t))$$

A local measure of stiffness is obtained by linearizing the system of ODEs at some point (of interest!) t_n, \vec{y}_n :

$$\begin{aligned} \vec{F}(t, \vec{y}) \approx & \vec{F}(t_n, \vec{y}_n) + \frac{\partial \vec{F}}{\partial t}(t_n, \vec{y}_n) \cdot (t - t_n) \\ & + \underbrace{\frac{\partial \vec{F}}{\partial \vec{y}}(t_n, \vec{y}_n)}_{\text{Jacobian matrix } J = \frac{\partial \vec{F}}{\partial \vec{y}}} \cdot (\vec{y} - \vec{y}_n) \end{aligned}$$

So one obtains the following inhomogeneous linear system of ODEs

$$\dot{\vec{y}}(t) = \underbrace{J(t_n, \vec{y}_n)}_A \vec{y}(t) + \underbrace{\left(\vec{F}(t_n, \vec{y}_n) + \frac{\partial \vec{F}}{\partial t}(t_n, \vec{y}_n) \cdot (t - t_n) - J(t_n, \vec{y}_n) \vec{y}_n \right)}_{\vec{b}(t)}$$

If this is stiff, then one says that the nonlinear system of ODEs is locally stiff around (t_n, \vec{y}_n) .

Ex.: (17) Stiff nonlinear system

→ slides (Stiff nonlinear IVP)

We conclude from Ex. (16) and (17) that explicit methods are inefficient for the numerical treatment of stiff problems:

because the step size is dictated by stability and NOT accuracy considerations

explicit		implicit
cheap per step		expensive per step
step size limited by fastest decaying component		step size only limited by desired accuracy

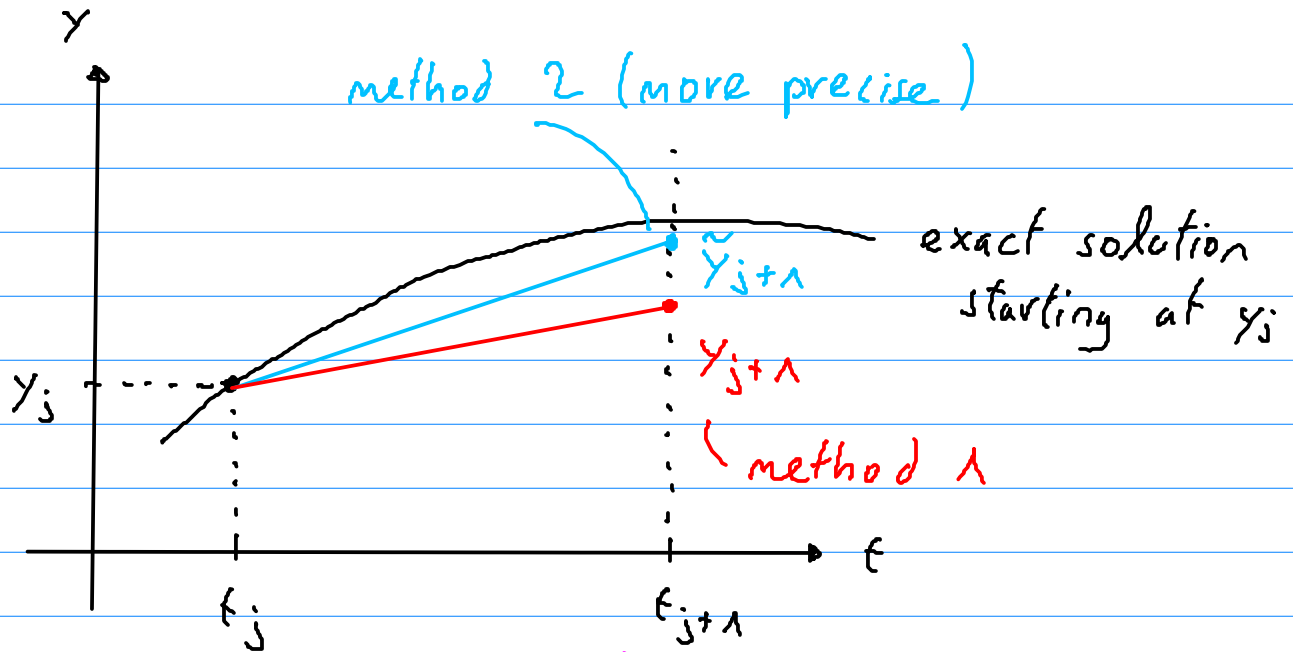
III.7 Error control and adaptive step size

In practice, one wants to choose the step size (the h !) sufficiently small to ensure some required precision and sufficiently large to avoid unnecessary computational work.

To achieve this, we need some local error estimate to indicate if the step size should be "small" or "large".

Idea: perform a step with two methods, where a first method is compared to the result of a second, more precise, method

↳ the digits that match are assumed to be correct !



Now: IF: $|y_{j+1} - \tilde{y}_{j+1}| \leq \text{tolerance}$

accept step

↳ AbsTol or
RelTol in
MATLAB

Else:

repeat step with smaller step size
and check again

Rem.: (i) the more precise method (method 2) could be of higher order than method 1 or method 2 could be the same as method 1 but with a smaller step size h (substeps!)

(ii) the above pseudo-code is far from complete as one also wishes to increase the step size if possible

(iii) no guarantee that the error estimate reflects the true error.
However, often works well in practice.

Ex.: (18) Adaptive step size method (MATLAB ode45) for the van der Pol ODE:

$$\dot{y}_1 = y_2$$

$$\dot{y}_2 = \mu(1 - y_1^2)y_2 - y_1$$

with $\mu = 1$ and $\vec{y} = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$, $t \in [0, 20]$

no slides

(locally)
makes the problem ^vstiff!

(19) " with $\mu = 1000$ and $t \in [0, 3000]$

no slides & exercises