

1 ► Order criteria

Prof. D. Kressner
M. Steinlechner

Consider an explicit two-stage Runge-Kutta method given by the Butcher table

$$\begin{array}{c|cc} 0 & 0 & 0 \\ c & a & 0 \\ \hline & b_1 & b_2 \end{array}$$

Write down equations for the parameters such that the method is consistent of order two. Show that the order of consistency cannot be higher.

This Runge-Kutta method reads

$$y_1(t_0+h) = y_0 + hb_1 f(t_0, y_0) + hb_2 f(t_0+hc, y_0 + hf(t_0, y_0))$$

let's choose $\boxed{a=c}$, otherwise it will be difficult. (cf. lemma 2.7)

Then $y_0 + hf(t_0, y_0) = y_0 + hc y'(t_0) = y(t_0+hc) + O(h^2)$, so that, if f is sufficiently smooth

$$\begin{aligned} hb_2 f(t_0+hc, y_0 + hf(t_0, y_0)) &= hb_2 (f(t_0+hc, y(t_0+hc)) + O(h^2)) \\ &= hb_2 y'(t_0+hc) + O(h^3) \\ &= hb_2 y'(t_0) + h^2 b_2 c y''(t_0) + O(h^3). \end{aligned}$$

Hence

$$y_1(t_0+h) = y_0 + (b_1 + b_2) y'(t_0) h + \frac{1}{2} b_2 c y''(t_0) h^2 + O(h^3).$$

The terms up to order two should equal the ones in the Taylor expansion of y , that is,

$$\boxed{b_2 c = \frac{1}{2}} \quad , \quad \boxed{b_1 + b_2 = 1}$$

For $a=c = \frac{1}{2}$ we get $b_2 = 1, b_1 = 0$, i.e. Runge's method.

Consistency of order three is impossible

Argument is like in Problem 1b): For the IVP

$$y' = y, \quad y(t_0) = y(0)$$

a two-step method will only produce a polynomial of degree two. \square

2 ► Simpler order conditions for simple ODEs

Let $(A, \mathbf{b}, \mathbf{c})$ be an explicit s -stage Runge Kutta method, $s \geq p$, which is invariant under autonomization, i.e., $\mathbf{c} = A\mathbf{e}$.

(a) Consider the simplest autonomous initial value problem

$$\mathbf{y}'(t) = B\mathbf{y}(t), \quad \mathbf{y}(t_0) = \mathbf{y}_0 \in \mathbb{R}^d, \quad B \in \mathbb{R}^{d \times d}.$$

Show that the Runge-Kutta method has consistency order p if and only if

$$\mathbf{b}^T A^{(\beta)} = \frac{1}{(\#\beta)!} \quad \text{for all linear trees } \beta = [[\dots[\odot]\dots]] \text{ with } \#\beta \leq p.$$

a) $y' = By, \quad y(t_0) = y_0$

In this case $f(y) = By$, so $f^{(k)}(y) = 0$ if $k \geq 2$.

Therefore $f^{(\beta)} = 0$ if β contain any branching.

Comparing the Taylor expansion of lemmas 2.10 and 2.11 it follows that we need

$$\mathbf{b}^T A^{(\beta)} = \frac{1}{\beta!} \quad \text{for all linear trees } [\dots[\odot]\dots] \text{ with } \#\beta \leq p.$$

For such trees $\beta! = (\#\beta)!$. □

(b) Consider the quadrature problem

$$\mathbf{y}'(t) = f(t), \quad \mathbf{y}(t_0) = \mathbf{y}_0 \in \mathbb{R}^d$$

of a function $f \in C^\infty([0, T])$. Show that the Runge-Kutta method has consistency order p if and only if

$$\mathbf{b}^T A^{(\beta)} = \frac{1}{\#\beta} \quad \text{for all one-level trees } \beta = [\odot, \dots, \odot] \text{ with } \#\beta \leq p.$$

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

The numerical solution satisfies (component-wise)

$$\begin{aligned} y_1(t_0+h) &= y_0 + h \sum_{i=1}^s b_i f(t_0 + c_i h) \\ &= y_0 + h \sum_{i=1}^s b_i \left(\sum_{k=1}^p \frac{y^{(k)}(t_0)}{(k-1)!} (c_i h)^{k-1} + \mathcal{O}(|h|^p) \right) \\ &= y_0 + \sum_{k=1}^p \frac{y^{(k)}(t_0)}{(k-1)!} \left(\sum_{i=1}^s b_i c_i^{k-1} \right) h^k + \mathcal{O}(|h|^{p+1}). \end{aligned}$$

This leads to the condition

$$\sum_{i=1}^s b_i c_i^{k-1} = \frac{1}{k} \quad \text{for } k=1, \dots, p.$$

But for one-level trees $\beta = [\odot, \dots, \odot]$ with $\#\beta = k \leq p$ we

have $A_i^{(\beta)} = (A\mathbf{e})_i^{k-1} = \left(\sum_j a_{ij} \right)^{k-1} = c_i^{k-1}$. Therefore we could write the condition as

$$\mathbf{b}^T A^{(\beta)} = \frac{1}{\#\beta} \quad \text{for all } \beta = [\odot, \dots, \odot] \text{ with } \#\beta \leq p.$$

□

3 ► Local vs. global error

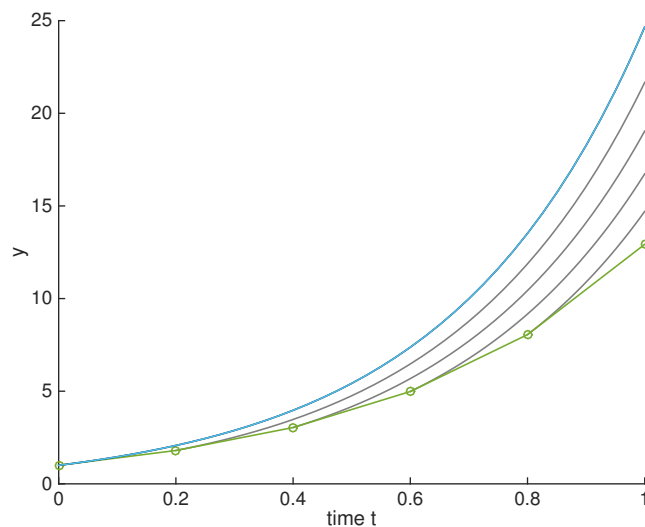
In this exercise, we will investigate the difference between the local and the global error. We consider the ODE

$$y'(t) = 1 - t + 3y, \quad y(t_0) = y_0$$

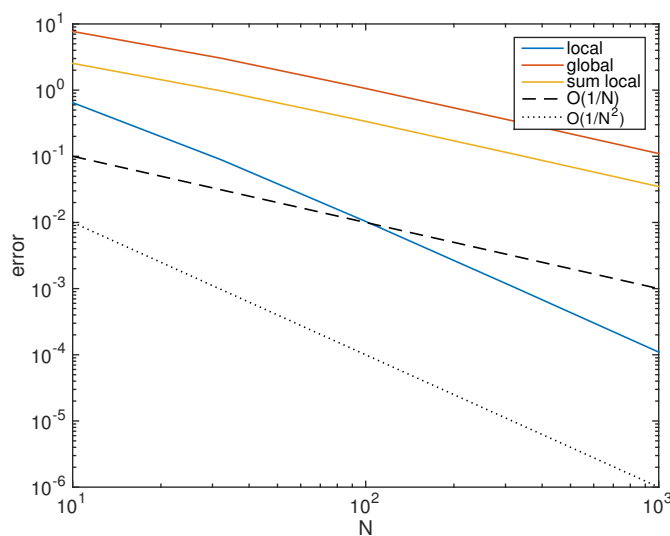
with exact solution

$$y(t) = Ce^{3t} - \frac{2}{9} + \frac{t}{3}, \quad C(t_0, y_0) = \frac{y_0 - \frac{t_0}{3} + \frac{2}{9}}{e^{3t_0}}.$$

1. Create a plot similar to Figure 2.1 (*Lady Windemere's fan*) using an explicit Euler method on the interval $T = [0, 1]$ with stepsize $h = 0.2$ and initial condition $y(0) = 1$. To do this, you have to calculate the exact solutions $y(t)$ with initial conditions $y(t_i) = y_i$ obtained from the i th step of the explicit Euler scheme.



2. Plot the global error $|y(1) - y_N|$, where y_N is the last step of the Euler scheme and compare it to the maximum local error, both as functions of the stepsize h . What is the convergence order that you obtain for the local and the global error?



```

1  function windemere()
3  f = @(t,y) 1 - t + 3*y;
5  sol = @(t) 1/9 * (3*t + 11*exp(3*t) - 2);
   T = [0,1];
7  y0 = 1;
   [t,y] = ode45(f, T, y0);
9  [t5,y5] = expeuler( f, T, y0, 5);

11 hold on
   for i = 2:5
13     [tt,yy] = ode45(f, [t5(i),1], y5(i));
        plot(tt,yy,'-', 'Color',[0.5,0.5,0.5])
15 end

17 set(0,'defaultlinelength',1.2)
   set(gca,'fontsize',14)
19 plot(t5,y5,'-o')

21 xlabel('time_t')
   ylabel('y')
23 plot(t,y,'-k')
   plot(t,sol(t))
25
27 ylim([0,25])

   solC = @(t,C) C*exp(3*t) - 2/9 + t/3;
29 N = round(logspace(1,3,5));
   for i=1:5
31     [t, y] = expeuler( f, T, y0, N(i) );
        y_local = zeros(length(t)-2,1);
33
        for j=2:length(t)-1
35             C = (y(j) - t(j)/3 + 2/9 )/exp(3*t(j));
                 y_local(j-1) = solC(t(j+1), C);
37         end

39         yex = sol(t);
           err_local(i) = max(abs(y(3:end) - y_local));
41         err_global(i) = abs(yex(end) - y(end));
           err_sum_local(i) = sum(abs(y(3:end) - y_local));
43 end

45 figure
   loglog( N, err_local)
47 set(gca,'fontsize',14)
   hold on
49 loglog( N, err_global)
   loglog( N, err_sum_local)
51
   loglog( N, 1./N, '--k')
53 loglog( N, 1./(N.^2), ':k')
   xlabel('N')
55 ylabel('error')
   legend('local','global','sum_local', '0(1/N)', '0(1/N^2)')
57 hold off
   end
59

61 function [t, y] = expeuler( f, T, y0, N );
   t = linspace( T(1), T(2), N+1 )';
63 h = t(2) - t(1);
   y = zeros( N+1, 1 );
65 y(1) = y0;
   for i = 1:N
67     y(i+1) = y(i) + h * f(t(i), y(i));
   end
69 end

```

4 ▶ Order isn't everything¹

Solve the initial value problem

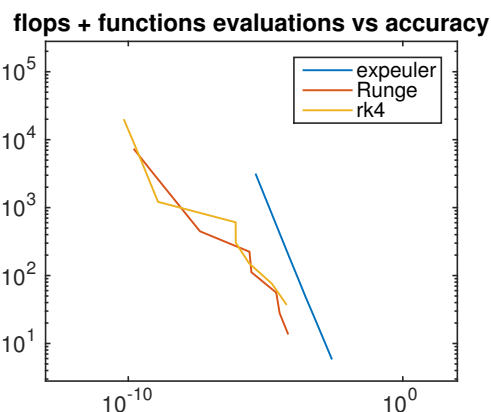
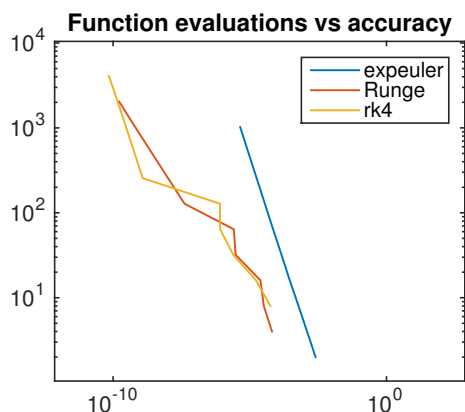
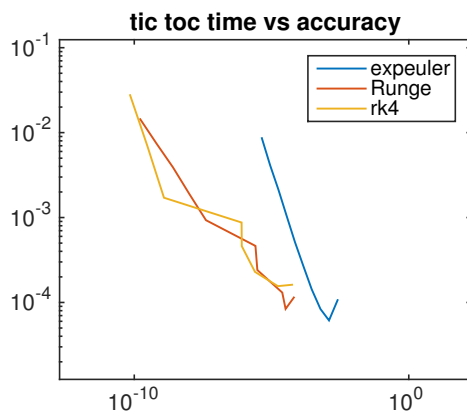
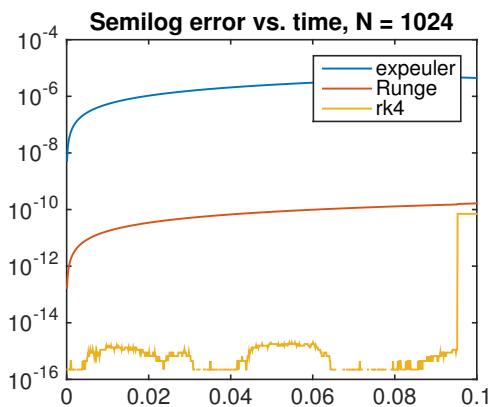
$$y'(t) = |1.1 - y| + 1, \quad t \in [0, 0.1], \quad y(0) = 1,$$

both analytically, and numerically using the explicit Euler method, Runge's method, and the classical RK4 method. Plot errors of all three methods against time.

Then write a script which does the same for various step sizes $N = 2^k$, $k = 1, 2, \dots, 10$, and measures the runtime of each method (use the commands `tic` and `toc`). For all three methods, plot the runtime versus the achieved final accuracy at $T = 0.1$ (use a double logarithmic scale). You may compare to a plot of the required number of function evaluations versus the final achieved accuracy.

Hint. The exact solution is

$$y(t) = \begin{cases} -1.1e^{-t} + 2.1, & 0 \leq t \leq \ln 1.1, \\ \frac{10}{11}e^t + 0.1, & \ln 1.1 \leq t \leq 0.1. \end{cases}$$



¹Title and exercise taken from Deuffhard, Bornemann: Scientific Computing with Ordinary Differential Equations, Springer, 2002.

```

1  % Exercise 4 Problem 3
2  % =====
3
4  function ex4problem3
5  close all;
6
7  % set parameters
8  f = @(t,y) abs(1.1 - y) + 1;
9  T = 0.1;
10  tspan = [0,T];
11  y0 = 1;
12  N = 1024;
13
14  % exact solution
15  y = @(t) (2.1-1.1*exp(-t)).*(t<log(1.1)) + (0.1+10/11*exp(t)).*(t>=log(1.1));
16
17  % run expeuler, runge, and RK4
18  [t1,y1] = expeuler(f,tspan,y0,N);
19  [t2,y2] = runge(f,tspan,y0,N);
20  [t3,y3] = rk4(f,tspan,y0,N);
21
22  % errors
23  e1 = abs(y(t1) - y1);
24  e2 = abs(y(t2) - y2);
25  e3 = abs(y(t3) - y3);
26
27  % plot all in one figure
28  scrsz = get(0,'ScreenSize');
29  figure('Position',[1 scrsz(4)/2-100 scrsz(3)/2 scrsz(4)/2])
30
31  subplot(2,2,1)
32  semilogy(t1,e1,t2,e2,t3,e3);
33  legend('expeuler','Runge','rk4')
34  title(sprintf('Semilog_error_vs_time,_N_=%i',N))
35
36  % performance test
37  for i=1:1:10
38      N = 2^i;
39
40      % exp euler
41      tic;
42      [t1,y1] = expeuler(f,tspan,y0,N);
43      runtime1(i) = toc;
44      funeval1(i) = N;
45      flops1(i) = 3*N;
46      final_error1(i) = abs(y(T) - y1(N + 1));
47
48      % Runge
49      tic;
50      [t2,y2] = runge(f,tspan,y0,N);
51      runtime2(i) = toc;
52      funeval2(i) = 2*N;
53      flops2(i) = 7*N;
54      final_error2(i) = abs(y(T) - y2(N + 1));
55
56      % RK4
57      tic;
58      [t3,y3] = rk4(f,tspan,y0,N);
59      runtime3(i) = toc;
60      funeval3(i) = 4*N;
61      flops3(i) = 19*N;
62      final_error3(i) = abs(y(T) - y3(N + 1));
63  end
64
65  subplot(2,2,2)
66  loglog(final_error1,runtime1,final_error2,runtime2,final_error3,runtime3);
67  legend('expeuler','Runge','rk4')
68  title('tic_toc_time_vs_accuracy')
69

```

```

subplot(2,2,3)
71 loglog(final_error1,funeval1,final_error2,funeval2,final_error3,funeval3);
    legend('expeuler','Runge','rk4')
73 title('Function_evaluations_vs_accuracy')

75 subplot(2,2,4)
loglog(final_error1,flops1,final_error2,flops2,final_error3,flops3);
77 legend('expeuler','Runge','rk4')
    title('flops_+_functions_evaluations_vs_accuracy')
79
end
81
% Required functions
83 % -----
85 % Explicit Euler
87 function [t,y] = expeuler(f,tspan,y0,N)
    y(:,1) = y0;
89     h = (tspan(2)-tspan(1))/N;
        t = tspan(1)*ones(1,N+1) + h*(0:N);
91     for i =2:N+1
        y(:,i) = y(:,i-1) + h*f(t(i-1),y(:,i-1));
93     end
end
95
% Runge's method
97
function [t,y] = runge(f,tspan,y0,N)
99     y(:,1) = y0;
        h = (tspan(2)-tspan(1))/N;
101    t = tspan(1)*ones(1,N+1) + h*(0:N);
        for i =2:N+1
103        k1 = f( t(i-1), y(:,i-1) );
            k2 = f(t(i-1) + h/2, y(:,i-1) + h/2*k1);
105
            y(:,i) = y(:,i-1) + h*k2;
107        end
end
109
% Classical RK4
111
function [t,y] = rk4(f,tspan,y0,N)
113
h = (tspan(2)-tspan(1))/N;
115 t = tspan(1)*ones(1,N+1) + h*(0:N);

117 t(1) = tspan(1);
y(:,1) = y0;
119
g=[1/6 1/3 1/3 1/6];
121
for i=2:N+1
123     k1 = f( t(i-1), y(:,i-1) );
        k2 = f( t(i-1)+h/2, y(:,i-1)+(h/2)*k1 );
125     k3 = f( t(i-1)+h/2, y(:,i-1)+(h/2)*k2 );
        k4 = f( t(i-1)+h, y(:,i-1)+h*k3 );
127
        y(:,i) = y(:,i-1) + h*(g(1)*k1+g(2)*k2+g(3)*k3+g(4)*k4);
129 end
131 end

```