

**ANSWERS OF THE TEST NUMERICAL METHODS FOR  
DIFFERENTIAL EQUATIONS ( WI3097 TU AESB2210 CTB2400 )**

**Thursday July 2 2015, 18:30-21:30**

1. (a) The local truncation error is given by

$$\tau_{n+1}(\Delta t) = \frac{y_{n+1} - z_{n+1}}{\Delta t} \quad (1)$$

where  $z_{n+1}$  is the result of applying the method once with starting solution  $y_n$ . Here we obtain  $y_{n+1}$  by a Taylor expansion around  $t_n$ :

$$y_{n+1} = y_n + \Delta t y'(t_n) + \frac{(\Delta t)^2}{2} y''(t_n) + O((\Delta t)^3). \quad (2)$$

For  $z_{n+1}$ , we obtain, after substitution of the predictor step for  $z_{n+1}^*$  into the corrector step

$$z_{n+1} = y_n + \Delta t ((1 - \mu)f(t_n, y_n) + \mu f(t_n + \Delta t, y_n + \Delta t f(t_n, y_n))) \quad (3)$$

After a Taylor expansion of  $f(t_n + \Delta t, y_n + \Delta t f(t_n, y_n))$  around  $(t_n, y_n)$  one obtains:

$$z_{n+1} = y_n + \Delta t \left( (1 - \mu)f(t_n, y_n) + \mu \left( f(t_n, y_n) + \Delta t \left( \frac{\partial f(t_n, y_n)}{\partial t} + f(t_n, y_n) \frac{\partial f(t_n, y_n)}{\partial y} \right) \right) + O((\Delta t)^2) \right). \quad (4)$$

From the differential equation we know that:

$$y'(t_n) = f(t_n, y_n) \quad (5)$$

From the Chain Rule of Differentiation, we derive

$$y''(t_n) = \frac{df(t_n, y_n)}{dt} = \frac{\partial f(t_n, y_n)}{\partial t} + \frac{\partial f(t_n, y_n)}{\partial y} y'(t_n) \quad (6)$$

after substitution of the differential equation one obtains:

$$y''(t_n) = \frac{\partial f(t_n, y_n)}{\partial t} + \frac{\partial f(t_n, y_n)}{\partial y} f(t_n, y_n) \quad (7)$$

This implies that  $z_{n+1} = y_n + \Delta t y'(t_n) + \mu(\Delta t)^2 y''(t_n) + O((\Delta t)^3)$ . Subsequently, it follows that

$$y_{n+1} - z_{n+1} = O((\Delta t)^2), \text{ and, hence } \tau_{n+1}(\Delta t) = \frac{O((\Delta t)^2)}{\Delta t} = O(\Delta t) \text{ for } 0 \leq \mu \leq 1, \quad (8)$$

$$y_{n+1} - z_{n+1} = O((\Delta t)^3), \text{ and, hence } \tau_{n+1}(\Delta t) = \frac{O((\Delta t)^3)}{\Delta t} = O((\Delta t)^2) \text{ for } \mu = \frac{1}{2}. \quad (9)$$

(b) Consider the test equation  $y' = \lambda y$ , then, herewith, one obtains

$$\begin{aligned} w_{n+1}^* &= w_n + \lambda \Delta t w_n = (1 + \lambda \Delta t) w_n, \\ w_{n+1} &= w_n + ((1 - \mu) \lambda \Delta t w_n + \mu \lambda \Delta t w_{n+1}^*) = \\ &= w_n + ((1 - \mu) \lambda \Delta t w_n + \mu \lambda \Delta t (w_n + \lambda \Delta t w_n)) = (1 + \lambda \Delta t + \mu (\lambda \Delta t)^2) w_n. \end{aligned} \quad (10)$$

Hence the amplification factor is given by

$$Q(\lambda \Delta t) = 1 + \lambda \Delta t + \mu (\lambda \Delta t)^2. \quad (11)$$

(c) Doing one step with the given method with  $\Delta t = \frac{1}{2}$  and  $\mu = \frac{1}{2}$  leads to the following steps:

Predictor:

$$\begin{pmatrix} w_1 \\ w_2 \end{pmatrix}^* = \begin{pmatrix} 0 \\ 1 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} -0 + \cos(0) + 2 + 0 \\ 0 - 1 \end{pmatrix} = \begin{pmatrix} \frac{3}{2} \\ \frac{1}{2} \end{pmatrix}$$

Corrector:

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \end{pmatrix} + \frac{1}{2} \left( \frac{1}{2} \begin{pmatrix} 3 \\ -1 \end{pmatrix} + \frac{1}{2} \begin{pmatrix} -\frac{3}{2} + \cos(\frac{3}{2}) + 2 \cdot \frac{1}{2} + \frac{1}{2} \\ \frac{3}{2} - (\frac{1}{2})^2 \end{pmatrix} \right)$$

which can be written as:

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 + \frac{3}{4} - \frac{3}{8} + \frac{1}{4} \cos(\frac{3}{2}) + \frac{3}{8} \\ 1 - \frac{1}{4} + \frac{3}{8} - \frac{1}{16} \end{pmatrix} = \begin{pmatrix} \frac{3}{4} + \frac{1}{4} \cos(\frac{3}{2}) \\ \frac{17}{16} \end{pmatrix} = \begin{pmatrix} 0.7677 \\ 1.0625 \end{pmatrix}$$

(d) In order to compute the Jacobian, we note that the right-hand side of the non linear system can be noted by:

$$f_1(x_1, x_2) = -x_1 + \cos x_1 + 2x_2 + t$$

$$f_2(x_1, x_2) = x_1 - x_2^2$$

From the definition of the Jacobian it follows that:

$$\begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} \end{pmatrix} = \begin{pmatrix} -1 - \sin x_1 & 2 \\ 1 & -2x_2 \end{pmatrix}.$$

Substitution of  $\begin{pmatrix} x_1(0) \\ x_2(0) \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$  shows that

$$J = \begin{pmatrix} -1 & 2 \\ 1 & -2 \end{pmatrix}.$$

- (e) For the stability it is sufficient to check that  $|Q(\lambda_i \Delta t)| \leq 1$  for all the eigenvalues of the Jacobian matrix. It is easy to see that the eigenvalues of the Jacobian matrix are  $\lambda_1 = -3$  and  $\lambda_2 = 0$ .

For the choice  $\mu = 0$  we note that the method is equal to the Euler Forward method. For real eigenvalues the Euler Forward method is stable if  $\Delta t \leq \frac{-2}{\lambda}$ . Since  $\lambda_1 = -3$  and  $\lambda_2 = 0$  we know that the method is stable if  $\Delta t \leq \frac{-2}{-3} = \frac{2}{3}$  (another option is to derive the values of  $\Delta t$  such that  $|Q(\lambda_i \Delta t)| \leq 1$  by using the description of  $Q(\lambda \Delta t)$ )

For the choice  $\mu = \frac{1}{2}$  we use the expression

$$Q(\lambda \Delta t) = 1 + \lambda \Delta t + \frac{1}{2}(\lambda \Delta t)^2$$

For  $\lambda_2 = 0$  it appears that  $Q(\lambda_2 \Delta t) = 1$  so the inequality is satisfied for all  $\Delta t$ . For  $\lambda_1 = -3$  we have to check the following inequalities:

$$-1 \leq 1 - 3\Delta t + \frac{9}{2}(\Delta t)^2 \leq 1$$

For the left-hand inequality we arrive at

$$0 \leq \frac{9}{2}(\Delta t)^2 - 3\Delta t + 2$$

It appears that the discriminant  $9 - 4 \cdot \frac{9}{2} \cdot 2$  is negative, so there are no real roots which implies that the inequality is satisfied for all  $\Delta t$ .

For the right-hand inequality we get

$$-3\Delta t + \frac{9}{2}(\Delta t)^2 \leq 0$$

$$\frac{9}{2}(\Delta t)^2 \leq 3\Delta t$$

so

$$\Delta t \leq \frac{2}{3}$$

(another option is to see that for  $\mu = \frac{1}{2}$  the method is equal to the modified Euler method, and remember that this method is stable for real eigenvalues if  $\Delta t \leq \frac{-2}{\lambda}$ )

2. (a) First, we check the boundary conditions:

$$u(0) = 0 - \frac{1 - e^0}{1 - e} = \frac{1 - 1}{1 - e} = 0, \quad u(1) = 1 - \frac{1 - e^1}{1 - e} = 0. \quad (12)$$

Further, we have

$$u'(x) = 1 + \frac{e^x}{1-e}, \quad (13)$$

$$u''(x) = \frac{e^x}{1-e}. \quad (14)$$

Hence, we immediately see

$$-u''(x) + u'(x) = -\frac{e^x}{1-e} + 1 + \frac{e^x}{1-e} = 1. \quad (15)$$

Hence, the solution  $u(x) = 1 - \frac{1-e^x}{1-e}$  satisfies the differential and the boundary conditions, and therewith  $u(x)$  is the solution to the boundary value problem (uniqueness can be demonstrated in a straightforward way, but this was not asked for).

- (b) The domain of computation, being  $(0, 1)$ , is divided into subintervals with mesh points, we set  $x_j = j\Delta x$ , where we use  $n$  unknowns, such that  $x_{n+1} = (n+1)\Delta x = 1$ . We are looking for a discretization with an error of second order,  $O((\Delta x)^2)$ . To this extent, we use the following central differences approximation at  $x_j$ :

$$u'(x_j) \approx \frac{u(x_{j+1}) - u(x_{j-1}))}{2\Delta x}, \text{ for } j \in \{1, \dots, n\}. \quad (16)$$

We note that the above formula can be derived formally by writing the derivative as

$$u'(x_j) = \frac{\alpha_0 u(x_{j-1}) + \alpha_1 u(x_j) + \alpha_2 u(x_{j+1}))}{\Delta x}, \quad (17)$$

and solve  $\alpha_0$ ,  $\alpha_1$  and  $\alpha_2$  from checking the zeroth, first and second order derivatives of  $u(x)$ . Further, the second order derivative is approximated by

$$u''(x_j) \approx \frac{u(x_{j+1}) - 2u(x_j) + u(x_{j-1}))}{(\Delta x)^2}. \quad (18)$$

Since we approximate the derivatives at the point  $x_j$ , we use Taylor series expansion about  $x_j$ , to obtain:

$$\begin{aligned} u(x_{j+1}) &= u(x_j + \Delta x) = u(x_j) + \Delta x u'(x_j) + \frac{(\Delta x)^2}{2} u''(x_j) + \frac{(\Delta x)^3}{6} u'''(x_j) + O((\Delta x)^4), \\ u(x_{j-1}) &= u(x_j - \Delta x) = u(x_j) - \Delta x u'(x_j) + \frac{(\Delta x)^2}{2} u''(x_j) - \frac{(\Delta x)^3}{6} u'''(x_j) + O((\Delta x)^4), \end{aligned} \quad (19)$$

This gives

$$\begin{aligned} -\frac{u(x_{j+1}) - 2u(x_j) + u(x_{j-1}))}{(\Delta x)^2} + \frac{u(x_{j+1}) - u(x_{j-1}))}{2\Delta x} &= -u''(x_j) + u'(x_j) \\ + \frac{O((\Delta x)^3)}{2\Delta x} + \frac{O((\Delta x)^4)}{(\Delta x)^2} &= -u''(x_j) + u'(x_j) + O((\Delta x)^2). \end{aligned} \quad (20)$$

Hence the error is second order, that is  $O((\Delta x)^2)$ . Next, we neglect the truncation error, and set  $w_j := u(x_j)$  to get

$$-\frac{w_{j+1} - 2w_j + w_{j-1}}{(\Delta x)^2} + \frac{w_{j+1} - w_{j-1}}{2\Delta x} = 1, \text{ for } j \in \{1, \dots, n\}. \quad (21)$$

At the boundaries, we see for  $j = 1$  and  $j = n$ , upon substituting  $w_0 = 0$  and  $w_{n+1} = 0$ , respectively:

$$\begin{aligned} -\frac{w_2 - 2w_1 + 0}{(\Delta x)^2} + \frac{w_2 - 0}{2\Delta x} &= 1, \\ -\frac{0 - 2w_n + w_{n-1}}{(\Delta x)^2} + \frac{0 - w_{n-1}}{2\Delta x} &= 1. \end{aligned} \quad (22)$$

This can be rewritten more neatly as follows:

$$\begin{aligned} \frac{-w_2 + 2w_1}{(\Delta x)^2} + \frac{w_2}{2\Delta x} &= 1, \\ \frac{2w_n - w_{n-1}}{(\Delta x)^2} - \frac{w_{n-1}}{2\Delta x} &= 1. \end{aligned} \quad (23)$$

(c) The real-valued exact solution and its first and second derivative are given by

$$u(x) = x - \frac{1 - e^x}{1 - e}, \quad (24)$$

$$u'(x) = 1 + \frac{e^x}{1 - e}, \quad (25)$$

$$u''(x) = \frac{e^x}{1 - e}. \quad (26)$$

First, we calculate the point  $x^* = \ln(1/(e - 1))$ , where  $u'(x^*) = 0$  and verify that  $u(x)$  attains its maximum value at  $x^*$  (since  $u''(x^*) = -1/(e - 1)^2 < 0$ ). Since  $u(0) = u(1) = 0$  we can conclude that the exact solution is monotonically increasing on  $[0, x^*]$  and monotonically decreasing on  $[x^*, 1]$ . Since the numerical solution should have the same characteristics as the exact solution, oscillatory solutions should be considered as not reflecting the analytic solution.

(d) Next, we use  $\Delta x = 1/4$ , then, from equations (21) and (23), one obtains the following system

$$32w_1 - 14w_2 = 1 \quad (27)$$

$$-18w_1 + 32w_2 - 14w_3 = 1 \quad (28)$$

$$-18w_2 + 32w_3 = 1 \quad (29)$$

- (e) The iteration process is a fixed point method. If the process converges we have:  $\lim_{n \rightarrow \infty} x_n = p$ . Using this in the iteration process yields:

$$\lim_{n \rightarrow \infty} x_{n+1} = \lim_{n \rightarrow \infty} [x_n + h(x_n)(x_n^3 - 27)]$$

Since  $h$  is a continuous function one obtains:

$$p = p + h(p)(p^3 - 27)$$

so

$$h(p)(p^3 - 27) = 0.$$

Since  $h(x) \neq 0$  for each  $x \neq 0$  it follows that  $p^3 - 27 = 0$  and thus  $p = 27^{\frac{1}{3}} = 3$ .

- (f) The convergence of a fixed point method  $x_{n+1} = g(x_n)$  is determined by  $g'(p)$ . If  $|g'(p)| < 1$  the method converges, whereas if  $|g'(p)| > 1$  the method diverges. For all choices we compute the first derivative in  $p$ . For the first method we elaborate all steps. For the other methods we only give the final result. For  $h_1$  we have  $g_1(x) = x - \frac{x^3 - 27}{x^4}$ . The first derivative is:

$$g'_1(x) = 1 - \frac{3x^2 \cdot x^4 - (x^3 - 27) \cdot 4x^3}{(x^4)^2}$$

Substitution of  $p$  yields:

$$g'_1(p) = 1 - \frac{3p^6 - (p^3 - 27) \cdot 4p^3}{p^8}.$$

Since  $p = 3$  the final term cancels:

$$g'_1(p) = 1 - \frac{3p^6}{p^8} = 1 - 3^{-1} = \frac{2}{3}.$$

This implies that the method is convergent with convergence factor  $\frac{2}{3}$ .

For the second method we have:

$$g'_2(p) = 1 - \frac{3p^4 - (p^3 - 27) \cdot 2p}{p^4} = 1 - \frac{3p^4}{p^4} = -2$$

Thus the method diverges.

For the third method we have:

$$g'_3(p) = 1 - \frac{9p^4 - (p^3 - 27) \cdot 6p}{9p^4} = 1 - \frac{9p^4}{9p^4} = 0$$

Thus the method is convergent with convergence factor 0.

Concluding we note that the third method is the fastest.

- (g) To estimate the error in  $p$  we first approximate the function  $f$  in the neighborhood of  $p$  by the first order Taylor polynomial:

$$P_1(x) = f(p) + (x - p)f'(p) = (x - p)f'(p).$$

Due to the measurement errors we know that

$$(x - p)f'(p) - \epsilon_{max} \leq \hat{P}_1(x) \leq (x - p)f'(p) + \epsilon_{max}.$$

This implies that the perturbed root  $\hat{p}$  is bounded by the roots of  $(x - p)f'(p) - \epsilon_{max}$  and  $(x - p)f'(p) + \epsilon_{max}$ , which leads to

$$p - \frac{\epsilon_{max}}{|f'(p)|} \leq \hat{p} \leq p + \frac{\epsilon_{max}}{|f'(p)|}.$$