

# Numerical Solutions for DEs HW3

YANG, Ze (5131209043)

April 1, 2017

## A Note to TA:

*Hi, this is the senior student from Antai College who did not register for this course. I would like to do all the assignments for practice, but feel free to just skip my homework if you don't have time. Thank you again for allowing me to access the assignments and other class material! : )*  
- Ze

**Problem 1.** (Iserles 4.4) Determine all values of  $\theta$  such that the theta method (1.13) is absolutely stable.

*Proof.* The theta method reads

$$\mathbf{y}_{n+1} = \mathbf{y}_n + h[\theta \mathbf{f}(t_n, \mathbf{y}_n) + (1 - \theta) \mathbf{f}(t_{n+1}, \mathbf{y}_{n+1})], \quad n = 0, 1, \dots$$

To find the a-stability region ( $\mathcal{D}$ ) we apply the method to scalar ode,

$$\begin{aligned} y_{n+1} &= y_n + h[\theta \lambda y_n + (1 - \theta) \lambda y_{n+1}] \\ y_{n+1} &= \frac{1 + h\lambda\theta}{1 - h\lambda(1 - \theta)} y_n = \frac{1 + z\theta}{1 - z(1 - \theta)} y_n \end{aligned} \quad (1)$$

Therefore, we want to find  $\theta$ , such that  $\mathbb{C}^- = \{z \in \mathbb{C} : \operatorname{Re}(z) < 0\} \subseteq \{z \in \mathbb{C} : \left| \frac{1+z\theta}{1-z(1-\theta)} \right| < 1\} = \mathcal{D}$  (By the definition of a-stability). I.e. we want  $|1 + z\theta| < |1 - z(1 - \theta)|$  (\*) holds for all  $z$  that  $\operatorname{Re}(z) < 0$ .

Write  $z = x + yi$ :

- $\theta = 1/2$ , the method is a-stable, since it's the trapezoid method.
- $\frac{1}{2} < \theta \leq 1$ , we write  $\theta = \frac{1}{2} + \delta$ . Let  $z = -1/\delta + 0i$ . Then  $LHS = 1/2\delta$ ,  $RHS = |1 + 1/\delta(-1/2 - \delta)| = 1/2\delta$ . Then the equation (\*) does not hold, so in this case the method is not a-stable.
- $0 \leq \theta < \frac{1}{2}$ , we want to show

$$|1 + \theta x + \theta yi| < |1 - (1 - \theta)x - (1 - \theta)yi|$$

for all  $(x, y) \in \mathbb{C}^-$ . Denote the complex number inside absolute value as  $r$  (right hand side) and  $l$  (left hand side).  $\operatorname{Im}(l) = \theta y$ ,  $\operatorname{Im}(r) = (1 - \theta)y$ , since  $0 \leq \theta < \frac{1}{2} \Rightarrow 1 - \theta > \theta > 0$ . So  $\operatorname{Im}(l) > \operatorname{Im}(r) > 0$ , Hence  $|\operatorname{Im}(l)| > |\operatorname{Im}(r)|$ .

For the real parts,  $\operatorname{Re}(l) = 1 + x\theta$ ,  $\operatorname{Re}(r) = 1 - (1 - \theta)x$ .

1.  $1 + x\theta \leq 0$ : since we must have  $x < 0$ .  $(2\theta - 1) < 0 \Rightarrow (2\theta - 1)x > 0$ ,  $2 + (2\theta - 1)x > 0 \Rightarrow 1 - (-\theta + 1)x > -1 - \theta x \geq 0$ , i.e.  $\operatorname{Re}(r) > -\operatorname{Re}(l) \geq 0$ .
2.  $1 + x\theta > 0$ :  $\theta x < -(1 - \theta)x \Rightarrow 0 < 1 + \theta x < 1 - (1 - \theta)x \Rightarrow \operatorname{Re}(r) > \operatorname{Re}(l) \geq 0$ .

In either case we have  $|\operatorname{Re}(l)| > |\operatorname{Re}(r)|$ . So finally, we have  $|r| < |l|$  for all  $z \in \mathbb{C}^-$ . And we conclude that for  $\theta \in [0, 1/2]$ , the method is a-stable.

□

**Problem 2.** (Iserles 4.9) The two-step method

$$\mathbf{y}_{n+2} - \mathbf{y}_n = 2h\mathbf{f}(t_{n+1}, \mathbf{y}_{n+1}), \quad n = 0, 1, \dots$$

is called the explicit midpoint rule.

- Denote by  $w_1(z)$  and  $w_2(z)$  the zeroes of the underlying function  $\eta(z, \cdot)$ , show that  $w_1(z)w_2(z) \equiv -1$  for all  $z \in \mathbb{C}$ .
- Show that  $\mathcal{D} = \emptyset$ .
- We say that  $\tilde{\mathcal{D}}$  is a *weak linear stability domain* of a numerical method if, when applied to the scalar linear equation, it produces a uniformly bounded solution sequence. (It is easy to see that  $\tilde{\mathcal{D}} = \text{cl}\mathcal{D}$  for most methods of interest.) Determine explicitly  $\tilde{\mathcal{D}}$  for the method.

*Proof.* (a.) We have  $\rho(w) = w^2 - 1$ ,  $\sigma(w) = 2w$ . So let  $z = \lambda h$ ,  $\eta(w, z) = \rho(w) - z\sigma(w) = w^2 - 2zw - 1$ . Suppose  $w_1(z)$  and  $w_2(z)$  solves  $\eta(\cdot, z) = 0$ , then by Vieta's formula, we have  $w_1(z)w_2(z) = \frac{-1}{1} = -1$ .  $\square$

(b.) We have  $w_1(z) = \frac{2z + \sqrt{4z^2 + 4}}{2} = z + \sqrt{z^2 + 1}$ .  $w_2(z) = z - \sqrt{z^2 + 1}$ . So  $\mathcal{D} = \{z \in \mathbb{C} : |w_1(z)| < 1, |w_2(z)| < 1\}$ . But

$$|w_1(z)| \cdot |w_2(z)| = |w_1(z)w_2(z)| = |-1| = 1$$

implies that  $|w_1(z)| = 1/|w_2(z)|$ . Hence either (1)  $|w_1| = |w_2| = 1$ , or (2) one of the two roots has modulus strictly smaller than 1, another one strictly greater than 1. And for all  $z \in \mathbb{C}$ , one of the two cases must be true.  $\Rightarrow \forall z \in \mathbb{C}, z \notin \mathcal{D}$ .  $\mathcal{D} = \emptyset$ .

(c) Apply to scalar linear equation  $y' = \lambda y$ ,

$$y_{n+2} - y_n = 2h\lambda y_{n+1}$$

$\square$

**Problem 3.** (Iserles 5.8) Show that the *Hénon-Heiles system* (using the definition from wikipedia)

$$\begin{aligned} q'_1 &= p_1 \\ p'_1 &= -q_1 - 2\lambda q_1 q_2 \\ q'_2 &= p_2 \\ p'_2 &= -q_2 - \lambda(q_1^2 - q_2^2) \end{aligned} \tag{2}$$

is Hamiltonian and identify explicitly the hamiltonian energy.

*Proof.* Let  $\mathbf{q} = (q_1, q_2)$ ,  $\mathbf{p} = (p_1, p_2)$ . Assume the system is Hamiltonian, then

$$\begin{aligned} \nabla_{\mathbf{p}} H &= \dot{\mathbf{q}} = (p_1 \quad p_2) \\ \nabla_{\mathbf{q}} H &= -\dot{\mathbf{p}} = (q_1 + 2\lambda q_1 q_2 \quad q_2 + \lambda(q_1^2 - q_2^2)) \end{aligned} \tag{3}$$

Hence

$$\begin{aligned} H(\mathbf{p}, \mathbf{q}) &= \frac{1}{2}(p_1^2 + p_2^2) + f(\mathbf{q}) \\ H(\mathbf{p}, \mathbf{q}) &= \int (q_1 + 2\lambda q_1 q_2) dq_1 + g(\mathbf{p}, \mathbf{q}) = \frac{1}{2}q_1^2 + \lambda q_2 q_1^2 + g(\mathbf{p}, \mathbf{q}) \\ H(\mathbf{p}, \mathbf{q}) &= \int [q_2 + \lambda(q_1^2 - q_2^2)] dq_2 + h(\mathbf{p}, \mathbf{q}) = \frac{1}{2}q_2^2 + \lambda q_1^2 q_2 - \frac{1}{3}\lambda q_2^3 + h(\mathbf{p}, \mathbf{q}) \end{aligned}$$

Combine the first integrals above we have:

$$H(\mathbf{p}, \mathbf{q}) = \frac{1}{2}(p_1^2 + p_2^2) + \frac{1}{2}(q_1^2 + q_2^2) + \lambda \left( q_1^2 q_2 - \frac{1}{3} q_2^3 \right)$$

Check the partial derivatives of  $H$ , we find that it gives the system indeed. Hence the system is Hamiltonian, and the hamiltonian energy is given by  $H(\mathbf{p}, \mathbf{q})$ .  $\square$

**Problem 4.** The symplectic Euler method for the Hamiltonian system reads

$$\begin{cases} \mathbf{p}_{n+1} = \mathbf{p}_n - h \frac{\partial H(\mathbf{p}_{n+1}, \mathbf{q}_n)}{\partial \mathbf{q}} \\ \mathbf{q}_{n+1} = \mathbf{q}_n + h \frac{\partial H(\mathbf{p}_{n+1}, \mathbf{q}_n)}{\partial \mathbf{p}} \end{cases} \quad (4)$$

- Show that this is a first order method.
- Prove from basic principles that, as implied by its name, the method is indeed symplectic. (Hint: let  $G = \nabla^2 H$ , and write  $G = [G_{11}, G_{12}; G_{21}, G_{22}]$ ).
- Assuming that the Hamiltonian is separable,  $H(\mathbf{p}, \mathbf{q}) = T(\mathbf{p}) + V(\mathbf{q})$ , where  $T$  and  $V$  correspond to kinetic and potential energy respectively. Show that the method can be implemented explicitly.
- Use symplectic Euler and explicit Euler to solve the problem of non-linear pendulum. The Hamiltonian is  $H(p, q) = \frac{1}{2}p^2 - \cos(q)$ , and the Hamiltonian equations are

$$\dot{p} = -\sin q, \quad \dot{q} = p$$

with initial condition  $p(0) = 0, q(0) = 1$ . Plot the error of the numerical methods in the Hamiltonian  $H$ .

*Proof.* (a.) We examine the truncation error  $T_n^{[1,2]}/h$  for  $\mathbf{q}$  and  $\mathbf{p}$  respectively. Denote  $\mathbf{f} = \dot{\mathbf{p}} = -\frac{\partial H}{\partial \mathbf{q}}$ ,  $\mathbf{g} = \dot{\mathbf{q}} = \frac{\partial H}{\partial \mathbf{p}}$ .

$$\begin{aligned} T_n^{[1]}/h &= \frac{\mathbf{p}(t_{n+1}) - \mathbf{p}(t_n)}{h} - \mathbf{f}(\mathbf{p}(t_{n+1}), \mathbf{q}(t_n)) \\ &= \frac{1}{h} [\dot{\mathbf{p}}(t_n)h + \frac{1}{2}\ddot{\mathbf{p}}(t_n)h^2 + O(h^3)] - (\mathbf{f}(t_n) + \mathbf{f}'(t_n)h + O(h^2)) \\ &= -\frac{1}{2}\mathbf{f}'(t_n)h + O(h^2) \end{aligned} \quad (5)$$

The calculation is same for  $T_n^{[2]}/h$ , we obtain  $T_n^{[2]}/h = -\frac{1}{2}\mathbf{g}'(t_n)h + O(h^2)$ . Hence we conclude that the method is of order 1.  $\square$

(b.) Let

$$\mathbf{G} = \nabla^2 H = \begin{pmatrix} \frac{\partial^2 H}{\partial \mathbf{p}^2} & \frac{\partial^2 H}{\partial \mathbf{p} \partial \mathbf{q}} \\ \frac{\partial^2 H}{\partial \mathbf{q} \partial \mathbf{p}} & \frac{\partial^2 H}{\partial \mathbf{q}^2} \end{pmatrix} =: \begin{pmatrix} \mathbf{G}_{pp} & \mathbf{G}_{qp}^\top \\ \mathbf{G}_{qp} & \mathbf{G}_{qq} \end{pmatrix}$$

Then calculate the partial derivatives from the method's formula:

$$\begin{cases} \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{p}_n} = \mathbf{I} - h\mathbf{G}_{qp} \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{p}_n} \\ \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{q}_n} = -h\mathbf{G}_{qq} - h\mathbf{G}_{qp} \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{q}_n} \\ \frac{\partial \mathbf{q}_{n+1}}{\partial \mathbf{p}_n} = h\mathbf{G}_{pp} \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{p}_n} \\ \frac{\partial \mathbf{q}_{n+1}}{\partial \mathbf{q}_n} = \mathbf{I} + h\mathbf{G}_{pq} + h\mathbf{G}_{pp} \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{q}_n} \end{cases} \quad (6)$$

I.e.

$$\begin{cases} (\mathbf{I} + h\mathbf{G}_{qp}) \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{p}_n} = \mathbf{I} \\ (\mathbf{I} + h\mathbf{G}_{qp}) \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{q}_n} = -h\mathbf{G}_{qq} \\ -h\mathbf{G}_{pp} \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{p}_n} + \frac{\partial \mathbf{q}_{n+1}}{\partial \mathbf{p}_n} = \mathbf{O} \\ -h\mathbf{G}_{pp} \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{q}_n} + \frac{\partial \mathbf{q}_{n+1}}{\partial \mathbf{q}_n} = \mathbf{I} + h\mathbf{G}_{pq} \end{cases} \quad (7)$$

rewrite in matrix form:

$$\begin{pmatrix} \mathbf{I} + h\mathbf{G}_{qp} & \mathbf{O} \\ -h\mathbf{G}_{pp} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{p}_n} & \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{q}_n} \\ \frac{\partial \mathbf{q}_{n+1}}{\partial \mathbf{p}_n} & \frac{\partial \mathbf{q}_{n+1}}{\partial \mathbf{q}_n} \end{pmatrix} = \begin{pmatrix} \mathbf{I} & -h\mathbf{G}_{qq} \\ \mathbf{O} & \mathbf{I} + h\mathbf{G}_{pq} \end{pmatrix} \quad (8)$$

And solve this linear system:

$$\Phi = \begin{pmatrix} \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{p}_n} & \frac{\partial \mathbf{p}_{n+1}}{\partial \mathbf{q}_n} \\ \frac{\partial \mathbf{q}_{n+1}}{\partial \mathbf{p}_n} & \frac{\partial \mathbf{q}_{n+1}}{\partial \mathbf{q}_n} \end{pmatrix} = \begin{pmatrix} (\mathbf{I} + h\mathbf{G}_{qp})^{-1} & -h\mathbf{G}_{qq}(\mathbf{I} + h\mathbf{G}_{qp})^{-1} \\ h\mathbf{G}_{pp}(\mathbf{I} + h\mathbf{G}_{qp})^{-1} & \mathbf{I} + h\mathbf{G}_{pq} - h^2\mathbf{G}_{pp}\mathbf{G}_{qq}(\mathbf{I} + h\mathbf{G}_{qp})^{-1} \end{pmatrix} \quad (9)$$

So:

$$\begin{aligned}\Phi^\top J &= \begin{pmatrix} (\mathbf{I} + h\mathbf{G}_{pq})^{-1} & h\mathbf{G}_{pp}(\mathbf{I} + h\mathbf{G}_{pq})^{-1} \\ -h\mathbf{G}_{qq}(\mathbf{I} + h\mathbf{G}_{pq})^{-1} & \mathbf{I} + h\mathbf{G}_{qp} - h^2\mathbf{G}_{pp}\mathbf{G}_{qq}(\mathbf{I} + h\mathbf{G}_{pq})^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{O} & \mathbf{I} \\ -\mathbf{I} & \mathbf{O} \end{pmatrix} \\ &= \begin{pmatrix} -h\mathbf{G}_{pp}(\mathbf{I} + h\mathbf{G}_{pq})^{-1} & (\mathbf{I} + h\mathbf{G}_{pq})^{-1} \\ -\mathbf{I} - h\mathbf{G}_{qp} + h^2\mathbf{G}_{pp}\mathbf{G}_{qq}(\mathbf{I} + h\mathbf{G}_{pq})^{-1} & -h\mathbf{G}_{qq}(\mathbf{I} + h\mathbf{G}_{pq})^{-1} \end{pmatrix}\end{aligned}\quad (10)$$

And then it is clear that  $(\Phi^\top J)\Phi = J$ . So the method is indeed symplectic.  $\square$

(c) If the method is separable,  $H(\mathbf{p}, \mathbf{q}) = T(\mathbf{p}) + V(\mathbf{q})$ , then it becomes:

$$\begin{cases} \mathbf{p}_{n+1} = \mathbf{p}_n - h \frac{\partial V(\mathbf{q}_n)}{\partial \mathbf{q}} \\ \mathbf{q}_{n+1} = \mathbf{q}_n + h \frac{\partial T(\mathbf{p}_{n+1})}{\partial \mathbf{p}} \end{cases}\quad (11)$$

is an explicit method. Denote  $\nabla_{\mathbf{q}}V = \mathbf{f}$  and  $\nabla_{\mathbf{p}}T = \mathbf{g}$ , and suppose we know them explicitly. The method is implemented as follows:

Given  $\mathbf{q}_0, \mathbf{p}_0$ .

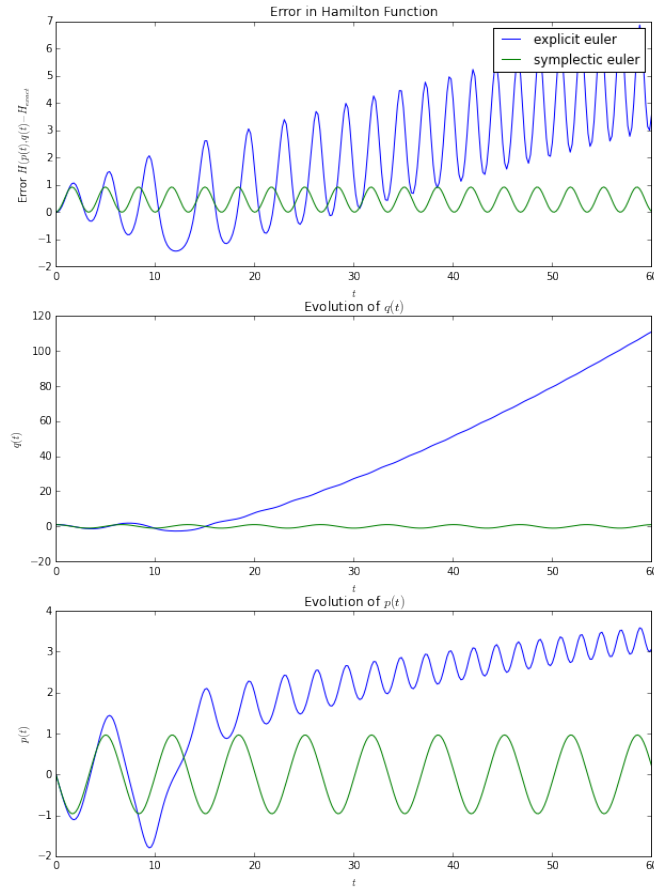
For  $k = 1$  to  $n$ :

$$\mathbf{p}_k \leftarrow \mathbf{p}_{k-1} - h\mathbf{f}(\mathbf{q}_k)$$

$$\mathbf{q}_k \leftarrow \mathbf{q}_{k-1} + h\mathbf{g}(\mathbf{p}_k) = \mathbf{q}_{k-1} + h\mathbf{g}(\mathbf{p}_{k-1} - h\mathbf{f}(\mathbf{q}_k))$$

which is explicit, and resembles the Runge-Kutta method.  $\square$

(d) The plots below gives the numerical solution of this problem choosing  $n = 300$ ,  $h = 0.2$ . I.e. solve the ode for  $t = 0.2, 0.4, \dots, 59.8, 60$  with both explicit euler and symplectic euler method.

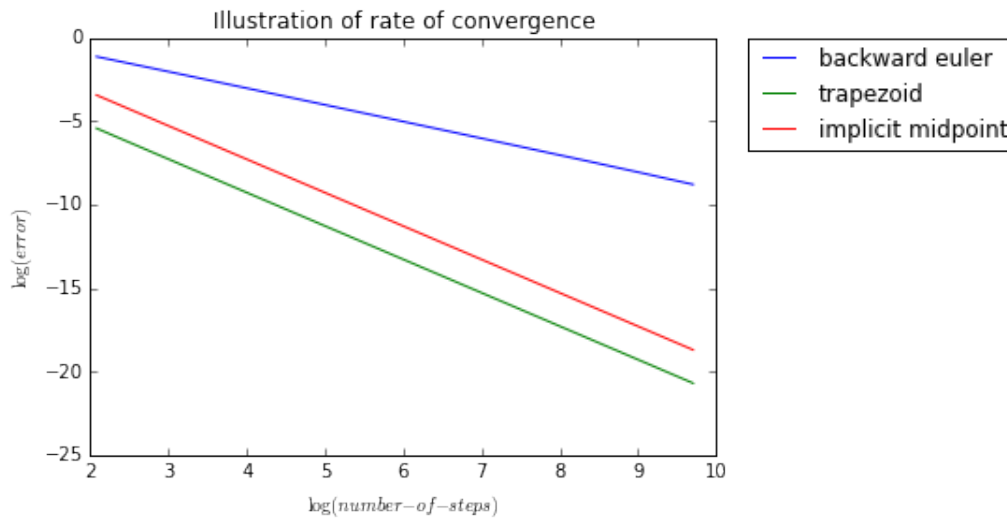


It is clear that for explicit euler method,  $H$  gradually grows with  $t$  as well as  $q(t)$ , and finally the object

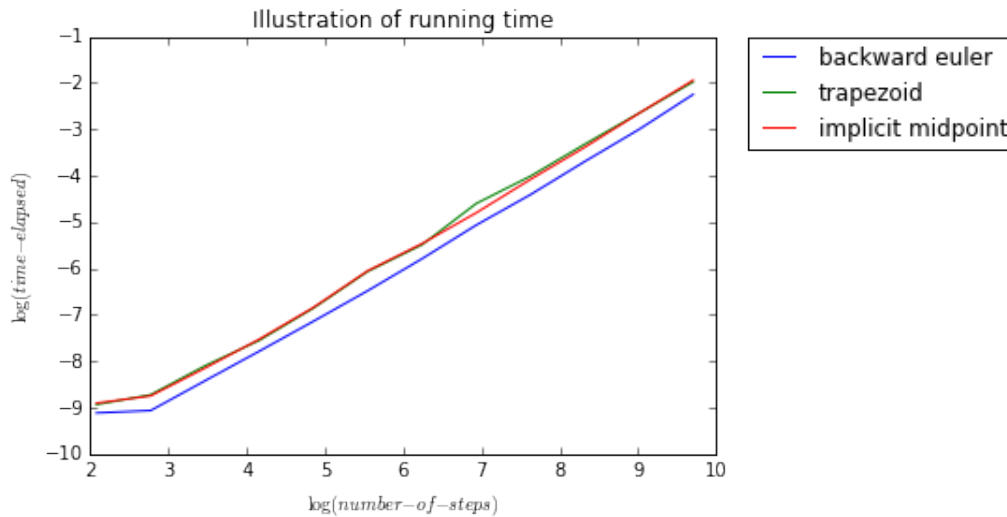
will go to infinitely faraway. For symplectic method  $H$  is bounded in a range, and the motion of the pendulum has a regular behavior.  $\square$

**Problem 5.** Implement implicit Euler, implicit midpoint method and trapezoid method, compare their rates of convergence and execution time.

**Solution.** We plot  $\log(n)$  against  $\log(\text{error})$  to check the rate of convergence. Theoretically, for a method of order  $k$ , we expect to get a line with slope  $-k$ .

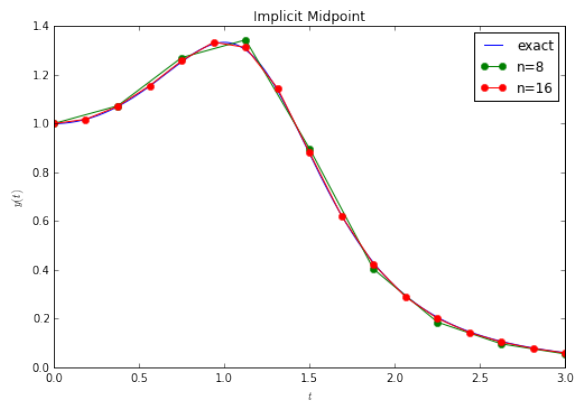
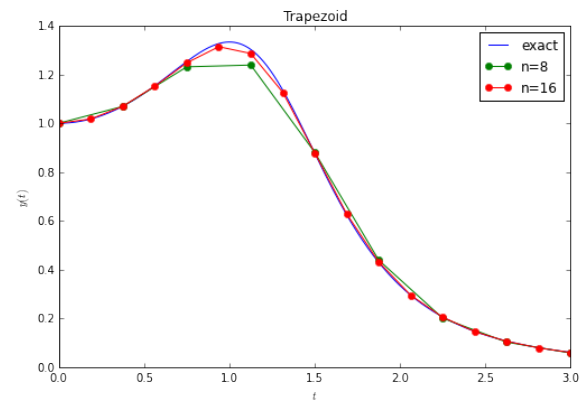
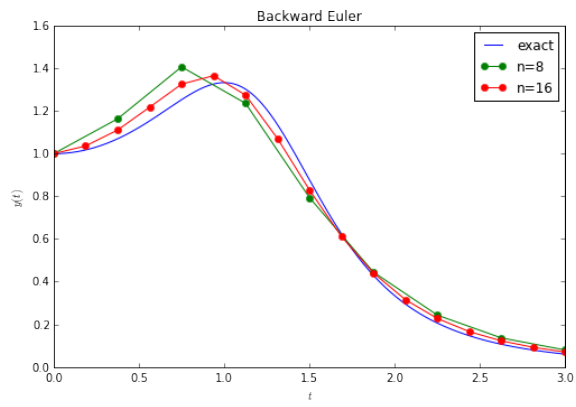


Which is indeed the case. We see that backward euler converges in order 1, other 2 methods have order 2.



We used the *Secant Method* find roots for these implicit methods. The execution time are all  $O(nm)$ . Where  $m$  is the average number of iterations it takes to find the zeros of implicit updating rules.

Finally we applied these methods to the testcase 5 in last homework. And plotted the numerical solution as below.



Please also see the attached code for details.

---