

MSAS – Assignment #1: Simulation

Lorenzo Cucchi, 221732

1 Implicit equations

Exercise 1

Let \mathbf{f} be a two-dimensional vector-valued function $\mathbf{f}(\mathbf{x}) = (x_2^2 - x_1 - 2, -x_1^2 + x_2 + 10)^\top$, where $\mathbf{x} = (x_1, x_2)^\top$. Find the zero(s) of \mathbf{f} by using Newton's method with $\partial\mathbf{f}/\partial\mathbf{x}$ 1) computed analytically, and 2) estimated through finite differences. Which version is more accurate?

(3 points)

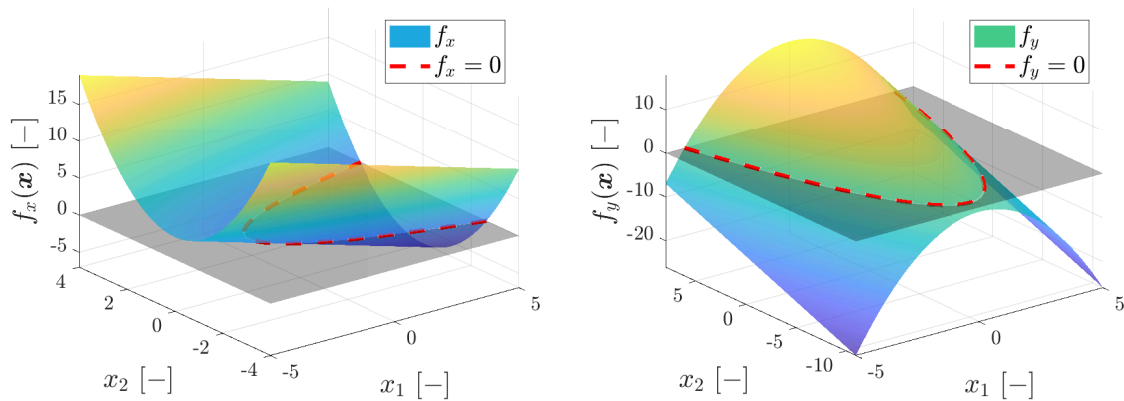


Figure 1: Surf plot of function $f_x(\mathbf{x}) = x_2^2 - x_1 - 2$ and function $f_y(\mathbf{x}) = -x_1^2 + x_2 + 10$.

The problem consists in finding the couple(s) $\{x_1, x_2\}$ which satisfy the relation $\mathbf{f}(x_1, x_2) = [0, 0]^\top$. In order to solve the problem using Newton's method, the analytical form or an approximation of the inverse Jacobian matrix \mathbf{J}^{-1} is needed. In fact, according to Newton's method, the i -th iteration \mathbf{x}_i is given by Equation 1.

$$\mathbf{x}_i = \mathbf{x}_{i-1} - \mathbf{J}^{-1}(\mathbf{x}_{i-1})\mathbf{f}(\mathbf{x}_{i-1}) \quad (1)$$

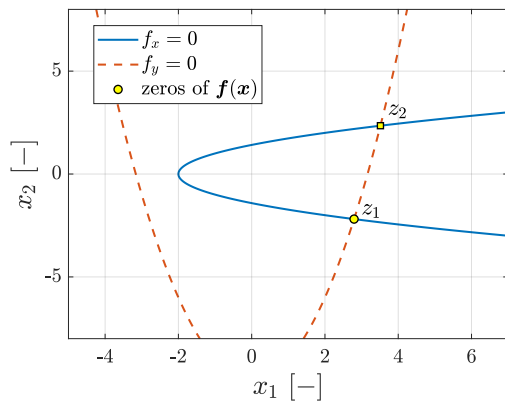


Figure 2: Zeros (z_1 and z_2) of the problem.

The analytical inverse Jacobian matrix can be efficiently obtained using the MATLAB symbolic toolbox through a few simple steps. Notably, for cases with a straightforward structure, like this one, manual calculation proves to be a faster alternative. On the other side, the approximated form of the Jacobian matrix is obtained through finite difference methods. In particular, first order *forward* and *centered* finite difference schemes are implemented. In this sense, the small increment δ which is applied in the finite difference methods follows the *rule*: $\delta = \sqrt{\text{eps}} \cdot \max(1, |\mathbf{x}|)$, where *eps* is the machine epsilon. Furthermore, to improve code efficiency and accuracy, the Jacobian calculated with finite differences is never inverted; instead, the MATLAB command `mldivide` or `\` is used.

As Figure 2 shows, there are two zeros of the function $\mathbf{f}(\mathbf{x})$. Indeed, two appropriate initial guesses \mathbf{x}_0 must be found to make the algorithms converge at the two zeros. The initial guesses are therefore $x_{0,1} = [1, -4]^T$ and $x_{0,2} = [6, 5]^T$. The stopping criteria chosen is the accuracy of the function evaluated in \mathbf{x}_i : when both the absolute values of $f_x(\mathbf{x}_i)$ and $f_y(\mathbf{x}_i)$ are lower than a tolerance set to $1e - 8$ the algorithm stops. Results reported in Table 1 show that the three algorithms converge at very close values and take the same number of iterations: the error $\|err\| = \|\mathbf{f}(\mathbf{x}_{end,analytical}) - \mathbf{f}(\mathbf{x}_{end,method})\|$ is very low, suggesting the high accuracy of both the forward and centered differences approximations.

Method	\mathbf{z}_1	Iterations	$\ err\ $
Analytical	$[2.794695112889339, -2.189679226029504]^T$	5	-
Forward differences	$[2.794695112889365, -2.189679226029472]^T$	5	4.1405e-14
Centered differences	$[2.794695112889314, -2.189679226029497]^T$	5	2.6291e-14
Method	\mathbf{z}_2	Iterations	$\ err\ $
Analytical	$[3.513999235947622, 2.348190630240421]^T$	5	-
Forward differences	$[3.513999235947627, 2.348190630240434]^T$	5	1.3773e-14
Centered differences	$[3.513999235947620, 2.348190630240424]^T$	5	3.1401e-15

Table 1: Zeros found by the three proposed methods.

2 Numerical solution of ODE

Exercise 2

The Initial Value Problem $\dot{x} = x - 2t^2 + 2$, $x(0) = 1$, has analytic solution $x(t) = 2t^2 + 4t - e^t + 2$.
 1) Implement a general-purpose, fixed-step Heun's method (RK2); 2) Solve the IVP in $t \in [0, 2]$ for $h_1 = 0.5$, $h_2 = 0.2$, $h_3 = 0.05$, $h_4 = 0.01$ and compare the numerical vs the analytical solution; 3) Repeat points 1)–2) with RK4; 4) Trade off between CPU time & integration error. (4 points)

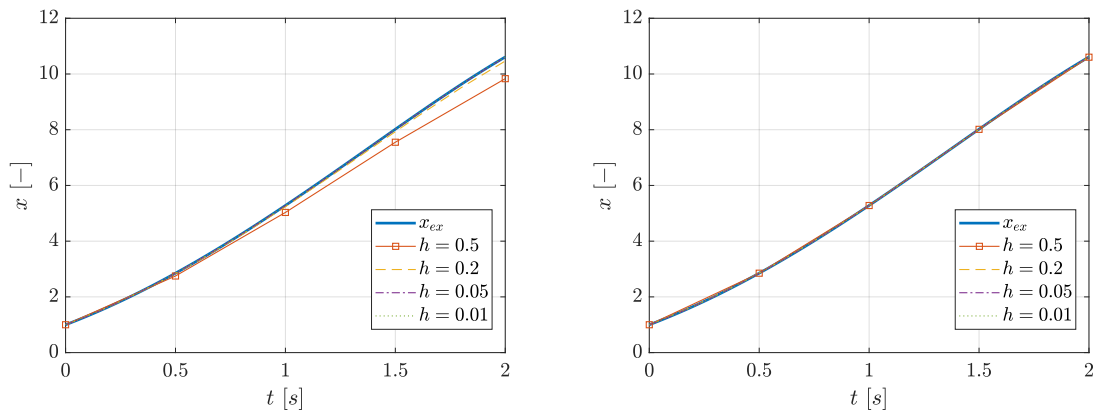


Figure 3: Solutions provided by with *RK2* (left) and *RK4* (right) methods by varying h value.

The initial value problem (IVP) over the interval $t \in [0, 2]$ is solved using step sizes $h_1 = 0.5$, $h_2 = 0.2$, $h_3 = 0.05$, and $h_4 = 0.01$ with both RK2 and RK4 methods. The results are illustrated in (Figure 3). Subsequently, the integration errors are compared in Figure 4. As evident from the figures, the RK4 method exhibits superior accuracy compared to RK2, even when using larger step sizes. The higher order of the RK4 method not only leads to increased accuracy but also results in a more substantial reduction in the global integration error, as demonstrated in Figure 5 (left). Considering the balance between computational time and integration error, Figure 5 (right) indicates that, to achieve the same level of accuracy in the final value, the RK4 method requires less time compared to the RK2 method. Therefore, the RK4 method is recommended irrespective of computational time considerations, as it proves to be the most efficient in terms of both accuracy and CPU time.

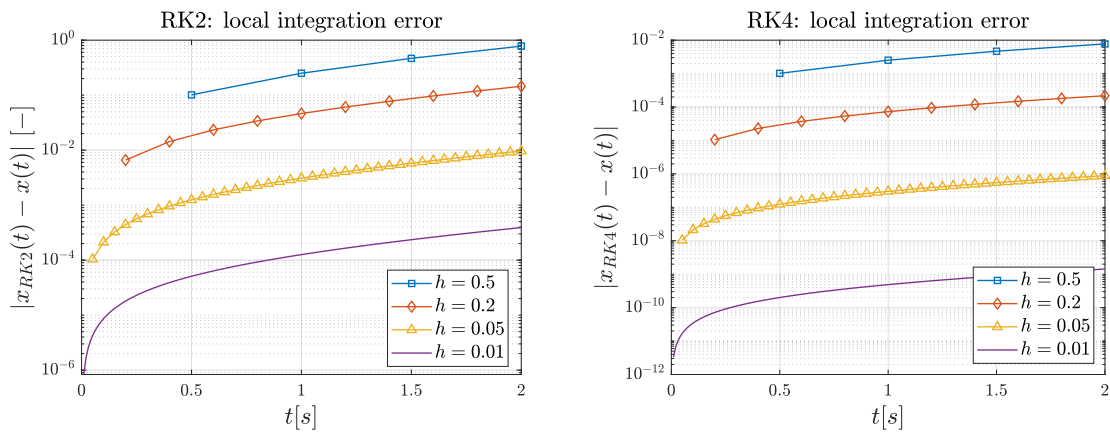


Figure 4: Local integration errors provided by with *RK2* (left) and *RK4* (right) methods by varying h value.

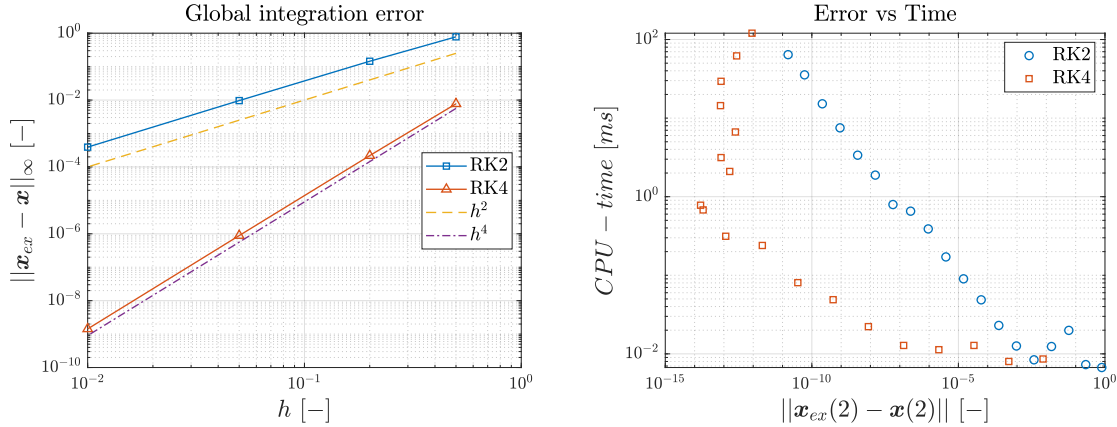


Figure 5: Global integration errors (left) and comparison between error and computational time of RK2 and RK4 methods (right).

Exercise 3

Let $\dot{\mathbf{x}} = A(\alpha)\mathbf{x}$ be a two-dimensional system with $A(\alpha) = [0, 1; -1, 2\cos\alpha]$. Notice that $A(\alpha)$ has a pair of complex conjugate eigenvalues on the unit circle; α denotes the angle from the $\text{Re}\{\lambda\}$ -axis. 1) Write the operator $F_{RK2}(h, \alpha)$ that maps \mathbf{x}_k into \mathbf{x}_{k+1} , namely $\mathbf{x}_{k+1} = F_{RK2}(h, \alpha)\mathbf{x}_k$. 2) With $\alpha = \pi$, solve the problem “Find $h \geq 0$ s.t. $\max(|\text{eig}(F(h, \alpha))|) = 1$ ”. 3) Repeat point 2) for $\alpha \in [0, \pi]$ and draw the solutions in the (h, λ) -plane. 4) Repeat points 1)–3) with RK4.

(5 points)

In order to retrieve the expression of the linear operator $F_{RK2}(h, \alpha)$ a generic RK2 iteration with step h is derived:

$$\begin{cases} \mathbf{x}_{k+1}^P = \mathbf{x}_k + h\mathbf{f}(\mathbf{x}_k, t_k) \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \frac{h}{2}[\mathbf{f}(\mathbf{x}_k, t_k) + \mathbf{f}(\mathbf{x}_{k+1}^P, t_{k+1})] \end{cases} \quad (2)$$

where, in our case, $\mathbf{f}(\mathbf{x}_k, t_k) = \mathbf{A}(\alpha)\mathbf{x}_k$. By substituting the first equation of Equation 2 in the second one it is obtained:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \frac{h}{2}\mathbf{A}(\alpha)\mathbf{x}_k + \frac{h}{2}\mathbf{A}(\alpha)\mathbf{x}_k + \frac{h}{2}h\mathbf{A}(\alpha)^2\mathbf{x}_k \quad (3)$$

By rearranging:

$$\mathbf{x}_{k+1} = (\mathbf{I} + h\mathbf{A}(\alpha) + \frac{h^2}{2}\mathbf{A}^2(\alpha))\mathbf{x}_k = \mathbf{F}_{RK2}(h, \alpha)\mathbf{x}_k \quad (4)$$

The same procedure can be performed to find F_{RK4} :

$$\mathbf{F}_{RK4}(h, \alpha) = \mathbf{I} + h\mathbf{A}(\alpha) + \frac{h^2}{2}\mathbf{A}^2(\alpha) + \frac{h^3}{6}\mathbf{A}^3(\alpha) + \frac{h^4}{24}\mathbf{A}^4(\alpha) \quad (5)$$

Table 2 shows the solutions of the statement “Find $h \geq 0$ s.t. $\max(|\text{eig}(F(h, \alpha))|) = 1$ ” imposing $\alpha = \pi$ with both F_{RK2} and F_{RK4} .

	RK2	RK4
h	2.0000000	2.7852935

Table 2: Results of statement “Find $h \geq 0$ s.t. $\max(|\text{eig}(F(h, \alpha))|) = 1$ ” where $\alpha = \pi$ with F_{RK2} and F_{RK4} functions.

Solving the problem for $\alpha \in [0, \pi]$ allows us to ascertain and visualize the numerical stability domains for both RK2 and RK4 methods, as depicted in Figure 6. As illustrated, the stability domains expand with higher approximation orders, necessitating larger step sizes for higher-order algorithms.

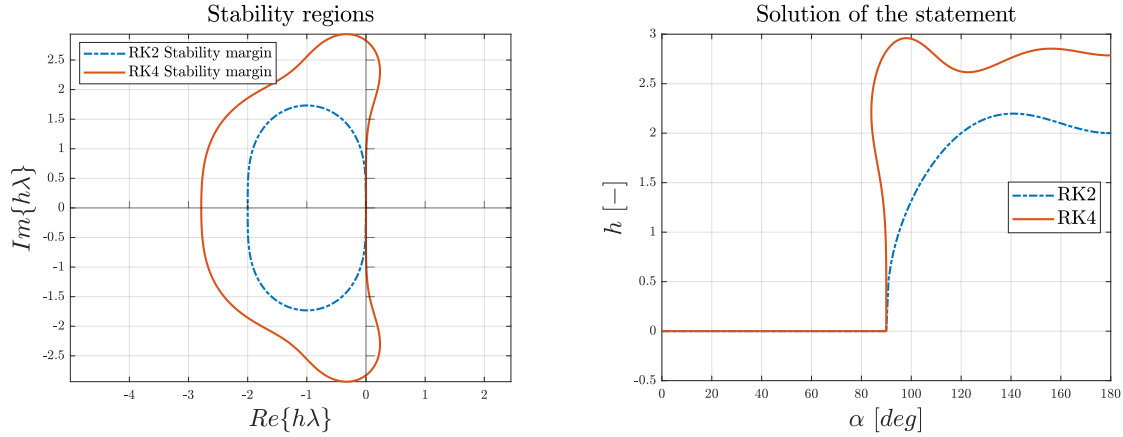


Figure 6: RK2 and RK4 method stability domain (left) and solution of the statement “Find $h \geq 0$ s.t. $\max(|\text{eig}(F(h, \alpha))|) = 1$ ” with RK2 and RK4 (right).

Exercise 4

Consider the IVP $\dot{\mathbf{x}} = A(\alpha)\mathbf{x}$, $\mathbf{x}(0) = [1, 1]^T$, to be integrated in $t \in [0, 1]$. 1) Take $\alpha \in [0, \pi]$ and solve the problem “Find $h \geq 0$ s.t. $\|\mathbf{x}_{\text{an}}(1) - \mathbf{x}_{\text{RK1}}(1)\|_{\infty} = \text{tol}$ ”, where $\mathbf{x}_{\text{an}}(1)$ and $\mathbf{x}_{\text{RK1}}(1)$ are the analytical and the numerical solution (with RK1) at the final time, respectively, and $\text{tol} = \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. 2) Plot the four locus of solutions in the $(h\lambda)$ -plane; plot also the function evaluations vs tol for $\alpha = \pi$. 3) Repeat points 1)–2) for RK2 and RK4.

(4 points)

The outcomes are illustrated in Figure 7 and Figure 8. As observed in the figures, the use of higher-order approximation methods allows for larger permissible values of the step size h . Consequently, to maintain the same error tolerance $\|\mathbf{x}_{\text{an}}(1) - \mathbf{x}_{\text{RK1}}(1)\|_{\infty}$, higher-order methods can accommodate larger h , resulting in fewer steps at which the function needs to be evaluated, this observation is further exemplified in Figure 8.

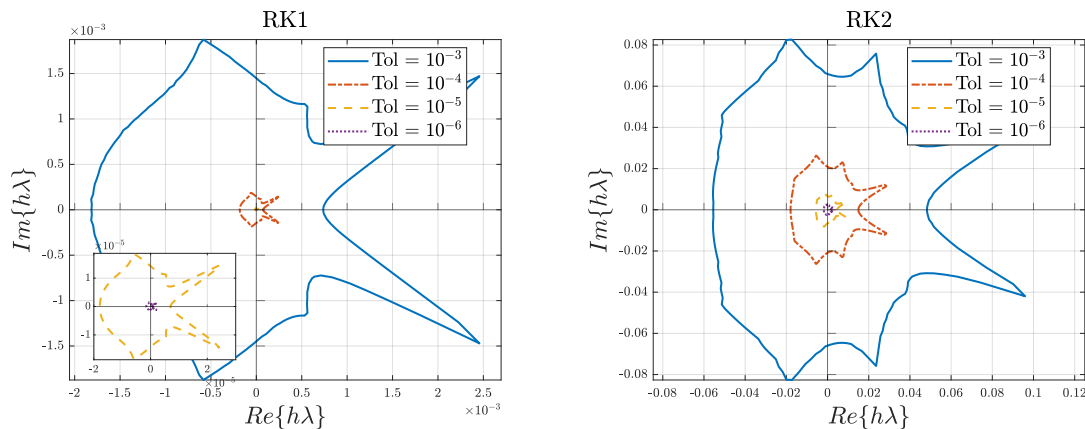


Figure 7: Five locus of solutions for RK1 (left) and RK2 (right) methods.

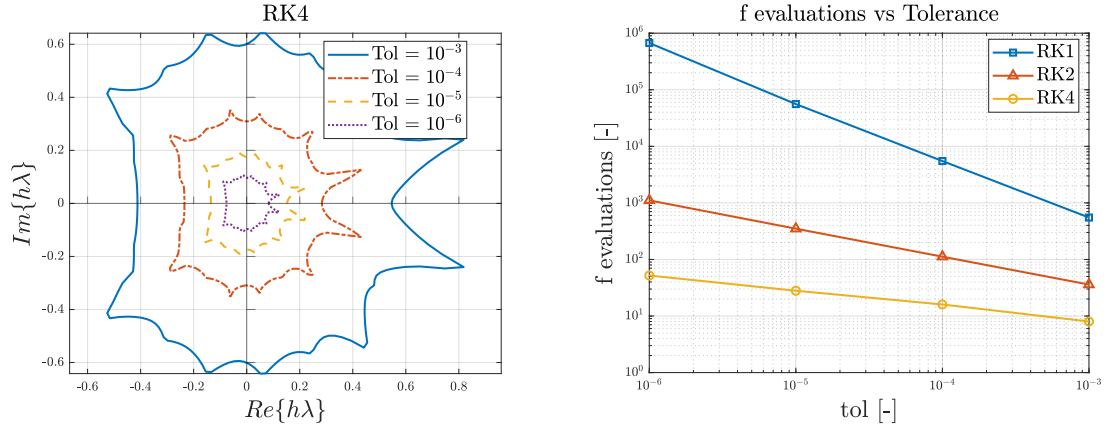


Figure 8: RK4 locus of solutions (left) and function evaluations vs. tolerance (right).

Exercise 5

Consider the backinterpolation method $BI_{2,0.4}$. 1) Derive the expression of the linear operator $B_{BI_{2,0.4}}(h, \alpha)$ such that $\mathbf{x}_{k+1} = B_{BI_{2,0.4}}(h, \alpha)\mathbf{x}_k$. 2) Following the approach of point 3) in Exercise 3, draw the stability domain of $BI_{2,0.4}$ in the $(h\lambda)$ -plane. 3) Derive the domain of numerical stability of $BI_{2,\theta}$ for the values of $\theta = [0.1, 0.3, 0.7, 0.9]$.

(5 points)

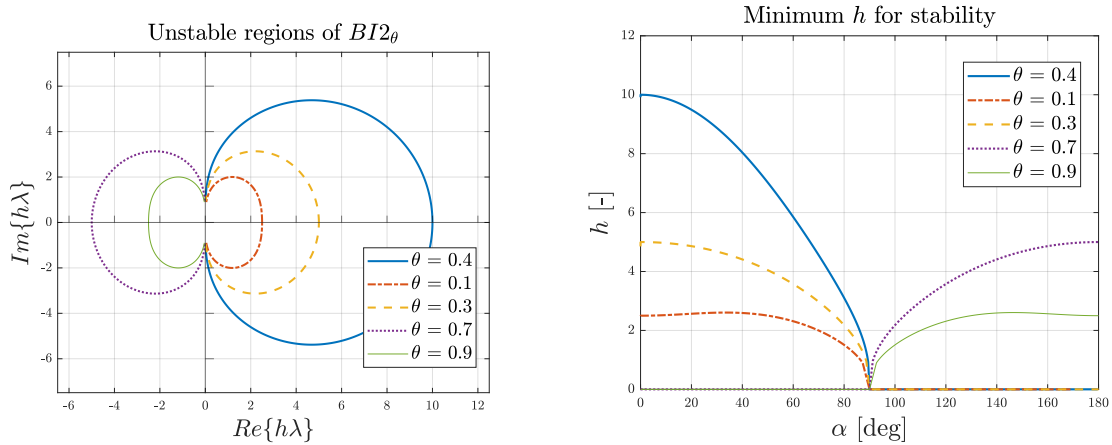


Figure 9: Unstable domains of $BI_{2,\theta}$ for different θ values (left). Minimum step h for stability (right).

BI_i backinterpolation methods are a special implicit Runge-Kutta (IRK) methods. In order to derive the expression of the linear operator $B_{BI_{2,0.4}}(h, \alpha)$ such that $\mathbf{x}_{k+1} = B_{BI_{2,0.4}}(h, \alpha)\mathbf{x}_k$, a generic iteration of RK2 with step $h\theta$ is derived:

$$\begin{cases} \mathbf{x}_{k+h\theta}^P = \mathbf{x}_k + \theta h \mathbf{f}(\mathbf{x}_k, t_k) \\ \mathbf{x}_{k+h\theta}^C = \mathbf{x}_k + \frac{\theta h}{2} [\mathbf{f}(\mathbf{x}_k, t_k) + \mathbf{f}(\mathbf{x}_{k+h\theta}^P, t_{k+h\theta})] \end{cases} \quad (6)$$

where, in our case, $\mathbf{f}(\mathbf{x}_k, t_k) = \mathbf{A}(\alpha)\mathbf{x}_k$. By substituting the latter and the first equation of Equation 6 in the second one, what is obtained is the following:

$$\mathbf{x}_{k+h\theta} = \mathbf{x}_k + \frac{\theta h}{2} \mathbf{A}(\alpha)\mathbf{x}_k + \frac{\theta h}{2} \mathbf{A}(\alpha)\mathbf{x}_k + \frac{\theta h}{2} \theta h \mathbf{A}^2(\alpha) \mathbf{x}_k \quad (7)$$

Which can be rewritten as:

$$\mathbf{x}_{k+h\theta} = (\mathbf{I} + h\theta \mathbf{A}(\alpha) + \frac{\theta^2 h^2}{2} \mathbf{A}^2(\alpha)) \mathbf{x}_k \quad (8)$$

The same procedure is performed to find the relation between $\mathbf{x}_{k+h\theta}$ and \mathbf{x}_{k+1} . In order to achieve that, the generic RK2 iteration with $-h(1-\theta)$ step is considered:

$$\begin{cases} \mathbf{x}_{k+h\theta}^P = \mathbf{x}_k - h(1-\theta)\mathbf{f}(\mathbf{x}_{k+1}, t_{k+1}) \\ \mathbf{x}_{k+h\theta}^C = \mathbf{x}_k - \frac{h(1-\theta)}{2}[\mathbf{f}(\mathbf{x}_{k+1}, t_{k+1}) + \mathbf{f}(\mathbf{x}_{k+h\theta}^P, t_{k+h\theta})] \end{cases} \quad (9)$$

Following the procedure that allowed to find Equation 8, Equation 9 becomes:

$$\mathbf{x}_{k+h\theta} = (\mathbf{I} - h(1-\theta)\mathbf{A}(\alpha) + \frac{h^2(1-\theta)^2}{2}\mathbf{A}^2(\alpha))\mathbf{x}_{k+1} \quad (10)$$

By comparing Equation 8 and Equation 9:

$$(\mathbf{I} + h\theta\mathbf{A}(\alpha) + \frac{\theta^2 h^2}{2}\mathbf{A}^2(\alpha))\mathbf{x}_k = (\mathbf{I} - h(1-\theta)\mathbf{A}(\alpha) + \frac{h^2(1-\theta)^2}{2}\mathbf{A}^2(\alpha))\mathbf{x}_{k+1} \quad (11)$$

By isolating \mathbf{x}_{k+1} :

$$\mathbf{x}_{k+1} = (\mathbf{I} - h(1-\theta)\mathbf{A}(\alpha) + \frac{h^2(1-\theta)^2}{2}\mathbf{A}^2(\alpha))^{-1}(\mathbf{I} + h\theta\mathbf{A}(\alpha) + \frac{\theta^2 h^2}{2}\mathbf{A}^2(\alpha))\mathbf{x}_k \quad (12)$$

As a result, the linear operator $B_{BI2_\theta}(h, \alpha)$ is obtained:

$$B_{BI2_{0.4}}(h, \alpha) = (\mathbf{I} - h(1-\theta)\mathbf{A}(\alpha) + \frac{h^2(1-\theta)^2}{2}\mathbf{A}^2(\alpha))^{-1}(\mathbf{I} + h\theta\mathbf{A}(\alpha) + \frac{\theta^2 h^2}{2}\mathbf{A}^2(\alpha)) \quad (13)$$

By imposing $\theta = 0.4$ the linear operator $B_{BI2_{0.4}}(h, \alpha)$ is obtained. Given the expression of the linear operator $B_{BI2_\theta}(h, \alpha)$, it is possible to derive the domain of numerical stability of $BI2_\theta$ method with different θ values by solving “Find $h \geq 0$ s.t. $\max(|\text{eig}(BI2_\theta(h, \alpha))|) = 1$ ”. Results are shown in Figure 9.

Exercise 6

Consider the IVP $\dot{\mathbf{x}} = \mathbf{B}\mathbf{x}$ with $\mathbf{B} = [-180.5, 219.5; 179.5, -220.5]$ and $\mathbf{x}(0) = [1, 1]^T$ to be integrated in $t \in [0, 5]$. Notice that $\mathbf{x}(t) = e^{\mathbf{B}t}\mathbf{x}(0)$. 1) Solve the IVP using RK4 with $h = 0.1$; 2) Repeat point 1) using implicit extrapolation technique IEX4; 3) Compare the numerical results in points 1) and 2) against the analytic solution; 4) Compute the eigenvalues associated to the IVP and represent them on the $(h\lambda)$ -plane both for RK4 and IEX4; 5) Discuss the results.

Imposing $h = 0.1$, the problem $\dot{\mathbf{x}} = \mathbf{B}\mathbf{x}$ is solved with both RK4 and IEX4 methods and the obtained solutions are compared to the analytical one ($x_{\text{analytical}} = e^{\mathbf{B}t}$) in Figure 11. While the fourth order implicit extrapolation technique IEX4 approximates the analytical solution with low errors, integration errors due to RK4 approximation diverges. This is due to the eigenvalues of the matrix \mathbf{B} : $\lambda_i = [-1, -400]$. In particular, the eigenvalue $\lambda_1 = -1$ multiplied by the step $h = 0.1$ lies in the unstable domain of the RK4 method, making the associated component of the solution diverges from the analytical one. Although the second eigenvalue lies in the stable domain of RK4, the solution still diverges since the state equations are coupled. A different behaviour is found with IEX4 method: as Figure 10 both the $h\lambda_i$ lie in the stable domain of the method, making the solution converges to the analytical one with low integration errors.

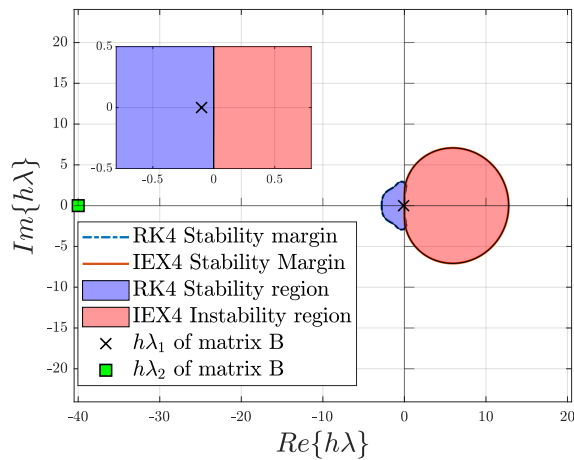


Figure 10: Stable and unstable domains of RK4 and IEX4 methods.

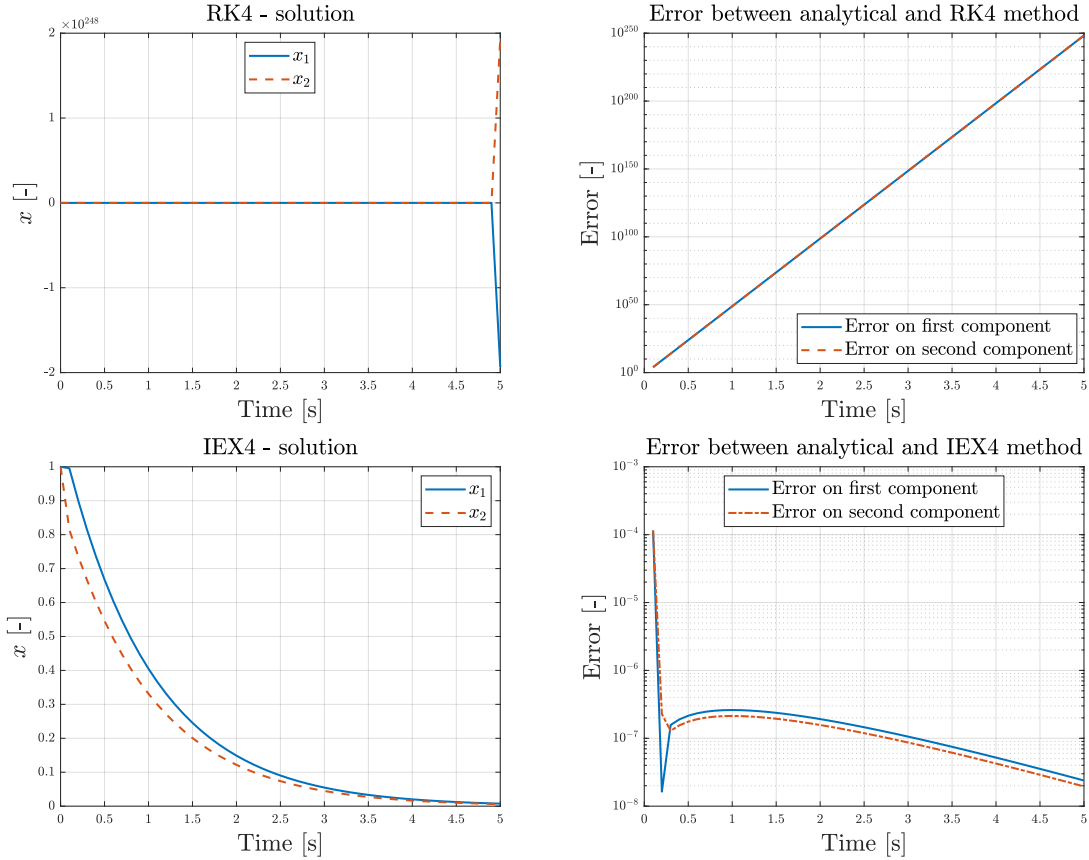


Figure 11: Solutions obtained with RK4 and IEX4 methods and their respective errors compared with analytic solution.

Exercise 7

Consider the two-dimensional IVP

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} -\frac{5}{2} [1 + 8 \sin(t)] x_1 \\ (1 - x_1)x_2 + x_1 \end{bmatrix}, \quad \begin{bmatrix} x_1(t_0) \\ x_2(t_0) \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

- 1) Solve the IVP using AB3 in $t \in [0, 3]$ for $h = 0.1$; 2) Repeat point 1) using AM3, ABM3, and BDF3; 3) Discuss the results.

(5 points)

Figure 12 shows the solution of the problem when a step size $h = 0.1$ is adopted with AB3, AM3, ABM3 and BDF3 methods. As the figure depicts, the AB3 method suffers from integration instability on both the components x and y . On the other side, AM3 and BDF3 methods can provide a reliable solution for the x component but not for y , while ABM3 method seems to suffer from instability for the x components between $t \simeq 1.5$ s and $t \simeq 2.5$ s and the same problem encountered with the other methods for the y component.

In order to study such a behaviour, the linearized system is studied. Regarded as \mathbf{x}_0 the equilibrium solution, the problem is formulated as follows:

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, t) \simeq \mathbf{f}(\mathbf{x}_0, t) + \left. \frac{\partial \mathbf{f}(\mathbf{x}, t)}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{x}_0} (\mathbf{x} - \mathbf{x}_0) \quad (14)$$

By definition, $\mathbf{f}(\mathbf{x}_0, t) = \mathbf{0}$. Furthermore, single \mathbf{x}_0 value is obtained by imposing $\dot{\mathbf{x}} = \mathbf{0}$: $\mathbf{x}_0 = [0, 0]^T$. As a result, Equation 14 becomes:

$$\mathbf{f}(\mathbf{x}, t) \simeq \left. \frac{\partial \mathbf{f}(\mathbf{x}, t)}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{x}_0} \mathbf{x} \quad (15)$$

The partial derivative term in Equation 15 is the Jacobian matrix of function $\mathbf{f}(\mathbf{x}, t)$:

$$\mathbf{J} = \mathbf{J}(\mathbf{x}, t) = \begin{bmatrix} \frac{5}{2} [1 + 8 \sin(t)] & 0 \\ 1 - y & 1 - x \end{bmatrix} \quad (16)$$

which, evaluated in \mathbf{x}_0 , reads:

$$\mathbf{J}(\mathbf{x}_0, t) = \begin{bmatrix} \frac{5}{2} [1 + 8 \sin(t)] & 0 \\ 1 & 1 \end{bmatrix} \quad (17)$$

As a result, the linearized problem reads:

$$\dot{\mathbf{x}} = \begin{bmatrix} \frac{5}{2} [1 + 8 \sin(t)] & 0 \\ 1 & 1 \end{bmatrix} \mathbf{x} \quad (18)$$

and, thus, the eigenvalues can be evaluated at each instant of time t . Left plot of Figure 13 shows the evolution in time of the $h\lambda_i$ values associated to the linearized problem of Equation 18. Looking at right plot of Figure 13, the behaviour of the four proposed methods is characterized:

- **AB3**: $h\lambda_y$ is never inside the method stability domain, so the y component diverges. With the exception of little intervals at the beginning and at the end of the time-span, the same situation is encountered for $h\lambda_x$;
- **AM3**: $h\lambda_y$ is never inside the method stability domain, so the y component diverges. On the other side, the $h\lambda_x$ is always in the stable domain, resulting in the convergence of the x component;
- **ABM3**: the situation is the same encountered with AM3 method, except for the fact that for a small interval between $t \simeq 1.5$ s and $t \simeq 2.5$ s the $h\lambda_x$ values are outside the stability domain;
- **ABM3**: the situation is the same encountered with AM3 method.

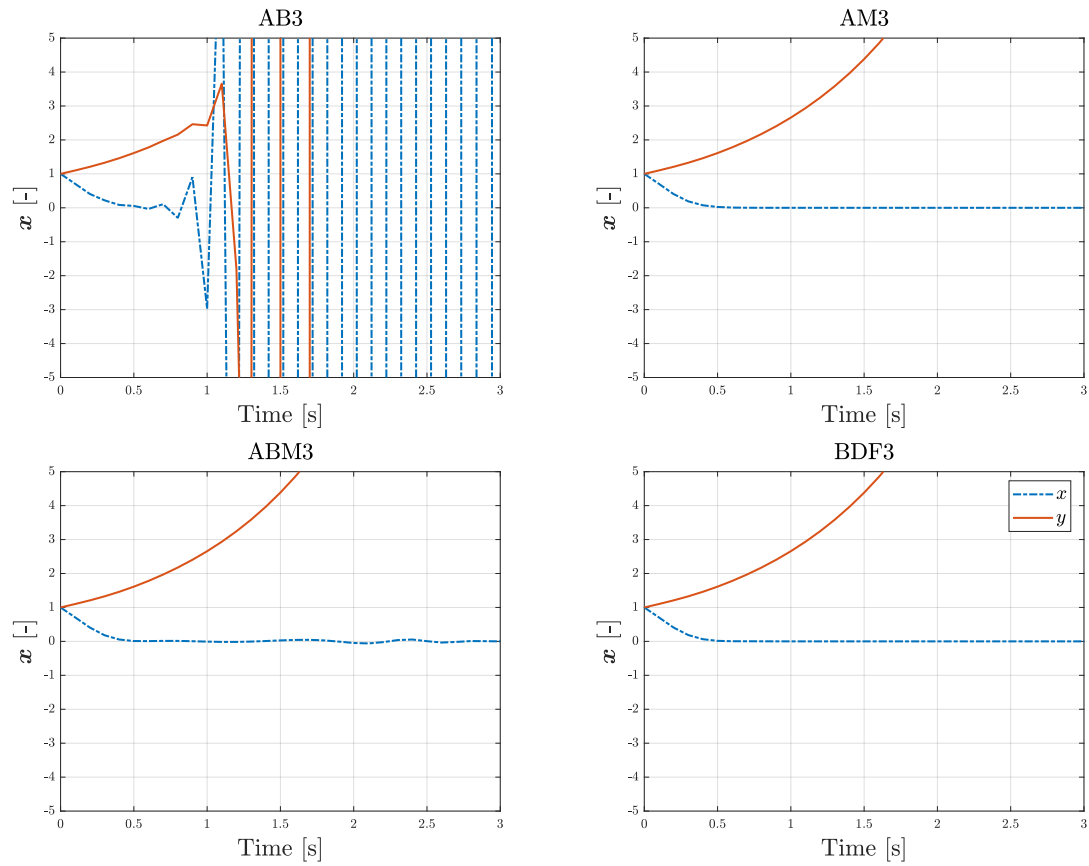


Figure 12: Solutions obtained with AB3, AM3, ABM3 and BDF3 methods.

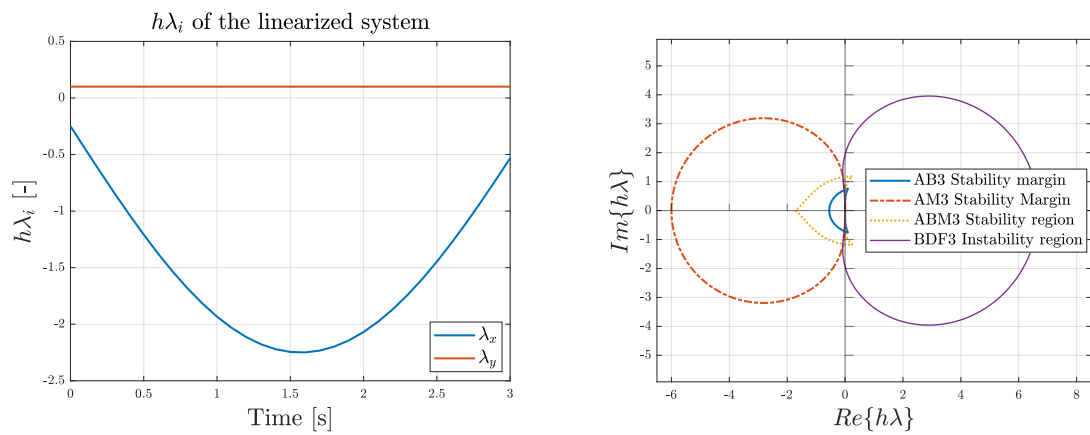


Figure 13: $h\lambda_i$ values of the linearized system (left). Stability and Instability domains of the proposed methods (right).