

Numerical methods for Partial Differential Equations

Adriano Festa
Politecnico of Turin, Italy
Undergraduate Lecture at KSU
2022



Table of Contents

- 1 The model of elastic string or heated beam
- 2 The model of elastic membrane or heated plate
- 3 Solution of large linear algebraic systems
- 4 Models of temporal evolution
- 5 Time-advancing schemes
- 6 Convection-diffusion and transport problems
- 7 Conservation Laws - Introduction to Finite Volumes

Consider the initial-value problem (also known as *Cauchy problem*)

$$\begin{cases} \mathbf{u}' = \mathbf{F}(\mathbf{u}, t) , & 0 < t \leq T , \\ \mathbf{u}(0) = \mathbf{u}_0 , \end{cases} \quad (62)$$

for a system of n first-order differential equations in n unknowns.

Here $\mathbf{u} : [0, T] \rightarrow \mathbb{R}^n$ is the vector, depending upon time in a differentiable way, that collects the unknowns, whereas $\mathbf{F} : \mathbb{R}^n \times [0, T] \rightarrow \mathbb{R}^n$ is a globally continuous function, which is *Lipschitz-continuous* in the variable \mathbf{u} uniformly in t , i.e., there exists a constant $L > 0$ such that

$$\|\mathbf{F}(\mathbf{u}, t) - \mathbf{F}(\mathbf{v}, t)\| \leq L \|\mathbf{u} - \mathbf{v}\| , \quad \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^n , \quad \forall t \in [0, T] .$$

For instance, this happens if

$$\mathbf{F}(\mathbf{u}, t) = \mathcal{A}\mathbf{u} + \mathbf{b}(t) ,$$

since

$$\|\mathbf{F}(\mathbf{u}, t) - \mathbf{F}(\mathbf{v}, t)\| = \|\mathcal{A}(\mathbf{u} - \mathbf{v})\| \leq \|\mathcal{A}\| \|\mathbf{u} - \mathbf{v}\| .$$

Under the assumptions made on \mathbf{F} , problem (62) admits one and only one solution for any choice of the initial data $\mathbf{u}_0 \in \mathbb{R}^n$.

For the numerical approximation of such a solution, the most popular strategy consists in resorting to a *time-advancing scheme*.

This means that one chooses (in various manners) $K > 0$ time instants t_k , with

$$0 = t_0 < t_1 < \cdots < t_k < t_{k+1} < \cdots < t_K = T ;$$

for each of them, one recursively defines an approximation

$$\mathbf{u}^k \simeq \mathbf{u}(t_k)$$

of the exact solution; precisely, at the time t_{k+1} , one defines \mathbf{u}^{k+1} using the knowledge of the approximations \mathbf{u}^j , $j \leq k$ already computed at one or more previous time instants.

The conceptually simplest situation occurs when one advances with a constant time step $\Delta t > 0$ (in this case, one sets $K = T/\Delta t$, assuming this ratio to be integer); the time instants are then defined by

$$t_k = k\Delta t , \quad \text{for } k = 0, 1, \dots, K .$$

Explicit Euler scheme (EE): (*explicit - one step - 1 order*)

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \Delta t \mathbf{F}(\mathbf{u}^k, t_k), \quad k \geq 0 ;$$

motivated by the *forward incremental quotient* formula

$$\frac{\mathbf{u}(t_{k+1}) - \mathbf{u}(t_k)}{\Delta t} \simeq \mathbf{u}'(t_k) = \mathbf{F}(\mathbf{u}(t_k), t_k) .$$

Implicit Euler scheme (IE): (*implicit - one step - 1 order*)

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \Delta t \mathbf{F}(\mathbf{u}^{k+1}, t_{k+1}), \quad k \geq 0 ;$$

motivated by the *backward incremental quotient* formula

$$\frac{\mathbf{u}(t_{k+1}) - \mathbf{u}(t_k)}{\Delta t} \simeq \mathbf{u}'(t_{k+1}) = \mathbf{F}(\mathbf{u}(t_{k+1}), t_{k+1})$$

Mid-point scheme (MP): (*explicit - two steps - II order*)

$$\mathbf{u}^{k+1} = \mathbf{u}^{k-1} + 2\Delta t \mathbf{F}(\mathbf{u}^k, t_k), \quad k \geq 1;$$

motivated by the *centered incremental quotient* formula

$$\frac{\mathbf{u}(t_{k+1}) - \mathbf{u}(t_{k-1}))}{2\Delta t} \simeq \mathbf{u}'(t_k) = \mathbf{F}(\mathbf{u}(t_k), t_k)$$

Trapezoidal (or Crank-Nicolson) scheme (CN): (*implicit - one step - II order*)

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \Delta t \left[\frac{1}{2} \mathbf{F}(\mathbf{u}^k, t_k) + \frac{1}{2} \mathbf{F}(\mathbf{u}^{k+1}, t_{k+1}) \right], \quad k \geq 0.$$

motivated by a *numerical quadrature formula* applied to the relation

$$\mathbf{u}(t_{k+1}) - \mathbf{u}(t_k) = \int_{t_k}^{t_{k+1}} \mathbf{u}'(t) dt = \int_{t_k}^{t_{k+1}} \mathbf{F}(\mathbf{u}(t), t) dt,$$

obtained by integrating the identity $\mathbf{u}'(t) = \mathbf{F}(\mathbf{u}(t), t)$ on the interval $[t_k, t_{k+1}]$. The integral on the right-hand side can be approximated by the *trapezoidal formula*

$$\int_{t_k}^{t_{k+1}} \mathbf{F}(\mathbf{u}(t), t) dt \simeq \Delta t \left[\frac{1}{2} \mathbf{F}(\mathbf{u}(t_k), t_k) + \frac{1}{2} \mathbf{F}(\mathbf{u}(t_{k+1}), t_{k+1}) \right].$$

Runge-Kutta scheme RK3: (*explicit - one-step - III order*)

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \Delta t \left[\frac{2}{9}\mathbf{F}_1 + \frac{1}{3}\mathbf{F}_2 + \frac{4}{9}\mathbf{F}_3 \right], \quad k \geq 0,$$

where the addends on the right-hand side are recursively defined by

$$\mathbf{F}_1 = \mathbf{F}(\mathbf{u}^k, t_k),$$

$$\mathbf{F}_2 = \mathbf{F}\left(\mathbf{u}^k + \frac{1}{2}\Delta t \mathbf{F}_1, t_k + \frac{1}{2}\Delta t\right),$$

$$\mathbf{F}_3 = \mathbf{F}\left(\mathbf{u}^k + \frac{3}{4}\Delta t \mathbf{F}_2, t_k + \frac{3}{4}\Delta t\right).$$

Runge-Kutta scheme RK4: (*explicit - one-step - IV order*)

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \Delta t \left[\frac{1}{6} \mathbf{F}_1 + \frac{1}{3} \mathbf{F}_2 + \frac{1}{3} \mathbf{F}_3 + \frac{1}{6} \mathbf{F}_4 \right], \quad k \geq 0,$$

where the addends on the right-hand side are recursively defined by

$$\begin{aligned} \mathbf{F}_1 &= \mathbf{F}(\mathbf{u}^k, t_k), \\ \mathbf{F}_2 &= \mathbf{F}(\mathbf{u}^k + \tfrac{1}{2} \Delta t \mathbf{F}_1, t_k + \tfrac{1}{2} \Delta t), \\ \mathbf{F}_3 &= \mathbf{F}(\mathbf{u}^k + \tfrac{1}{2} \Delta t \mathbf{F}_2, t_k + \tfrac{1}{2} \Delta t), \\ \mathbf{F}_4 &= \mathbf{F}(\mathbf{u}^k + \Delta t \mathbf{F}_3, t_k + \Delta t). \end{aligned}$$

“Backward Difference Formula” scheme BDF2: (*implicit - two-steps - II order*)

$$\frac{3}{2}\mathbf{u}^{k+1} - 2\mathbf{u}^k + \frac{1}{2}\mathbf{u}^{k-1} = \Delta t \mathbf{F}(\mathbf{u}^{k+1}, t_{k+1}), \quad k \geq 1,$$

“Backward Difference Formula” scheme BDF3: (*implicit - three-steps - III order*)

$$\frac{11}{6}\mathbf{u}^{k+1} - 3\mathbf{u}^k + \frac{3}{2}\mathbf{u}^{k-1} - \frac{1}{3}\mathbf{u}^{k-2} = \Delta t \mathbf{F}(\mathbf{u}^{k+1}, t_{k+1}), \quad k \geq 2.$$

An explicit one-step scheme has the general form

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \Delta t \Phi(\mathbf{u}^k, t_k, \Delta t), \quad k \geq 0,$$

whereas an implicit one-step scheme has the general form

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \Delta t \Psi(\mathbf{u}^k, \mathbf{u}^{k+1}, t_k, \Delta t), \quad k \geq 0.$$

The connection with the differential equation is given by the consistency condition:

Definition

A one-step scheme, defined by one of the previous formulas, is said to be **consistent** if the condition

$$\Phi(\mathbf{u}, t, 0) = F(\mathbf{u}, t) \quad \text{or} \quad \Psi(\mathbf{u}, \mathbf{u}, t, 0) = F(\mathbf{u}, t)$$

is satisfied for any $\mathbf{u} \in \mathbb{R}^n$ and $t \in [0, T]$.

All the previously presented one-step schemes are consistent.

Definition

A consistent one-step scheme is said to be **of order** $p > 0$ if, considering the solution \mathbf{u}^1 produced by the scheme at the first instant $t_1 = \Delta t$, one has

$$\|\mathbf{u}(t_1) - \mathbf{u}^1\| = O(\Delta t^{p+1}) \quad \text{for } \Delta t \rightarrow 0 ,$$

for all solutions \mathbf{u} that are $(p+1)$ -times differentiable at $t = 0$.

Theorem

If a one-step scheme is consistent and of order p , and if the exact solution $\mathbf{u}(t)$ is $(p+1)$ -times differentiable in $[0, T]$, one has

$$\max_{1 \leq k \leq K} \|\mathbf{u}(t_k) - \mathbf{u}^k\| \leq C_{L,T} \Delta t^p \max_{t \in [0,T]} \left\| \frac{d^{p+1} \mathbf{u}}{dt^{p+1}}(t) \right\| ,$$

where $C_{L,T}$ denotes a constant only depending on L , T and the numerical scheme.

The theorem guarantees the *convergence* of the scheme, i.e., the fact that the discrete solution generated by the numerical scheme converges towards the exact solution as $\Delta t \rightarrow 0$; it also predicts the behaviour of the error as the time step decreases.

In applications, it is important that a time-advancing scheme generates discrete solutions \mathbf{u}^k which stay bounded as $t_k \rightarrow +\infty$, whenever the exact solutions of the initial-value problem (62) stay bounded as $t \rightarrow +\infty$.

A quite relevant situation occurs when the differential system is

- *linear* and *autonomous*, i.e., one has $\mathbf{F}(\mathbf{u}, t) = \mathbf{F}(\mathbf{u}) = \mathcal{A}\mathbf{u} + \mathbf{b}$ with \mathcal{A} square matrix of order n and $\mathbf{b} \in \mathbb{R}^n$ independent of time, and
- *dissipative*, i.e., \mathcal{A} is diagonalizable with eigenvalues all having real part < 0 .

In this case, whichever is the initial datum \mathbf{u}_0 , the solution $\mathbf{u}(t)$ of the problem

$$\begin{cases} \mathbf{u}' = \mathcal{A}\mathbf{u} + \mathbf{b}, & t > 0, \\ \mathbf{u}(0) = \mathbf{u}_0, \end{cases} \quad (63)$$

converges, as $t \rightarrow +\infty$, towards the *steady state* $\mathbf{u}_\infty \in \mathbb{R}^n$, solution of the linear system

$$\mathcal{A}\mathbf{u}_\infty + \mathbf{b} = \mathbf{0};$$

hence, in particular, it stays bounded as $t \rightarrow +\infty$.

Indeed, applying the change of dependent variable $v(t) = u(t) - u_\infty$ and setting $v_0 = u_0 - u_\infty$, one immediately checks that v is the solution of the homogeneous Cauchy problem

$$\begin{cases} v' = \mathcal{A}v, & t > 0, \\ v(0) = v_0. \end{cases} \quad (64)$$

Let $\mathcal{A} = \mathbf{W}\mathbf{\Lambda}\mathbf{W}^{-1}$ be the diagonalization of the matrix \mathcal{A} , where $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$ is the diagonal matrix collecting the eigenvalues and $\mathbf{W} = (w_1, \dots, w_n)$ the matrix collecting the eigenvectors (placed column-wise).

Let us substitute such expression in (64) and let us multiply on the left by \mathbf{W}^{-1} . Setting $z(t) = \mathbf{W}^{-1}v(t)$ and $z_0 = \mathbf{W}^{-1}v_0$, problem (64) is equivalent to

$$\begin{cases} z' = \mathbf{\Lambda}z, & t > 0, \\ z(0) = z_0, \end{cases}$$

namely, to the n independent scalar problems in the components z_p , $1 \leq p \leq n$, of the vector z

$$\begin{cases} z'_p = \lambda_p z_p, & t > 0, \\ z_p(0) = z_{0p}. \end{cases}$$

The solutions of the n independent scalar problems

$$\begin{cases} z'_p = \lambda_p z_p, & t > 0, \\ z_p(0) = z_{0p}, \end{cases} \quad (65)$$

are given by

$$z_p(t) = e^{\lambda_p t} z_{0p}$$

and since by assumption $\operatorname{Re} \lambda_p < 0$, one has

$$|z_p(t)| \rightarrow 0 \quad \text{for } t \rightarrow +\infty \quad \text{and for each } p.$$

As a consequence, for $t \rightarrow +\infty$ one has $\|\mathbf{z}(t)\| \rightarrow 0$, hence, also $\|\mathbf{u}(t) - \mathbf{u}_\infty\| = \|\mathbf{v}(t)\| = \|\mathbf{W}\mathbf{z}(t)\| \rightarrow 0$, which means, as anticipated, $\mathbf{u}(t) \rightarrow \mathbf{u}_\infty$.

Thus, we are led to consider the generic scalar problem introduced above, i.e.,

$$\begin{cases} z' = \lambda z, & t > 0, \\ z(0) = z_0, \end{cases} \quad (66)$$

and to ask ourselves under which conditions a time-advancing scheme applied to such a problem produces discrete solutions z^k which stay bounded as $k \rightarrow \infty$.

The asymptotic stability of the Explicit Euler scheme (I)

The explicit Euler scheme applied to the equation $z' = \lambda z$ yields

$$z^{k+1} = z^k + \Delta t \lambda z^k = (1 + \Delta t \lambda) z^k, \quad k \geq 0,$$

hence by recursion, and keeping into account the initial condition, one gets the explicit expression

$$z^k = (1 + \Delta t \lambda)^k z_0,$$

where the symbol k means raising the basis to the k -th power.

The condition $|z^k| \rightarrow 0$ as $k \rightarrow +\infty$ is then equivalent to the condition

$$|1 + \Delta t \lambda| < 1.$$

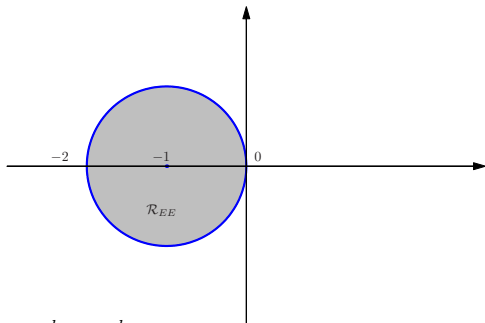
On the other hand, if instead one has $|1 + \Delta t \lambda| = 1$, then the solution stays bounded as $k \rightarrow +\infty$, although it does not tend to 0; indeed, one has $|z^k| = |z_0|$ for each $k \geq 0$. At last, if $|1 + \Delta t \lambda| > 1$, then necessarily one has $|z^k| \rightarrow +\infty$ as $k \rightarrow +\infty$.

In conclusion, setting $\alpha = \Delta t \lambda \in \mathbb{C}$, boundedness of the discrete solutions is equivalent to

$$|1 + \alpha| \leq 1.$$

The asymptotic stability of the Explicit Euler scheme (II)

The inequality $|1 + \alpha| \leq 1$ defines, in the complex plane, a circle of center -1 and radius 1 . This region is termed the *region of asymptotic stability* of the explicit Euler scheme, and denoted by \mathcal{R}_{EE} . Its internal part will be denoted by $\text{int}(\mathcal{R}_{EE})$.



Thus, if we want all the components of $\mathbf{z}^k = (z_p^k)$ to decay to 0, we are forced to choose Δt in such a way that the condition

$$\Delta t \lambda_p \in \text{int}(\mathcal{R}_{EE}) \quad \text{for each } p = 1, \dots, n, \quad (67)$$

is satisfied; this is always possible, thanks to the assumption $\text{Re} \lambda_p < 0$ on the eigenvalues of \mathcal{A} .

The explicit Euler scheme is therefore **conditionally asymptotically stable**, i.e., it is necessary to satisfy this condition on Δt , termed the **asymptotic stability condition**, in order to get the desired behaviour of the discrete solutions. It may occur, however, that such a condition is overly restrictive in practice.

The asymptotic stability of the Explicit Euler scheme (III)

In the important case where the matrix \mathcal{A} is symmetric, and consequently all its eigenvalues are real and negative, we have $\Delta t \lambda_p \in \text{int}(\mathcal{R}_{EE})$ if and only if

$$-2 < \Delta t \lambda_p < 0 ;$$

the second inequality is always fulfilled, whereas the first one is equivalent to

$$\Delta t < \frac{2}{|\lambda_p|} .$$

Thus, the asymptotic stability condition (67) becomes

$$\Delta t < \frac{2}{\max_p |\lambda_p|} . \quad (68)$$

If the matrix \mathcal{A} has eigenvalues with orders of magnitude quite different from each other (i.e., if it is *ill-conditioned* - in the language of systems of differential equations one says that the system is *stiff*), we are obliged to advance with a time-step dictated by the eigenvalue of largest absolute value, even if the actual behaviour of the exact solution would not require such a restriction.

(See next example.)

The asymptotic stability of the Implicit Euler scheme (I)

The implicit Euler scheme applied to the equation $z' = \lambda z$ yields

$$z^{k+1} = z^k + \Delta t \lambda z^{k+1}, \quad k \geq 0,$$

i.e.,

$$(1 - \Delta t \lambda) z^{k+1} = z^k, \quad k \geq 0.$$

Thus, the scheme generates the sequence

$$z^k = \left(\frac{1}{1 - \Delta t \lambda} \right)^{\wedge k} z_0, \quad k \geq 0,$$

which converges to 0 as $k \rightarrow +\infty$ if and only if the quantity in parenthesis is smaller than 1 in absolute value.

Setting $\alpha = \Delta t \lambda$, this is equivalent to the condition

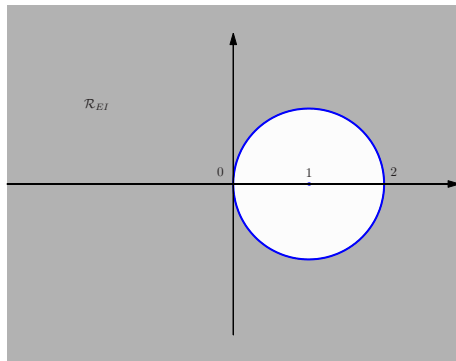
$$\frac{1}{|1 - \alpha|} < 1,$$

i.e.,

$$|1 - \alpha| > 1.$$

The asymptotic stability of the Implicit Euler scheme (II)

Therefore, the *region of asymptotic stability* \mathcal{R}_{EI} of the implicit Euler scheme is defined as the set of all complex numbers α such that $|1 - \alpha| \geq 1$.



Since the real part of $\Delta t \lambda_p$ is strictly negative, we have

$$\Delta t \lambda_p \in \text{int}(\mathcal{R}_{EI}) \quad \text{for each } p = 1, \dots, n \quad \text{and whichever } \Delta t > 0 \text{ is.}$$

Thus, there are no restrictions on the choice of the time step. We conclude that the implicit Euler scheme is **unconditionally asymptotically stable**.

The asymptotic stability of the trapezoidal scheme (I)

The trapezoidal (or Crank-Nicolson) scheme applied to the equation $z' = \lambda z$ yields

$$z^{k+1} = z^k + \Delta t \lambda \left(\frac{1}{2} z^k + \frac{1}{2} z^{k+1} \right), \quad k \geq 0,$$

namely,

$$\left(1 - \frac{\Delta t}{2} \lambda\right) z^{k+1} = \left(1 + \frac{\Delta t}{2} \lambda\right) z^k, \quad k \geq 0,$$

whose exact solution is given by

$$z^k = \left(\frac{1 + \frac{\Delta t}{2} \lambda}{1 - \frac{\Delta t}{2} \lambda} \right)^{\wedge k} z_0, \quad k \geq 0.$$

Setting $\alpha = \Delta t \lambda$, such a solution tends to 0 as $k \rightarrow +\infty$ if and only if one has

$$\left| 1 + \frac{\alpha}{2} \right| < \left| 1 - \frac{\alpha}{2} \right|.$$

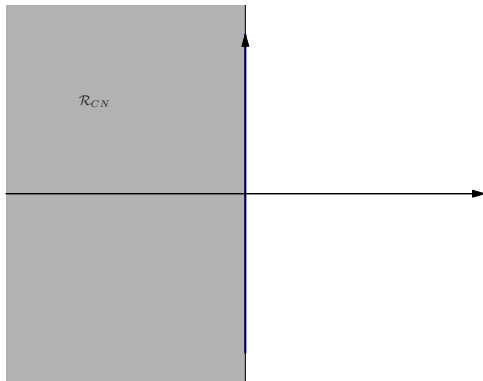
Now, it is easily seen that if α is any complex number, the inequality

$$\left| 1 + \frac{\alpha}{2} \right| \leq \left| 1 - \frac{\alpha}{2} \right|$$

is equivalent to $\operatorname{Re} \alpha \leq 0$.

The asymptotic stability of the trapezoidal scheme (II)

Therefore, the *region of asymptotic stability* \mathcal{R}_{CN} of the trapezoidal scheme is give by the half-plane on the left of the imaginary axis; such a region coincides indeed with the region in which all the eigenvalues of the matrix \mathcal{A} have to fall in order for the system (64) to have all its solutions bounded as $t \rightarrow +\infty$.



In this case, too, one has

$$\Delta t \lambda_p \in \text{int}(\mathcal{R}_{CN}) \quad \text{for each } p = 1, \dots, n \quad \text{and whichever } \Delta t > 0 \text{ is ;}$$

hence, the trapezoidal scheme is **unconditionally asymptotically stable**.

The asymptotic stability of other schemes

- The Runge Kutta schemes presented above are all **conditionally asymptotically stable**. Their regions of asymptotic stability are larger than the one for the Explicit Euler scheme.
(See plots in the Notes of the Course).
- The BDF schemes presented above have regions of asymptotic stability which contain a corner in the half-plane $\operatorname{Re} \alpha \leq 0$, with center at the origin and symmetrically placed around the real negative semi-axis.
Therefore, such methods are particularly suited for the discretization of *stiff* differential systems, as they turn out to be **unconditionally asymptotically stable** if the matrix of the system has all real and negative eigenvalues.

An example of “stiff” system (I)

Consider the problem

$$\begin{cases} \mathbf{v}' = \mathcal{A}\mathbf{v}, & t > 0, \\ \mathbf{v}(0) = \mathbf{v}_0, \end{cases}$$

with

$$\mathcal{A} = - \begin{pmatrix} \frac{1001}{2} & \frac{999}{2} \\ \frac{999}{2} & \frac{1001}{2} \end{pmatrix} \quad \text{and} \quad \mathbf{v}_0 = \begin{pmatrix} 2 \\ 0 \end{pmatrix}.$$

The eigenvalues of the matrix are $\lambda_1 = -1$ and $\lambda_2 = -1000$, with corresponding eigenvectors given by $\mathbf{w}_1 = (1, 1)^T$ and $\mathbf{w}_2 = (1, -1)^T$. Thus, the exact solution can be written as

$$\begin{pmatrix} v_1(t) \\ v_2(t) \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} z_1(t) \\ z_2(t) \end{pmatrix},$$

with $z_1(t) = e^{-t}$ and $z_2(t) = e^{-1000t}$, i.e., we obtain

$$\begin{cases} v_1(t) = e^{-t} + e^{-1000t}, \\ v_2(t) = e^{-t} - e^{-1000t}. \end{cases}$$

We thus have a very small transient, of the order of 1/1000 seconds, followed by a much slower evolution, in the time scale of seconds.

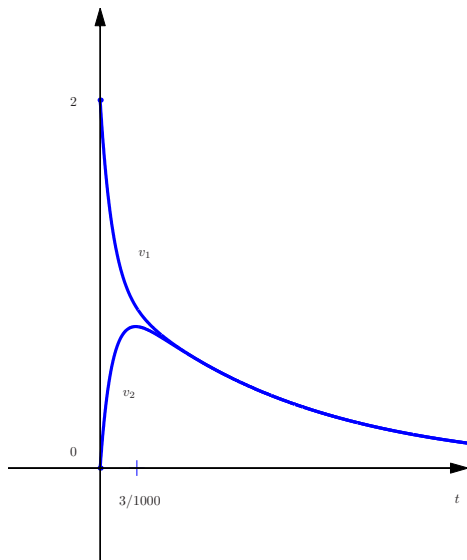
An example of “stiff” system (II)

If we advance in time with the explicit Euler scheme, we are forced to respect the asymptotic stability condition (68), namely

$$\Delta t < \frac{2}{\max(|-1|, |-1000|)} = \frac{1}{500} ,$$

not only during the transient period, but throughout the whole simulation, otherwise the computation quickly yields an overflow.

Using instead the implicit Euler scheme, or the more accurate trapezoidal scheme, allows us to vary the time-advancing step, choosing it in the order of $1/500$ seconds during the transient period, but subsequently letting it to increase as the solution goes towards a steady state.



Let us now discretize in time the problem

$$\begin{cases} \mathbf{B}\mathbf{u}' + \mathbf{A}\mathbf{u} = \mathbf{f}(t) , & 0 < t \leq T , \\ \mathbf{u}(0) = \mathbf{u}_0 . \end{cases}$$

using one of the previously presented schemes, namely the Explicit Euler scheme, the Implicit Euler scheme and the Trapezoidal (or Crank-Nicolson) scheme.

To this end, it is convenient to write the differential system in the equivalent (normal) form

$$\mathbf{u}' = \mathbf{F}(\mathbf{u}, t) = -\mathbf{B}^{-1}\mathbf{A}\mathbf{u} + \mathbf{B}^{-1}\mathbf{f}(t) \quad (= \mathcal{A}\mathbf{u} + \mathbf{b}(t)) .$$

Let us stress that such a transformation is only useful at the conceptual level, in order to apply the abstract time-advancing schemes in our specific setting. At the implementation level, the multiplication by the inverse of \mathbf{B} is almost invariably not efficient, hence, to be avoided; it is surely preferable to leave such matrix on the left-hand side of the equation, as shown in the sequel.

- The *Explicit Euler* scheme yields

$$\mathbf{u}^{k+1} = \mathbf{u}^k - \Delta t \mathbf{B}^{-1} \mathbf{A} \mathbf{u}^k + \Delta t \mathbf{B}^{-1} \mathbf{f}(t_k), \quad k \geq 0,$$

which we re-formulate as

$$\mathbf{B} \mathbf{u}^{k+1} = (\mathbf{B} - \Delta t \mathbf{A}) \mathbf{u}^k + \Delta t \mathbf{f}(t_k), \quad k \geq 0.$$

In the latter form, the scheme is not explicit, since at each iteration it requires the solution of a linear system with matrix \mathbf{B} . However, if we approximate \mathbf{B} with the “lumped” mass matrix $\tilde{\mathbf{B}}$, which is diagonal, then the cost of computing \mathbf{u}^{k+1} is essentially comparable to that of an explicit method.

In the sequel, we will see that the *asymptotic stability condition* for the Explicit Euler scheme becomes, in this case,

$$\Delta t \leq C h^2$$

(with C proportional to $\mu = \kappa/c$); this poses a restriction on the choice of the time-step, which often is not acceptable in practice.

The scheme is first-order accurate in time.

- The *Implicit Euler* scheme yields

$$\mathbf{u}^{k+1} = \mathbf{u}^k - \Delta t \mathbf{B}^{-1} \mathbf{A} \mathbf{u}^{k+1} + \Delta t \mathbf{B}^{-1} \mathbf{f}(t_{k+1}) , \quad k \geq 0 ,$$

which we re-formulate as

$$(\mathbf{B} + \Delta t \mathbf{A}) \mathbf{u}^{k+1} = \mathbf{B} \mathbf{u}^k + \Delta t \mathbf{f}(t_{k+1}) , \quad k \geq 0 .$$

In this case, one has to solve a linear system at each iteration, with matrix $\mathbf{B} + \Delta t \mathbf{A}$ symmetric and positive definite (since it is the sum of two matrices with these properties).

The gain over Explicit Euler is that the scheme is unconditionally stable.

The scheme is again first-order accurate in time.

- The *Trapezoidal* (or *Crank-Nicolson*) scheme yields, for our problem, the update

$$(B + \frac{\Delta t}{2}A)u^{k+1} = (B - \frac{\Delta t}{2}A)u^k + \frac{\Delta t}{2}(f(t_k) + f(t_{k+1})), \quad k \geq 0.$$

The scheme is again unconditionally stable, but now second-order accurate in time.

Hence, with a moderate increment in the cost of computing the right-hand side as compared to the Implicit Euler scheme, one significantly gains in precision, as the order of the scheme is increased by one.

If we keep the time-step Δt constant and we use a direct method to solve the linear system, then it is convenient to factorize the matrix

$$B + \frac{\Delta t}{2}A$$

(for instance by computing its Choleski factorization) once and for all at the beginning of the time loop: at each time instant, we will only have to perform a forward-backward substitution.

On the contrary, if one uses an iterative method, such as Conjugate Gradient, then the computed solution u^k at the previous time instant will provide the initial guess for the new iteration.

The asymptotic stability condition for the Explicit Euler scheme

Recall that the asymptotic stability condition for the Explicit Euler scheme

$$\Delta t \theta_h \in \text{int}(\mathcal{R}_{EE}) \quad \text{for any eigenvalue } \theta_h \text{ of the matrix } \mathcal{A} = -\mathbf{B}^{-1}\mathbf{A}.$$

Let us observe that θ_h is an eigenvalue of $\mathcal{A} = -\mathbf{B}^{-1}\mathbf{A}$ if and only if there exists a non-zero vector \mathbf{w} such that

$$\mathcal{A}\mathbf{w} = \theta_h \mathbf{w} ,$$

i.e.,

$$\mathbf{B}^{-1}\mathbf{A}\mathbf{w} = -\theta_h \mathbf{w} ,$$

namely, if and only if $\theta_h = -\lambda_h$, with λ_h eigenvalue of the matrix $\mathbf{B}^{-1}\mathbf{A}$.

In turns, this is true if and only if λ_h is a solution of the *generalized eigenvalue problem*

$$\mathbf{A}\mathbf{w} = \lambda_h \mathbf{B}\mathbf{w} .$$

We have seen that such a problem arises in the discretization of the *modal analysis problem* for an elastic membrane.

From this relation we derive

$$\mathbf{w}^T \mathbf{A} \mathbf{w} = \lambda_h \mathbf{w}^T \mathbf{B} \mathbf{w} ,$$

which provides the expression of λ_h as a **Rayleigh quotient**

$$\lambda_h = \frac{\mathbf{w}^T \mathbf{A} \mathbf{w}}{\mathbf{w}^T \mathbf{B} \mathbf{w}} .$$

This shows that all the eigenvalues λ_h are *real and strictly positive* (since the matrices \mathbf{A} and \mathbf{B} are both symmetric and positive definite), hence, all the eigenvalues θ_h of the matrix \mathcal{A} are real and negative.

Consequently, the asymptotic stability condition of the Explicit Euler scheme

$$\Delta t \leq \frac{2}{\max |\theta_h|}$$

is equivalent to

$$\Delta t \leq \frac{2}{\max \lambda_h} .$$

Since we have stated above that $\max \lambda_h \sim ch^{-2}$, the asymptotic stability condition becomes

$$\Delta t \leq C h^2 ,$$

as anticipated.

Time discretization of the elastic model: the Newmark scheme

In order to discretize in time the second-order system

$$\mathbf{u}'' = \mathbf{F}(\mathbf{u}, t) = -\mathbf{B}^{-1}\mathbf{A}\mathbf{u} + \mathbf{B}^{-1}\mathbf{f}(t) ,$$

it is customary to resort to the *Newmark scheme*.

It is written as

$$\mathbf{u}^{k+1} - 2\mathbf{u}^k + \mathbf{u}^{k-1} = \Delta t^2 (\beta \mathbf{F}(\mathbf{u}^{k+1}, t_{k+1}) + (1 - 2\beta) \mathbf{F}(\mathbf{u}^k, t_k) + \beta \mathbf{F}(\mathbf{u}^{k-1}, t_{k-1})) ,$$

where $\beta \geq 0$ is a parameter to be chosen. Equivalently, this can be written as

$$\begin{aligned} \mathbf{B}(\mathbf{u}^{k+1} - 2\mathbf{u}^k + \mathbf{u}^{k-1}) + \Delta t^2 \mathbf{A}(\beta \mathbf{u}^{k+1} + (1 - 2\beta) \mathbf{u}^k + \beta \mathbf{u}^{k-1}) = \\ = \Delta t^2 (\beta \mathbf{f}(t_{k+1}) + (1 - 2\beta) \mathbf{f}(t_k) + \beta \mathbf{f}(t_{k-1})) . \end{aligned}$$

The scheme defines the time approximations starting from $k = 1$. Hence, it requires two initial values, \mathbf{u}^0 and \mathbf{u}^1 . The first value is obviously chosen equal to \mathbf{u}_0 , whereas the Taylor expansion of $\mathbf{u}(t)$ at the origin, $\mathbf{u}(t_1) = \mathbf{u}(0) + \Delta t \mathbf{u}'(0) + O(\Delta t^2)$, suggests the choice $\mathbf{u}^1 = \mathbf{u}_0 + \Delta t \mathbf{v}_0$.

Let us remark that if $\beta = 0$ and if the mass matrix \mathbf{B} is replaced by the lumped mass matrix $\hat{\mathbf{B}}$, then the previous scheme is explicit. In any other case, the scheme is implicit, and at each time instant \mathbf{u}^{k+1} is determined by solving a linear system with matrix $\mathbf{B} + \Delta t^2 \beta \mathbf{A}$.

As far as accuracy is concerned, the scheme turns out to be *second-order accurate* for each value of β .

On the other hand, as far as asymptotic stability is concerned, the scheme turns out to be

- *unconditionally stable* if $\beta \geq \frac{1}{4}$.
- only *conditionally stable* if $\beta < \frac{1}{4}$.

In the latter case, the asymptotic stability condition is

$$\Delta t^2 \lambda_{h,max} < \frac{4}{1 - 4\beta} ,$$

where $\lambda_{h,max}$ denotes the maximum eigenvalue of the matrix $\mathbf{B}^{-1}\mathbf{A}$. We have already noticed that the order of magnitude of this eigenvalue is $O(h^{-2})$.

As a consequence, the previous conditions enforces a restriction of the type

$$\Delta t \leq C h$$

on the choice of the time step. Such a condition is by far less stringent than the stability condition for a conditionally stable scheme applied to the heat equation, which is of the type $\Delta t \leq C h^2$.

In most cases, such a condition turns out to be fully acceptable, since it is anyway required by the need of guaranteeing enough time accuracy on the discretization.