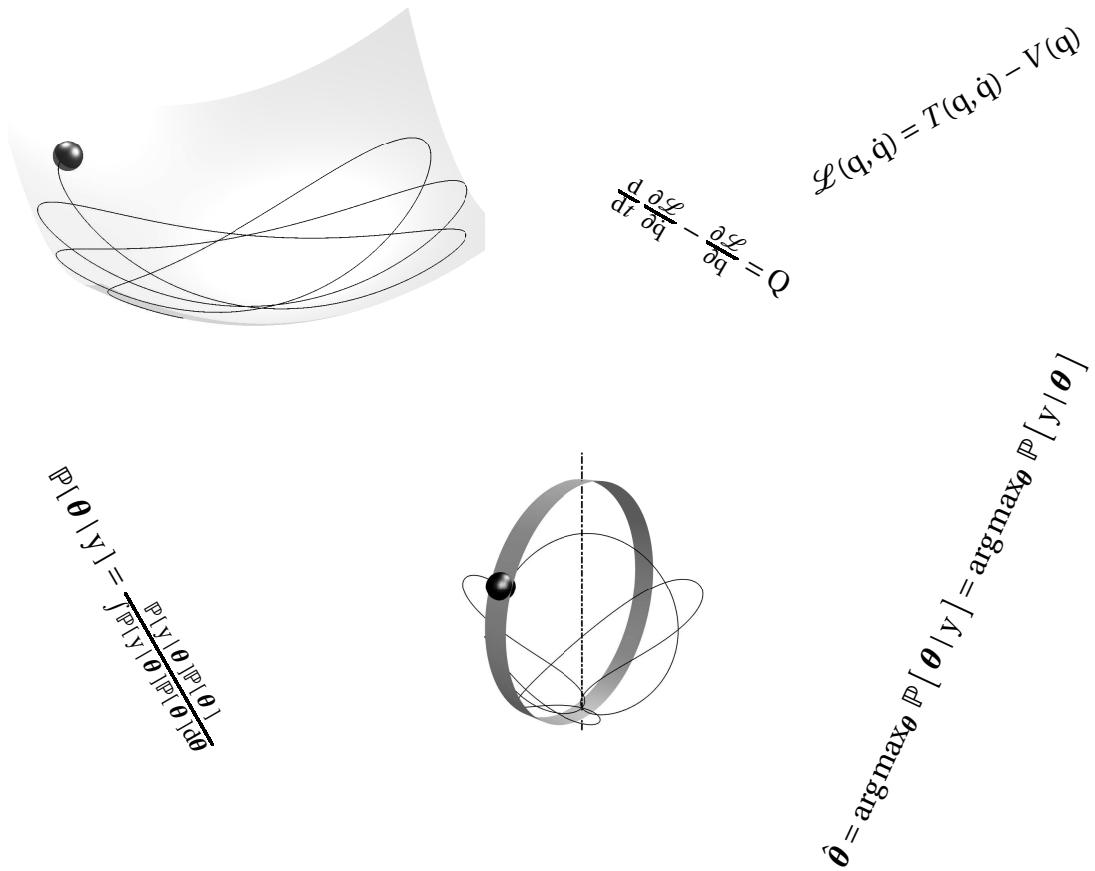

Modelling And Simulation

Lecture notes for the NTNU/ITK course TTK4130

Sebastien Gros

Edition: January 19, 2021



CONTENTS

1 Notation & Background material	4
1.1 Norms and scalar products	4
1.2 Probability	5
1.3 Some basics from Linear Algebra	6
1.4 Some basics from Calculus	6
1.5 Ordinary Differential Equations (ODE)	8
1.6 Linear ODE	9
1.7 Linearization	9
1.8 Solution of ODEs	11
1.9 Solution to LTI ODEs	13
1.10 Z transform	14
2 Lagrange mechanics	16
2.1 Kinetic Energy	16
2.2 Potential Energy	20
2.3 Lagrange Equation	22
2.4 External forces	37
2.5 Constrained Lagrange Mechanics	40
2.5.1 Handling Models from Constrained Lagrange	46
2.6 Consistency conditions	51
2.7 Constraints drift	52
3 Newton Method	54
3.1 Basic idea of Newton method	54
3.2 Convergence of the Newton method	56
3.2.1 Convergence rate	56
3.2.2 Reduced Newton steps	60
3.3 Implicit Function Theorem	62
3.4 Jacobian Approximation	64
3.5 Newton Methods for Unconstrained Optimization	65
3.5.1 Gauss-Newton Hessian approximation	67
3.5.2 Convex optimization	68
3.6 Summary	68
4 System Identification (SYSID)	70
4.1 Bayesian reasoning & Bayesian Estimation	74
4.2 Maximum Likelihood Estimation	76
4.3 Least-Squares Estimation	81
4.4 Solution to the least-squares problem	85
4.5 Covariance, Bias & Consistency	86
4.6 Estimation for Linear Dynamic Systems	89
4.6.1 Estimation of Finite Impulse Response models	90

4.6.2	Auto-regressive models - Simulation error	91
4.6.3	Auto-regressive models - Prediction error	92
4.6.4	Prediction error & Noise model	93
5	Differential Algebraic Equations (DAEs)	97
5.1	What are DAEs?	97
5.2	Different forms of DAEs	100
5.3	Tikhonov Theorem	102
5.4	Differential Index of DAEs	107
5.5	Connection to Lagrange mechanics	112
5.6	Index reduction	113
6	Explicit Integration Methods - Runge-Kutta	116
6.1	Explicit Euler	116
6.1.1	Accuracy of the explicit Euler method	117
6.1.2	Stability of the explicit Euler method	120
6.2	Explicit Runge-Kutta 2 methods	122
6.3	General RK methods	126
6.3.1	Butcher tableau	126
6.3.2	The RK4 method	128
6.3.3	Stages, order & efficiency of explicit RK methods	129
6.4	Stability of explicit RK schemes	133
6.4.1	Stiff systems	134
6.5	Error control & Adaptive integrators	136
7	Implicit Integration Methods - Runge-Kutta	139
7.1	Implicit Euler method	140
7.1.1	Stability of the implicit Euler method	141
7.2	Implicit Runge-Kutta methods	142
7.3	Stability & order of IRK methods	144
7.4	Collocation methods	147
7.4.1	Polynomial interpolation	147
7.4.2	Interpolation of the trajectories	149
7.5	RK methods for implicit ODEs	152
7.6	RK methods for implicit DAEs	154
7.6.1	RK method for semi-explicit DAE models	156
8	Sensitivity of Simulations	158
8.1	Variational approach	158
8.2	Algorithmic Differentiation of the explicit Euler scheme	160
8.3	Algorithmic Differentiation of explicit Runge-Kutta methods	161
8.4	Sensitivity of implicit Runge-Kutta steps	163
8.5	Sensitivity with respect to inputs	165

1 NOTATION & BACKGROUND MATERIAL

Most often, we will note vectors in the bold math roman format, i.e. $\mathbf{a} \in \mathbb{R}^n$. Scalars will be noted in the standard math format, i.e. $a \in \mathbb{R}$. Matrices will be noted using capital letters, i.e. $A \in \mathbb{R}^{n \times m}$. There will be a few exceptions to these rules.

We will use subscript to denote the elements of vectors and matrices, i.e. $a_i \in \mathbb{R}$ will be the i :th element of vector \mathbf{a} , and A_{ij} the element in the i :th row and j :th column of matrix A . Generic functions will obey the same notation rules. We will reserve the capital letter I for the identity matrix.

1.1 NORMS AND SCALAR PRODUCTS

In this course we will often use the notion of norm, noted as $\|\cdot\|$. A norm associates to a vector a positive scalar. It is often used to simply “measure the length” of a vector, but it has a deeper meaning than this. An example of norm is:

$$\|\mathbf{x}\| = \sqrt{\mathbf{x}^\top \mathbf{x}} \quad (1.1)$$

But any operator ρ taking an element of the vector space under consideration into \mathbb{R}_+ and satisfying:

1. $\rho(\mathbf{x} + \mathbf{y}) \leq \rho(\mathbf{x}) + \rho(\mathbf{y})$
2. $\rho(a\mathbf{x}) = |a|\rho(\mathbf{x})$ for any $a \in \mathbb{R}$
3. $\rho(\mathbf{x}) = 0$ entails that $\mathbf{x} = 0$

There are several norms that differ from (1.1), e.g. the 1-norm

$$\|\mathbf{x}\|_1 = \sum_{k=1}^n |\mathbf{x}_k| \quad (1.2)$$

and the infinity norm

$$\|\mathbf{x}\|_\infty = \max_{k=1,\dots,n} |\mathbf{x}_k| \quad (1.3)$$

Note that since all norms are equivalent (i.e. boundedness in one norm implies boundedness in the others), we often omit (when irrelevant) which type of norm we use.

We will use the notation

$$\langle \cdot, \cdot \rangle \quad (1.4)$$

to denote scalar products on any vector space. On \mathbb{R}^n , the scalar product between two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ reads as:

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y} \quad (1.5)$$

1.2 PROBABILITY

We will use probabilities in one of the sections of this course. Let us review some basic principles and notations here. We will note

$$\mathbb{P}[A] \quad (1.6)$$

the probability that statement A is true. We will abuse this notation and use it for both probabilities and probability distributions. When \mathbb{P} is used for a probability distribution, then $\mathbb{P}[x]$ ought to be understood as “probability that variable x takes a specific value”, and $\mathbb{P}[x]$ will be a function of x . E.g.

$$\mathbb{P}[x] = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}x^2} \quad (1.7)$$

is a normal, centred probability distribution describing the probability that a variable takes a specific value x . We will additionally need the notion of conditional probability (or conditional probability distribution), that we will note as:

$$\mathbb{P}[A|B] \quad (1.8)$$

meaning “probability that A is true knowing that B is true”. We will also use the same notation conditional probability distributions, i.e. $\mathbb{P}[x|y]$ will be the “probability that variable x takes a specific value knowing that variable y has taken the specific value”. $\mathbb{P}[x|y]$ will then be a function of x and y , e.g. it could be:

$$\mathbb{P}[x|y] = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-y)^2} \quad (1.9)$$

We will need the notion of expected value, noted and defined as:

$$\mathbb{E}[x] = \int x \mathbb{P}[x] \quad (1.10)$$

for continuous distributions. Note that $\mathbb{E}[x]$ is a number or vector of the same size as x . In practice, one can construe the expected value as “the average of variable x over many experiments”, i.e. if one were to draw the random variable x many times and average the result, one would get something close to the expected value. Similarly, we can define the conditional expected value as:

$$\mathbb{E}[x|y] = \int x \mathbb{P}[x|y] \quad (1.11)$$

In practice, one can construe the conditional expected value as “the average of variable x over many experiments provided that y takes a specific value”. Note that $\mathbb{E}[x|y]$ is a function of y alone, and returns a scalar if x is scalar or a vector of the size of x .

1.3 SOME BASICS FROM LINEAR ALGEBRA

Let us review here some basic principles of linear algebra that can be useful in this course.

- The eigenvalues λ and eigenvectors v satisfy the equation

$$Av = \lambda v \quad (1.12)$$

i.e. the eigenvectors are transformed into scaled versions of themselves by the application of matrix A . The eigenvalues can be computed by solving the polynomial equation

$$\det(A - \lambda I) = 0 \quad (1.13)$$

where \det is the matrix determinant.

- A matrix A is invertible (i.e. A^{-1} exists) iff
 - $\det(A) \neq 0$ or
 - All the eigenvalues of A are nonzero

An invertible matrix is said "full rank".

- A matrix A is said symmetric if $A = A^\top$
- A matrix A is said positive-definite if all its eigenvalues are real positive.
- A quadratic form is the operation:

$$x^\top Ax \in \mathbb{R} \quad (1.14)$$

where A is a square matrix, and x is a vector of adequate dimension. If A is a positive-definite matrix, then

$$x^\top Ax \geq 0 \quad (1.15)$$

for any x .

1.4 SOME BASICS FROM CALCULUS

We will use a lot of calculus in this course, and you need to be comfortable with it. We recall some basic calculus here.

- The Jacobian of a multi-variate function $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$

$$f(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_n(x) \end{bmatrix}, \quad x \in \mathbb{R}^m \quad (1.16)$$

is a $\mathbb{R}^{n \times m}$ matrix given by:

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_m} \\ \vdots & & \vdots \\ \frac{\partial f_n}{\partial x_1} & \dots & \frac{\partial f_n}{\partial x_m} \end{bmatrix} \quad (1.17)$$

- It will be convenient certain times to use the gradient operator ∇ , which is essentially the transpose of the Jacobian operator, i.e.

$$\nabla_x f = \frac{\partial f^\top}{\partial x} \quad (1.18)$$

The gradient is most often used for scalar functions $f \in \mathbb{R}$, but the notion of gradient can be readily applied to all functions.

- Chain ruling will be used a lot in this course. For a composition of functions $f(g(x))$, chain ruling reads as:

$$\frac{\partial f(g(x))}{\partial x} = \frac{\partial f(y)}{\partial y} \Big|_{y=g(x)} \frac{\partial g(x)}{\partial x} \quad (1.19)$$

Note that the order in the right-hand-side of the equation matters for multi-variate functions, as the Jacobians are matrices. This order is reversed under the gradient operator, i.e.

$$\nabla_x f(g(x)) = \nabla_x g(x) \nabla_y f(y) \Big|_{y=g(x)} \quad (1.20)$$

- We will need a few times to distinguish between total and partial derivatives. The notion can be tricky and even ambiguous sometimes, but let us recall here the basic principle. Consider a function:

$$f(x, y) \quad (1.21)$$

where $y = g(x)$. Then the partial derivative ignores that y is intrinsically a function of x , i.e.

$$\frac{\partial f(x, y)}{\partial x} \quad (1.22)$$

disregards this dependency, while the total derivatives takes it into account and does:

$$\frac{df(x, y)}{dx} = \frac{\partial f(x, y)}{\partial x} + \frac{\partial f(x, y)}{\partial y} \frac{\partial y}{\partial x} = \frac{\partial f(x, y)}{\partial x} + \frac{\partial f(x, y)}{\partial y} \frac{\partial g}{\partial x} \quad (1.23)$$

We will briefly review next some basic notions that will be useful in this course, and introduce some new material that you may have not seen before.

1.5 ORDINARY DIFFERENTIAL EQUATIONS (ODE)

Ordinary Differential Equations (ODEs) will play a crucial role in this course. ODEs are a way to relate the outputs y of a system to the inputs u acting on the system. Generally speaking, an ODE reads as

$$\varphi(y^{(m)}, y^{(m-1)}, \dots, y, u^{(m-1)}, u^{(m-2)}, \dots, u) = 0 \quad (1.24)$$

for some function φ and where we use the notation

$$y^{(k)} = \frac{d^k}{dt^k} y \quad (1.25)$$

a trivial example of such an equation is the motion of a mass m attached to a spring of constant K , and subject to an external force u and viscous friction of parameter ξ . The ODE describing such a system reads as:

$$\varphi = m\ddot{y} + \xi\dot{y} + Ky - u = 0 \quad (1.26)$$

Most often we prefer to work with ODEs in their state-space form rather than in the form (1.24). An ODE in the state-space is most often written as:

$$\dot{x} = f(x, u) \quad (1.27)$$

for some state $x \in \mathbb{R}^n$. E.g. the spring-mass example (1.26) in its state-space form reads as:

$$\dot{x} = \begin{bmatrix} x_2 \\ \frac{1}{m}(u - \xi\dot{y} - Ky) \end{bmatrix}, \quad x = \begin{bmatrix} y \\ \dot{y} \end{bmatrix} \quad (1.28)$$

A complete state-space model is then made of the dynamics (1.27) and an output function $y = h(x, u)$ telling us what we can "measure" or alternatively what we "care about" in the system. The full state-space model then reads as:

$$\dot{x} = f(x, u) \quad (1.29a)$$

$$y = h(x, u) \quad (1.29b)$$

We can make a few remarks here concerning (1.29).

- For $x \in \mathbb{R}^n$, n provides the order of the system (state-space dimension)
- Functions f, h are generally nonlinear, often smooth
- The model is said *time-invariant* if the functions f, h do not *explicitly* depend on time

We will see later in the course that it can be sometimes useful to treat the dynamics of a system via an implicit ODE, which has the generic form:

$$F(\dot{x}, x, u) = 0 \quad (1.30)$$

i.e. it delivers the derivatives of the state space not via a function f that can be evaluated directly, but implicitly via the equations (1.30).

1.6 LINEAR ODE

A very important class of models is the class of *linear* models, where functions f, h are linear. We can then rewrite (1.29) in the form:

$$\dot{x} = Ax + Bu \quad (1.31a)$$

$$y = Cx + Du \quad (1.31b)$$

for some matrices $A \in \mathbb{R}^{n_x \times n_x}$, $B \in \mathbb{R}^{n_x \times n_u}$, $C \in \mathbb{R}^{n_y \times n_x}$, $D \in \mathbb{R}^{n_y \times n_u}$. Linear systems have the superposition property, i.e. if $x_1(t)$ and $x_2(t)$ are solution of (1.31) for the input profiles $u_1(t), u_2(t)$ respectively, then $\alpha_1 x_1(t) + \alpha_2 x_2(t)$ is solution corresponding to the input profile $\alpha_1 u_1(t) + \alpha_2 u_2(t)$. Linear ODEs are "nice" to work with because one can treat them exploiting the powerful theorems provided by functional analysis, which - among other things - give us the Fourier and Laplace transforms.

When matrices A, B, C, D are fixed, then system (1.31) is said Linear Time Invariant (LTI). If (any of) the matrices $A(t), B(t), C(t), D(t)$ is functions of time, then system (1.31) is said Linear Time Varying (LTV).

1.7 LINEARIZATION

Linearization consists in forming *locally valid* linear approximations of nonlinear ODEs (1.29). The linearization of an ODE is nothing more than an application of the Taylor expansion, to a first order. I.e. the functions f, h are replaced by their first-order approximation with respect to all arguments. In that context, we will consider deviations in the states and inputs from a given reference, and establish how this deviation will evolve in time, to a first-order approximation. Let us build this step-by-step. Consider first the first-order expansion of function f . Consider a *reference* trajectory $x(t)$ solution of (1.29a) for a given *reference* input profile u and *reference* initial conditions x_0 . Consider then a deviation $\Delta u(t)$ and Δx_0 in the input profile and/or initial conditions, and consider $\Delta x(t)$ the resulting deviation in the ODE trajectories. We observe that:

$$\dot{x}(t) + \dot{\Delta x}(t) = f(x(t) + \Delta x(t), u(t) + \Delta u(t)), \quad \Delta x(0) = \Delta x_0 \quad (1.32a)$$

$$y(t) + \Delta y(t) = h(x(t) + \Delta x(t), u(t) + \Delta u(t)) \quad (1.32b)$$

must hold. For small deviations $\Delta u(t), \Delta x(t)$, we can form the first-order Taylor expansion of the right-hand-side of (1.32), and obtaining:

$$\dot{x}(t) + \dot{\Delta x}(t) = f(x(t), u(t)) + \frac{\partial f}{\partial x} \Big|_{x(t), u(t)} \Delta x(t) + \frac{\partial f}{\partial u} \Big|_{x(t), u(t)} \Delta u(t) + \mathcal{O}(\|\Delta x\|^2, \|\Delta u\|^2) \quad (1.33a)$$

$$y(t) + \Delta y(t) = h(x(t), u(t)) + \frac{\partial h}{\partial x} \Big|_{x(t), u(t)} \Delta x(t) + \frac{\partial h}{\partial u} \Big|_{x(t), u(t)} \Delta u(t) + \mathcal{O}(\|\Delta x\|^2, \|\Delta u\|^2) \quad (1.33b)$$

Since $\dot{x}(t) = f(x(t), u(t))$, we equivalently get:

$$\dot{\Delta x}(t) \approx A(t)\Delta x(t) + B(t)\Delta u(t) \quad (1.34a)$$

$$\Delta y(t) \approx C(t)\Delta x(t) + D(t)\Delta u(t) \quad (1.34b)$$

for the time-varying matrices:

$$A(t) = \frac{\partial f}{\partial x} \Big|_{x(t), u(t)}, \quad B(t) = \frac{\partial f}{\partial u} \Big|_{x(t), u(t)} \quad (1.35a)$$

$$C(t) = \frac{\partial h}{\partial x} \Big|_{x(t), u(t)}, \quad D(t) = \frac{\partial h}{\partial u} \Big|_{x(t), u(t)} \quad (1.35b)$$

A very commonly used special case of (1.34)-(1.35) is when the reference inputs and trajectories $u(t), x(t)$ are constant, in which case matrices $A(t), B(t), C(t), D(t)$ are also constant (because they come from Jacobians evaluated on fixed arguments). Such trajectories are in fact a stationary point of the dynamics, and can be found by solving:

$$f(x_0, u_0) = 0 \quad (1.36)$$

and using $x = x_0, u(t) = u_0$ such that $\dot{x} = 0$. We can then use these reference trajectory in (1.35).

Example Consider the nonlinear dynamics

$$\dot{x}_1 = (1 - x_2^2)x_1 - x_2 + u, \quad (1.37a)$$

$$\dot{x}_2 = x_1, \quad (1.37b)$$

with the output function:

$$h(x) = x_1 + x_2^3 \quad (1.38)$$

One can compute:

$$A = \begin{bmatrix} 1 - x_2(t)^2 & -2x_1x_2 - 1 \\ 1 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad C = [1 \ 3x_2(t)^2], \quad D = 0 \quad (1.39)$$

where $x(t)$ is a trajectory of the system. A steady state \bar{x} of the system is obtained by solving:

$$\dot{x}_1 = (1 - x_2^2)x_1 - x_2 + u = 0, \quad (1.40a)$$

$$\dot{x}_2 = x_1 = 0, L \quad (1.40b)$$

yielding

$$\bar{x} = \begin{bmatrix} 0 \\ \bar{u} \end{bmatrix} \quad (1.41)$$

for a given \bar{u} constant. A linearization at steady-state then reads as:

$$A = \begin{bmatrix} 1 - \bar{u}^2 & -1 \\ 1 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad C = [1 \ 3\bar{u}^2], \quad D = 0 \quad (1.42)$$

It is important to observe here that the linear system approximating the nonlinear dynamics depends in the input \bar{u} selected.

□

1.8 SOLUTION OF ODES

Let us now turn to a crucial question in this course. Consider an ODE in its state-space form

$$\dot{x} = f(x, u) \quad (1.43)$$

An important question to ask here is “provided the initial conditions $x(0) = x_0$, is there a trajectory $x(t)$ solution of (1.43), and is it unique?”.

Before treating that question, let us consider two simple examples showing that it is not trivial.

Example

- Consider the ODE:

$$\dot{x}(t) = x(t)^2, \quad x(0) = 1 \quad (1.44)$$

One can verify that (1.44) admits

$$x(t) = \frac{1}{1-t} \quad (1.45)$$

as solution. We illustrate the trajectory $x(t)$ in the next Figure. One can easily observe from (1.45) that the state $x(t)$ becomes arbitrarily large as t approaches 1. The trajectories, in fact, do not exist at and beyond $t = 1$.

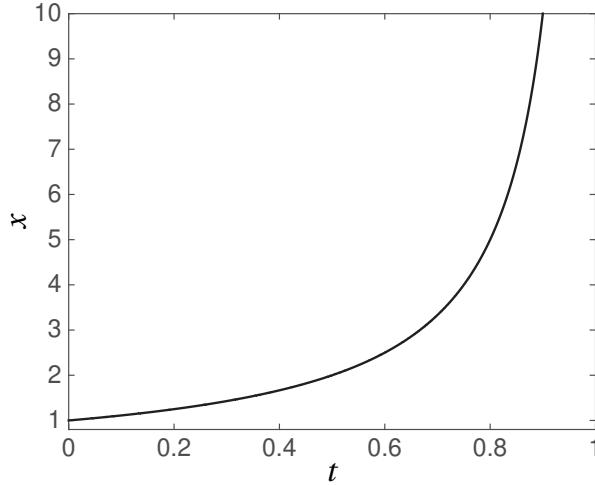


Figure 1.1: Illustration of the solution (1.45). The trajectories reach ∞ in finite time (when getting close to $t = 1$).

This simple example illustrates that the solution of an ODE can “explode” in finite time, and cease to exist beyond a limited time interval.

- Consider the ODE

$$\dot{x} = \sqrt{|x|}, \quad x(0) = 0 \quad (1.46)$$

One can verify that:

$$x(t) = \begin{cases} \frac{(t-t_0)^2}{4} & t \geq t_0 \\ 0 & t < t_0 \end{cases} \quad (1.47)$$

is an admissible solution for (1.46) for any $t_0 \geq 0$. The issue here is that it is an admissible solution regardless of t_0 . In other words, there are infinitely many trajectories $x(t)$ (for infinitely many choice of t_0) that are solution of (1.46). We illustrate these solutions in the next Figure.

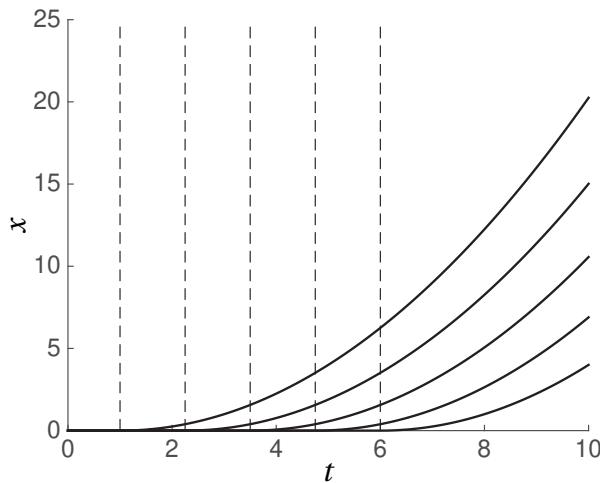


Figure 1.2: Illustration of the solutions (1.47) for different values of t_0 (dashed lines).

Let us introduce two theorems that unpack the two issues (existence and unicity) raised in the examples above. The first theorem deals with the existence of the ODE solution.

Theorem 1. *Consider the ODE*

$$\dot{x} = f(x) \quad (1.48)$$

where f is continuous. Then if¹

$$\|f(x) - f(y)\| \leq c \cdot \|x - y\|, \quad \forall x, y \quad (1.49)$$

for some finite constant c , then the solution of (1.48) exists and is unique for all t .

¹the following property is labelled “Lipschitz continuity”

Let us return to the first example above. One can observe that

$$f(x) = x^2 \quad (1.50)$$

does not satisfy condition (1.49). Indeed, in this case,

$$\|f(x) - f(y)\| = x^2 - y^2 = (x + y)(x - y) \quad (1.51)$$

we can observe that the term $x + y$ can be arbitrarily large, such that we cannot find a constant c for which (1.49) holds.

Unfortunately, property (1.49) can be difficult to verify for non-trivial ODEs. Let us introduce a second theorem that makes our life easier.

Theorem 2. *Consider the ODE*

$$\dot{x} = f(x) \quad (1.52)$$

Then if f is continuously differentiable (i.e. the Jacobian $\frac{\partial f}{\partial x}$ exists and is continuous), then the solution to the ODE exists and is unique on some time interval.

Let us return to the second example above. One can observe that

$$\frac{\partial f}{\partial x} = \frac{1}{2\sqrt{x}} \quad (1.53)$$

such that f is not differentiable at $x = 0$. This results in the non-uniqueness of the solution of (1.46). Note that the first theorem does not apply to the ODE (1.46) as function f has “infinite slopes” at $x = 0$, such that condition (1.49) does not hold around $x = 0$. Note that the second theorem is a “weaker” version of the first one, and requires weaker conditions.

Before closing, let us consider the case of linear ODEs. We observe that the theorems above readily apply to any linear ODE as linear functions readily satisfy (1.49), and are continuously differentiable. That is, linear ODEs always have a unique solution over the time interval $[0, \infty[$. However, the solution may be unbounded, i.e. it can grow forever if the ODE is unstable, but there is not specific time where the solution becomes infinite (unlike (1.44)).

1.9 SOLUTION TO LTI ODES

Consider the LTI ODE:

$$\dot{x} = Ax + Bu, \quad x(t) = x_0 \quad (1.54)$$

One can verify that the solution to (1.54) reads as:

$$x(t) = e^{At}x_0 + \int_0^t e^{A(t-\tau)}Bu(\tau) d\tau \quad (1.55)$$

We ought to recall here the definition of a “matrix exponential” like e^{At} . Here we simply extend the series corresponding to the exponential function to matrices. More specifically, similarly to:

$$e^t = 1 + t + \frac{t^2}{2!} + \dots = \sum_{k=0}^{\infty} \frac{t^k}{k!} \quad (1.56)$$

we define:

$$e^{At} = I + At + \frac{(At)^2}{2!} + \dots = \sum_{k=0}^{\infty} \frac{(At)^k}{k!} \quad (1.57)$$

Note that here the power of a matrix A^k does not correspond to taking the power of the matrix entries, but rather to multiplying the matrix by itself k times. We additionally observe that the differentiation rule:

$$\frac{d}{dt} e^{at} = ae^{at} \quad (1.58)$$

becomes for the matrix exponential function:

$$\frac{d}{dt} e^{At} = Ae^{At} \quad (1.59)$$

Note that the matrix exponential function is “expm.m” in Matlab, and does not (necessarily) correspond to the classic exponential function “exp.m”.

We can conclude here that LTI ODEs do not need to be “simulated” as one can compute their solution explicitly from (1.55). For nonlinear ODEs (and some LTV ODEs), however, the best way of building their trajectories is by using computer-based simulations.

1.10 Z TRANSFORM

In one part of this course we will briefly use the Z -transform. We briefly recall here its principle. The Z -transform applies to LTI **discrete**-time systems, i.e. systems of e.g. the form:

$$y_k = \sum_{i=0}^{N_b} b_i u_{k-i} + \sum_{i=1}^{N_a} a_i y_{k-i} \quad (1.60)$$

We observe that the system output y is here described as a discrete sequence y_0, \dots, y_∞ . Similarly to the Laplace transform the Z -transform allows one to treat dynamics in the form (1.60) in the “polynomial world”.

The Z -transform of a signal y_k is given by:

$$Y(z) = Z(y) = \sum_{k=0}^{\infty} y_k z^{-k} \quad (1.61)$$

The Z -transform has a number of useful properties (such as e.g. linearity), among which the time-shift property:

$$Z(y_{k-i}) = z^{-i} Z(y) \quad (1.62)$$

allows to write the Z -transform of the dynamics (1.60) as:

$$Y(z) = \sum_{i=0}^{N_b} b_i z^{-i} U(z) + \sum_{i=1}^{N_a} a_i z^{-i} Y(z) \quad (1.63)$$

i.e.

$$Y(z) = \frac{\sum_{i=0}^{N_b} b_i z^{-i}}{1 - \sum_{i=1}^{N_a} a_i z^{-i}} U(z) \quad (1.64)$$

2 LAGRANGE MECHANICS

In this Chapter we will study the basics of Lagrange mechanics. Lagrange mechanics is a powerful tool to build mathematical models for complex mechanical systems, and is used in many applications. Lagrange mechanics

- Allows to describe arbitrarily complex mechanical system, including relative moving parts, accelerated frames, gyroscopic and centrifugal effects without hassle.
- Often allows to build simple model (in terms of complexity of the resulting model equations) even for mechanical systems that are normally described by very complex models.

We define the generalized coordinates of a mechanical system as a vector of time-varying "coordinates" $\mathbf{q}(t) \in \mathbb{R}^{n_c}$ that must be able to describe the "configuration" of the system at a given time t . One can construe these coordinates as forming a "snapshot" of the system. The snapshot does not tell us how the system will evolve in time, but it can tell us in what configuration the system is at a given time. The generalized coordinates typically gather positions and angles, but can also include more abstract representations of the system. For the sake of simplicity, in the following we will omit the explicit dependence of the generalized coordinates and other time-varying objects in the notation, i.e. we will e.g. note the time-varying generalized coordinates $\mathbf{q}(t)$ simply as \mathbf{q} .

Lagrange mechanics is based on a description of the mechanical system in terms of energy. In order to build the Lagrange framework and models using the Lagrange equations, we need to compute the kinetic and potential energy functions of the system, conventionally labelled T and V , respectively. Building these objects is not necessarily straightforward, but it can be done fairly systematically, and with a very limited risk of error. Let us discuss the kinetic and potential energy functions next.

2.1 KINETIC ENERGY

We recall here that the kinetic energy of a "punctual mass" (i.e. of a dimensional particle of mass m) whose position with respect to a fixed (inertial) reference frame is given by $\mathbf{p} \in \mathbb{R}^D$ (where the number of dimensions can be $D = 1, 2, 3$) reads as:

$$T = \frac{1}{2}m\|\dot{\mathbf{p}}\|^2 = \frac{1}{2}m\dot{\mathbf{p}}^\top \dot{\mathbf{p}} \quad (2.1)$$

A few useful observation ought to be made here:

- The kinetic energy function is a *quadratic form* of the mass velocity $\dot{\mathbf{p}}$. It is always positive (as $m > 0$) and takes the minimum value $T = 0$. This observation is generally true for any mechanical system
- The generalized coordinates \mathbf{q} describe the configuration of the entire mechanical system we consider. I.e. if our punctual mass is part of a mechanical system having

the generalized coordinates q , then its position p is a function of q (in the sense that p can be calculated from q). More formally, one could write $p(q)$. The mass velocity \dot{p} is then a direct result of the chain rule:

$$\dot{p} = \frac{\partial p}{\partial q} \dot{q} \quad (2.2)$$

- Imagine that our mechanical system is made of N punctual masses of position p_i , fully determined by the generalized coordinates q . The kinetic energy of the system is then an addition of all the individual kinetic energies, i.e.

$$T = \frac{1}{2} \sum_{i=1}^N m_i \|\dot{p}_i\|^2 = \frac{1}{2} \sum_{i=1}^N m_i \dot{p}_i^\top \dot{p}_i \quad (2.3)$$

The generalized coordinates describe the configuration of the whole system, i.e. it describes the position p_i of each mass. Each velocity is given by the chain rule:

$$\dot{p}_i = \frac{\partial p_i}{\partial q} \dot{q} \quad (2.4)$$

The kinetic energy of the N masses is then given by:

$$T = \frac{1}{2} \sum_{i=1}^N m_i \dot{p}_i^\top \dot{p}_i = \frac{1}{2} \sum_{i=1}^N m_i \dot{q}^\top \frac{\partial p_i}{\partial q}^\top \frac{\partial p_i}{\partial q} \dot{q} = \frac{1}{2} \dot{q}^\top \left(\sum_{i=1}^N m_i \frac{\partial p_i}{\partial q}^\top \frac{\partial p_i}{\partial q} \right) \dot{q} \quad (2.5)$$

Here we can define:

$$W = \sum_{i=1}^N m_i \frac{\partial p_i}{\partial q}^\top \frac{\partial p_i}{\partial q} \quad (2.6)$$

and we observe that W is a square, symmetric, (semi-)positive-definite matrix of size $\mathbb{R}^{n_q \times n_q}$. This matrix is in general a function of q , but can also be constant for “smart” choices of generalized coordinates.

- The observation above holds for systems having infinitely many “particles” (e.g. distributed mass systems). Generally, the kinetic energy function will take the form:

$$T(q, \dot{q}) = \frac{1}{2} \dot{q}^\top W(q) \dot{q} \quad (2.7)$$

where each of these variables (i.e. q and \dot{q}) are time-varying, and where W is a square, symmetric, positive-definite matrix of size $\mathbb{R}^{n_q \times n_q}$, possibly function of q .

- Note that since the kinetic energy function T is given by (2.7), it is essentially defined via matrix $W(q)$, hence in our constructions in the following, matrix $W(q)$ will play a central role. In particular, we will see that $W(q)$ actually being a function of q or being constant plays a central role in the complexity of the model describing the corresponding mechanical system.

Example - mass on a rigid rod: imagine a punctual mass m attached to the origin by a rigid, massless rod of length L , and moving in the plane $x - y$. We can describe the system by the angle θ between the x axis and the rod. This angle alone describes entirely this simple system, as knowing specifies in what “configuration” (or position) the system is. Our (unique) generalized coordinate here is $q = \theta \in \mathbb{R}$. The position of the mass in the $x - y$ plane is then given by:

$$\mathbf{p} = L \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} \quad (2.8)$$

The chain rule then tells us that:

$$\dot{\mathbf{p}} = L \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix} \dot{\theta} \quad (2.9)$$

Our kinetic energy is then:

$$T = \frac{1}{2} m \dot{\mathbf{p}}^\top \dot{\mathbf{p}} = \frac{1}{2} mL^2 \dot{\theta}^2 \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix}^\top \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix} = \frac{1}{2} mL^2 \dot{\theta}^2 (\sin^2 \theta + \cos^2 \theta) = \frac{1}{2} mL^2 \dot{\theta}^2 \quad (2.10)$$

Here “matrix” W is simply

$$W = mL^2 \quad (2.11)$$

and is constant. **Example - mass on an elastic rod:** imagine a punctual mass m attached to the origin by an elastic, massless rod of varying length L , and moving in the plane $x - y$. To describe the system we now need the angle θ between the x axis and the rod, and the length L (as the latter is also changing and therefore not a simple constant). Our generalized coordinates are then:

$$\mathbf{q} = \begin{bmatrix} \theta \\ L \end{bmatrix} \quad (2.12)$$

as both are needed at a given time t to describe the “configuration” (or position) the system. The position of the mass in the $x - y$ plane is still given by:

$$\mathbf{p} = L \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} \quad (2.13)$$

but now L is no longer a constant and is part of our generalized coordinates. The chain rule then tells us that:

$$\dot{\mathbf{p}} = L \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix} \dot{\theta} + \dot{L} \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} \quad (2.14)$$

Our kinetic energy is then:

$$\begin{aligned}
T &= \frac{1}{2} m \dot{\mathbf{p}}^\top \dot{\mathbf{p}} \\
&= \frac{1}{2} mL^2 \dot{\theta}^2 \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix}^\top \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix} + \frac{1}{2} m \dot{L}^2 \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix}^\top \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} \\
&\quad + \frac{1}{2} m L \dot{L} \dot{\theta} \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix}^\top \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} \\
&= \frac{1}{2} mL^2 \dot{\theta}^2 + \frac{1}{2} m \dot{L}^2 = \frac{1}{2} \dot{\mathbf{q}}^\top \begin{bmatrix} mL^2 & 0 \\ 0 & m \end{bmatrix} \dot{\mathbf{q}}
\end{aligned} \tag{2.15}$$

Here matrix W is diagonal:

$$W = \begin{bmatrix} mL^2 & 0 \\ 0 & m \end{bmatrix} \tag{2.16}$$

and not constant as L is part of the state.

Example - rigid rod with a mass: imagine a uniform rod of length L and of negligible thickness attached at the origin, having a mass m (uniformly distributed). The rod position is again described by the angle θ alone (we consider the length L fixed as the rod is rigid), hence $\mathbf{q} = \theta$. We can consider our rod as made of infinitely many punctual masses distributed along the rod axis. If we locate these masses via their position $v \in [0, L]$ along the rod axis, we can describe their position in the plane $x - y$ by:

$$\mathbf{p}(v) = v \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} \tag{2.17}$$

and their velocity is given by the chain rule:

$$\dot{\mathbf{p}}(v) = v \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix} \dot{\theta} \tag{2.18}$$

The overall kinetic energy of the rod is given by the summation of the kinetic energy of all the infinitesimal masses (each having a mass $\frac{m}{L} dv$). It can be computed by the integral:

$$\begin{aligned}
T &= \frac{m}{L} \int_0^L \left(v \dot{\theta} \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix} \right)^\top \left(v \dot{\theta} \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix} \right) dv \\
&= \frac{m}{L} \int_0^L v^2 \dot{\theta}^2 dv = \frac{m}{3L} \dot{\theta}^2 v^3 \Big|_{v=0}^{v=L} = \frac{m}{3L} \dot{\theta}^2 L^3 = \frac{1}{3} mL^2 \dot{\theta}^2
\end{aligned} \tag{2.19}$$

Our "matrix" W is now reduced to $2/3$ of the one we had for the mass concentrated at the end of the rod, i.e.

$$W = \frac{2}{3} mL^2 \tag{2.20}$$

(one can compare (2.19) and (2.10) to see that).

□

2.2 POTENTIAL ENERGY

Let us provide a few examples of potential energy, some of which will be extensively used in this course.

- The potential energy due to gravity in most “standard” mechanical applications derives from:

$$V = mgz \quad (2.21)$$

which gives the potential energy of a punctual mass m , where “ z ” is the “height” of the mass in the field of gravity. Consider e.g. in the examples above that the mass m is concentrated at the end of the rigid rod. The position of the mass is given by (2.17), such that its vertical position is given by:

$$p_z = L \sin \theta \quad (2.22)$$

Its potential energy is then given by:

$$V = mgL \sin \theta \quad (2.23)$$

In contrast, consider that the mass is distributed throughout the rod. Computing the potential energy then follows similar lines as (2.19), i.e. we need to “sum up” the potential energy of every “particle of the rod”. A “particle” is described here as a infinitesimal piece of the rod, of length dv , having a mass $\frac{m}{L}dv$ and a vertical position $v \sin \theta$. We can then do:

$$V = \frac{mg}{L} \int_0^L v \sin \theta dv = \frac{mg}{L} \frac{1}{2} v^2 \sin \theta \Big|_{v=0}^{v=L} = \frac{1}{2} mgL \sin \theta \quad (2.24)$$

hence the potential energy then corresponds to considering that the whole mass of the rod is concentrated at the half-length (this is to put in contrast with the kinetic energy computation, where one can consider that the whole mass of the rod is concentrated at the 2/3 of the length of the rod!).

- The considerations above are valid for mechanical systems evolving “at a small scale” in the field of gravity, for which one can consider the field as “straight” and uniform. For mechanical systems evolving at a “large scale” such as e.g. a satellite orbiting a planet, one needs to apply (at least) the genuine formula for classic gravity, i.e.:

$$V = -G \frac{m}{r} \quad (2.25)$$

where r is the distance to the center of the planet.

Note that the Lagrange formalism can handle potential energy from other sources (magnetic field, electric field, etc.), but we will not consider these cases in this course. Another source of potential energy that is commonly present in mechanical systems is potential energy stored in flexible components. We will often simply consider springs in this course, but

any part of the system that undergoes elastic deformations does in principle store kinetic energy. Consider the spring illustrated below.

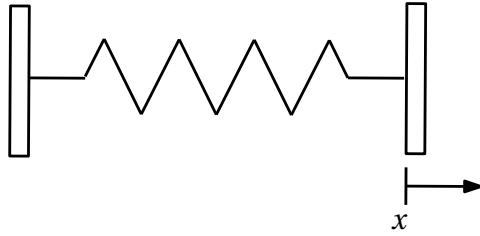


Figure 2.1: Illustration of a spring. The position $x = 0$ is set at the spring rest length. The rigidity constant of the spring is k

The force delivered by the spring is given by

$$F = -kx \quad (2.26)$$

where x is the position of the spring end, relative to its rest position (position at which the spring yields no force). The potential energy stored in the spring is given by:

$$V = \int_x^0 -k\nu d\nu = -\frac{1}{2}k\nu^2 \Big|_x^0 = \frac{1}{2}kx^2 \quad (2.27)$$

Note that if the reference frame is chosen such that the rest length of the spring is not at $x = 0$ but rather at an arbitrary position $x = x_0$, then the potential energy is given by:

$$V = \frac{1}{2}k(x - x_0)^2 \quad (2.28)$$

Let us consider now a slightly more complex example illustrated in the figure below.

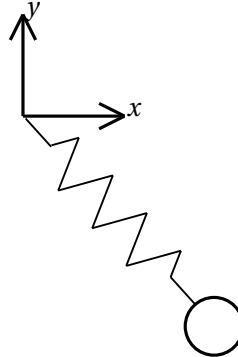


Figure 2.2: Illustration of a spring moving in the plane. The rigidity constant of the spring is k . We suppose that the rest length of the spring is L_0 .

We consider that the position of the end of the spring is given by

$$\mathbf{p} = \begin{bmatrix} x \\ y \end{bmatrix} \quad (2.29)$$

such that the elongation of the spring with respect to its rest length is given by:

$$\Delta = (\mathbf{p}^\top \mathbf{p})^{\frac{1}{2}} - L_0 \quad (2.30)$$

The potential energy stored in the spring is then given by:

$$V = \frac{1}{2} k \Delta^2 = \frac{1}{2} k \left((\mathbf{p}^\top \mathbf{p})^{\frac{1}{2}} - L_0 \right)^2 \quad (2.31)$$

The latter expression is fairly complex, as the $\cdot^{\frac{1}{2}}$ does not simplify with the \cdot^2 unless $L_0 = 0$. If $L_0 = 0$, then the potential energy simplifies to:

$$V = \frac{1}{2} k \mathbf{p}^\top \mathbf{p} \quad (2.32)$$

In many exercises in this course, will use the simplifying assumption that $L_0 = 0$ in order to get simpler expressions.

2.3 LAGRANGE EQUATION

Before detailing the Euler-Lagrange equation, we need to define the Lagrange function

$$\mathcal{L}(\mathbf{q}, \dot{\mathbf{q}}) = T(\mathbf{q}, \dot{\mathbf{q}}) - V(\mathbf{q}) \in \mathbb{R} \quad (2.33)$$

simply made from subtracting the potential energy to the kinetic energy. In general, the Lagrange function takes the generalized coordinates \mathbf{q} and their time derivative $\dot{\mathbf{q}}$ as arguments. Using the observations of the Sections above, we observe that the Lagrange function takes the form:

$$\mathcal{L}(\mathbf{q}, \dot{\mathbf{q}}) = \frac{1}{2} \dot{\mathbf{q}}^\top W(\mathbf{q}) \dot{\mathbf{q}} - V(\mathbf{q}) \quad (2.34)$$

The Euler-Lagrange equation then reads as:

$$\frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{\mathbf{q}}} - \frac{\partial \mathcal{L}}{\partial \mathbf{q}} = 0 \quad (2.35)$$

and defines the model of the mechanical system. It is useful to observe here that (2.35) defines equations as a *lying vector*. Indeed, e.g. $\frac{\partial \mathcal{L}}{\partial \mathbf{q}} \in \mathbb{R}^{1 \times n_q}$ is a line-vector. One can verify that the first term is (of course) also a line vector such that the subtraction between the two terms is well defined. As we tend to prefer working with equations in a “standing-vector form”, it can be useful when needed to rewrite (2.35) in its “transposed” version, i.e.:

$$\frac{d}{dt} \nabla_{\dot{\mathbf{q}}} \mathcal{L} - \nabla_{\mathbf{q}} \mathcal{L} = 0 \quad (2.36)$$

where we use the “gradient” notation:

$$\nabla_{\dot{\mathbf{q}}} \mathcal{L} = \frac{\partial \mathcal{L}^\top}{\partial \dot{\mathbf{q}}} , \quad \nabla_{\mathbf{q}} \mathcal{L} = \frac{\partial \mathcal{L}^\top}{\partial \mathbf{q}} \quad (2.37)$$

This transformation is cosmetic and of secondary importance, but it is useful to point it out. Remember this observation in the following whenever we transpose the terms appearing in (2.35), as we will do it whenever it will help us write things in a neat and structured form.

Before investigating how (2.35) delivers the model of the system, it is worth here pausing to understand the mathematical meaning of the first term in (2.35). Indeed, the partial derivative of the Lagrange function \mathcal{L} with respect to the time derivative of the generalized coordinates \dot{q} , i.e. $\frac{\partial \mathcal{L}}{\partial \dot{q}}$ may appear surprising or ambiguous. In fact, it is a very straightforward operation. In order to perform it correctly, one has to consider \dot{q} as a variable in itself, independent of all other variables, and take the classical differential operations accordingly. More specifically, considering (2.34), the operation $\frac{\partial \mathcal{L}}{\partial \dot{q}}$ simply yields:

$$\nabla_{\dot{q}} \mathcal{L} = \nabla_{\dot{q}} \left(\frac{1}{2} \dot{q}^T W(q) \dot{q} - V(q) \right) = \nabla_{\dot{q}} \left(\frac{1}{2} \dot{q}^T W(q) \dot{q} \right) = W(q) \dot{q}. \quad (2.38)$$

Forming the term $\frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{q}}$ in the Lagrange equation requires a *total differentiation* with respect to time (operator $\frac{d}{dt}$). Using the chain rule, this yields:

$$\frac{d}{dt} (W(q) \dot{q}) = \frac{\partial}{\partial \dot{q}} (W(q) \dot{q}) \ddot{q} + \frac{\partial}{\partial q} (W(q) \dot{q}) \dot{q} = W(q) \ddot{q} + \frac{\partial}{\partial q} (W(q) \dot{q}) \dot{q} \quad (2.39)$$

We can then pack the developments above to observe that the Lagrange equation yields:

$$W(q) \ddot{q} + \frac{\partial}{\partial q} (W(q) \dot{q}) \dot{q} - \underbrace{\nabla_q \mathcal{L}}_{= \nabla_q T - \nabla_q V} = 0 \quad (2.40)$$

In order to understand the “mechanics” behind using the Lagrange equations, it is useful to run through a handful of simple and classical examples. It is really useful to underline here that **the computations performed hereafter are best carried out in a Computer Algebra System (CAS) such as the Matlab Symbolic Toolbox**. Indeed, the computations required to deploy Lagrange mechanics are very systematic, but become quickly involved. It is best to perform them in the computer. In the exercises associated to this course, we recommend you to use this approach to generate your models.

We can make a few observations here.

1. The Lagrange equation delivers the differential equation in an *implicit form*, i.e. it is a set of equations relating the generalized coordinates and their first and second-order time derivatives, i.e. q, \dot{q}, \ddot{q} .
2. The Lagrange equation allows one to compute the system acceleration \ddot{q} for a given system configuration q and its time derivative \dot{q} if matrix $W(q)$ is invertible. They can be used for simulating the system, if the initial conditions q, \dot{q} are provided (i.e. we need to specify the “position” and “velocity” of the system in order to predict its future trajectory). For q, \dot{q} given, then the Lagrange equations can be solved for the acceleration \ddot{q} , from which the time evolution of q, \dot{q} can be computed (i.e. q is obtained from integrating \dot{q} and \dot{q} is obtained by integrating \ddot{q}).

3. The second term

$$\nabla_q \mathcal{L} = \nabla_q T - \nabla_q V \quad (2.41)$$

in the Lagrange equation is generally speaking introducing forces that are intrinsic to the system (as opposed to forces applied "externally" to the system), such as e.g. forces deriving from potentials (coming from the "V" part, see remark following (2.51)) and centrifugal forces (coming from the "T" part).

4. The Lagrange equation is linear in the accelerations \ddot{q} . This is true for Lagrange functions in the form (2.34), and stem from (2.39), which yields the first term in the Lagrange equation in the form $W(q)\ddot{q}$, which is linear in \ddot{q} .
5. Because the accelerations enter linearly in the Lagrange equation, we can solve it for \ddot{q} . More specifically, if one write the Lagrange equation in the form (2.40), the accelerations \ddot{q} can be explicitly expressed as:

$$\ddot{q} = W(q)^{-1} \left[\nabla_q \mathcal{L} - \frac{\partial}{\partial q} (W(q) \dot{q}) \dot{q} \right] \quad (2.42)$$

Unfortunately, writing the model in this explicit form is not always a very good move. Indeed, for systems that do not have a vector of generalized coordinates q of very low dimension (e.g. > 2), the inverse $W(q)^{-1}$ can be very complex, even if matrix $W(q)$ is fairly simple. An exception to this observation is the case where matrix $W(q)$ is constant (i.e. not actually function of q). In this case

$$\nabla_q T = 0 \Rightarrow \nabla_q \mathcal{L} = -\nabla_q V \quad \text{and} \quad \frac{\partial}{\partial q} (W(q) \dot{q}) \dot{q} = 0 \quad (2.43)$$

such that (2.42) becomes:

$$\ddot{q} = -W^{-1} \nabla_q V \quad (2.44)$$

where matrix W is purely "numerical" (i.e. it contains only numbers, not expressions), and is therefore easy to invert.

Let us try to nail down these remarks via a series of examples, starting simple and ramping up to complex systems.

Examples

1. Consider the spring-mass system depicted in Fig. 2.3 below.

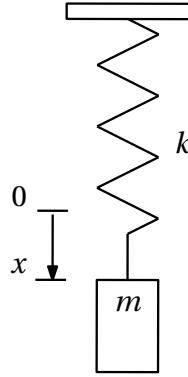


Figure 2.3: Vertical spring-mass system. The mass “0” position is set at the spring rest length. The rigidity constant of the spring is k , and the hanging mass is m .

Let us deploy the principles of Lagrange modelling on this simple system. Our goal is to unpack what the Lagrange equation does, and what the different terms “physically” correspond to. Let us start with setting up our generalized coordinates. We need to describe the mass position. Fig. 2.3 proposes to set the “0” position at the rest length of the spring, and consider that x increases when the mass goes down. This choice is not unique: we could set this “0” position anywhere we want and decide that x increases when the mass goes up.

The kinetic energy function is simple to calculate for this kind of system. We can use (2.1), where our position $p \equiv x$ is a scalar (we work in 1D). We then simply get:

$$T = \frac{1}{2}m\dot{x}^2 \quad (2.45)$$

One can observe here that our kinetic energy function is in the form (2.7), with $W(x) = m$. The potential energy is composed of the sum of gravity and of the energy stored in the spring (note that energies always add). These two terms read as:

$$V_{\text{gravity}} = -mgx \quad \text{and} \quad V_{\text{spring}} = \frac{1}{2}kx^2 \quad (2.46)$$

Note the minus sign on V_{gravity} . One can easily guess it by observing that if x increases the mass goes down and therefore the potential energy decreases. This requires the minus sign. The potential energy of a spring is always given by the quadratic form used here, i.e. it is given by the elongation of the spring (from rest length) squared, multiplied by the rigidity and divided by two.

We can now assemble the Lagrange function:

$$\mathcal{L} = T - (V_{\text{gravity}} + V_{\text{spring}}) = \frac{1}{2}m\dot{x}^2 + mgx - \frac{1}{2}kx^2 \quad (2.47)$$

Once the Lagrange function is assembled, the modelling does not require “intelligence” anymore, as the rest is just calculus and algebraic manipulations in order to

extract a differential equation from (2.35). This procedure can actually be automated, and its deployment entrusted to computer tools. If we do the exercise of computing (2.35) here. Let us start with:

$$\nabla_{\dot{q}} \mathcal{L} = \nabla_{\dot{x}} \mathcal{L} = m\ddot{x} \quad (2.48)$$

We observe that this term is a *momentum*. As a matter of fact, the term $\frac{\partial \mathcal{L}}{\partial \dot{q}}$ will always describe the momentum present in the system. We can also compare this expression with (2.38).

Furthermore, we can compute:

$$\frac{d}{dt} \nabla_{\dot{q}} \mathcal{L} = \frac{d}{dt} (m\dot{x}) = m\ddot{x} \quad (2.49)$$

We observe that this term is equivalent to the “ $m \cdot a$ ” of the $F = m \cdot a$ of Newton, i.e. it corresponds to a mass multiplied by accelerations. As in Newton, the product of a mass by an acceleration will equate forces. This observation will also hold true for all system, though “masses” and “forces” will take a somewhat more general meaning in Lagrange mechanics. We can also compare this expression to (2.39), and observe that

$$\frac{\partial}{\partial q} (W(q)\dot{q})\dot{q} = \underbrace{\frac{\partial W}{\partial x}}_{=0} \dot{x}^2 = 0 \quad (2.50)$$

because W is constant. We will see later on that having a matrix $W(q)$ that is constant simplifies significantly the resulting model equations.

We can turn to computing:

$$\nabla_q \mathcal{L} = \nabla_x \mathcal{L} = mg - kx \quad (2.51)$$

We ought to observe that this term delivers the forces acting on the system resulting from the potentials present in the system. We can now assemble the Lagrange equations, i.e. (2.35) for this simple example reads as:

$$m\ddot{x} - mg + kx = 0 \quad (2.52)$$

Because the acceleration \ddot{x} enters linearly in the Lagrange equation, we can solve it for \ddot{x} . Here we can trivially manipulate (2.52) to get:

$$m\ddot{x} = mg - kx. \quad (2.53)$$

This equation is essentially the Newton equation “ $F = m \cdot a$ ” for the spring-mass system. Indeed, we observe that:

$$\underbrace{m\ddot{x}}_{\equiv m \cdot a} = \underbrace{mg - kx}_{\equiv F} \quad (2.54)$$

2. We will now consider a linear crane as depicted in Fig. 2.4. This kind of system is e.g. important for loading/unloading cargo ships in large harbors. The proposed generalized coordinates are visible in Fig. 2.4, i.e.

$$\mathbf{q} = \begin{bmatrix} x \\ \theta \end{bmatrix} \in \mathbb{R}^2 \quad (2.55)$$

is made of the position of the cart on the rail, and the angle θ is the angle of the rod linking the cart to the hanging mass with respect to the vertical. We observe here that there is no problem mixing positions and angles (and any other physical quantity describing the "configuration" of a mechanical system) in the generalized coordinates.

We can now compute the kinetic and potential energy functions. As energies are additive, we can compute them separately for the two masses.

The kinetic energy of the cart is simply $T_{\text{cart}} = \frac{1}{2}M\dot{x}^2$. The energy of the hanging mass is a bit more complex to compute here, but can still be simply derived from formula (2.1). Indeed, we observe that the position of the hanging mass is given by:

$$\mathbf{p}_m = \begin{bmatrix} x + L\sin\theta \\ -L\cos\theta \end{bmatrix} \quad (2.56)$$

Note that the positive sign in the first line will specify in which direction a positive angle will "rotate" the hanging mass. We can then compute:

$$\dot{\mathbf{p}}_m = \frac{\partial \mathbf{p}_m}{\partial \mathbf{q}} \dot{\mathbf{q}} = \begin{bmatrix} \dot{x} + L\dot{\theta}\cos\theta \\ L\dot{\theta}\sin\theta \end{bmatrix} \quad (2.57)$$

and

$$\|\dot{\mathbf{p}}_m\|^2 = \dot{\mathbf{p}}_m^\top \dot{\mathbf{p}}_m = \dots = L^2\dot{\theta}^2 + \dot{x}^2 + 2L\dot{\theta}\dot{x}\cos\theta \quad (2.58)$$

The kinetic energy of the hanging crane is therefore given by:

$$T = \frac{1}{2}M\dot{x}^2 + \frac{1}{2}m\dot{\mathbf{p}}_m^\top \dot{\mathbf{p}}_m = \frac{1}{2}(m+M)\dot{x}^2 + \frac{1}{2}mL^2\dot{\theta}^2 + mL\dot{\theta}\dot{x}\cos\theta \quad (2.59)$$

We can deduce that the matrix $W(\mathbf{q})$ is now:

$$W(\mathbf{q}) = \begin{bmatrix} m+M & mL\cos\theta \\ mL\cos\theta & mL^2 \end{bmatrix} \quad (2.60)$$

and is not constant in this example (it is function of θ). It can be useful to observe here that the matrix $W(\mathbf{q})$ can be readily computed from the kinetic energy function by computing its Hessian with respect to $\dot{\mathbf{q}}$, i.e.

$$W(\mathbf{q}) = \frac{\partial^2 T(\mathbf{q}, \dot{\mathbf{q}})}{\partial \dot{\mathbf{q}}^2} \quad (2.61)$$

The potential energy of the hanging crane is just gravity, and is related to the "z" position of the hanging mass, i.e.:

$$V = mg [\begin{array}{cc} 0 & 1 \end{array}] p_m = -mgL\cos\theta \quad (2.62)$$

Our Lagrange function can now be assembled, and reads as:

$$\mathcal{L} = \frac{1}{2}(m+M)\dot{x}^2 + \frac{1}{2}mL^2\dot{\theta}^2 + mL\dot{\theta}\dot{x}\cos\theta + mgL\cos\theta \quad (2.63)$$

We can now develop the terms of the Lagrange equation:

$$\nabla_{\dot{q}} \mathcal{L} = \left[\begin{array}{c} (M+m)\dot{x} + mL\dot{\theta}\cos\theta \\ mL^2\ddot{\theta} + mL\dot{x}\cos\theta \end{array} \right] = W(q)\dot{q} \quad (2.64)$$

and

$$\frac{d}{dt} \nabla_{\dot{q}} \mathcal{L} = \left[\begin{array}{c} (M+m)\ddot{x} + mL\cos\theta\ddot{\theta} - mL\sin\theta\dot{\theta}^2 \\ mL\cos\theta\ddot{x} + mL^2\ddot{\theta} - mL\dot{\theta}\dot{x}\sin\theta \end{array} \right] \quad (2.65)$$

Moreover,

$$\nabla_q \mathcal{L} = -mL\sin\theta \left[\begin{array}{c} 0 \\ g + \dot{\theta}\dot{x} \end{array} \right] \quad (2.66)$$

The Lagrange equation then reads as:

$$\left[\begin{array}{c} (M+m)\ddot{x} + mL\cos\theta\ddot{\theta} - mL\sin\theta\dot{\theta}^2 \\ mL\cos\theta\ddot{x} + mL^2\ddot{\theta} + mgL\sin\theta \end{array} \right] = 0 \quad (2.67)$$

The acceleration \ddot{q} enters linearly in (2.67). We can actually observe that (2.67) can be written as:

$$W(q)\ddot{q} + mL\sin\theta \left[\begin{array}{c} -\dot{\theta}^2 \\ g \end{array} \right] = 0 \quad (2.68)$$

And inverting $W(q)$, we can write:

$$\ddot{q} = mL\sin\theta W(q)^{-1} \left[\begin{array}{c} \dot{\theta}^2 \\ -g \end{array} \right] \quad (2.69)$$

Unfortunately, this construction yields fairly complex expressions. We can verify that:

$$\ddot{q} = \frac{\sin\theta}{(M+m-m\cos\theta^2)} \left[\begin{array}{c} m(L\dot{\theta}^2 + g\cos\theta) \\ -\frac{1}{L}((M+m)g + mL\cos\theta\dot{\theta}^2) \end{array} \right] \quad (2.70)$$

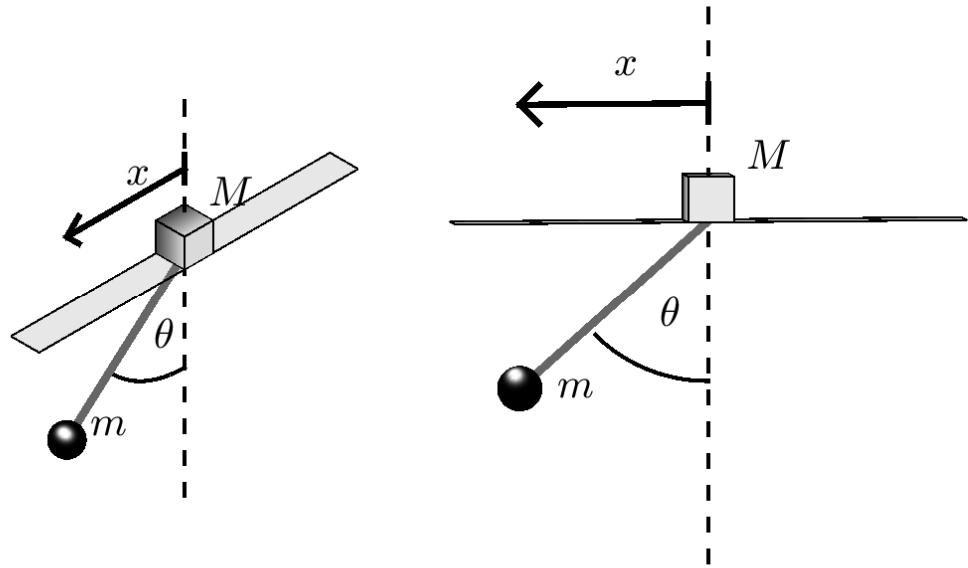


Figure 2.4: Simple hanging crane. The mass M moves on a rail (position x) and mass m is linked to it via a massless rod of length L . We describe the position of the hanging mass via the angle θ .

3. Let us consider an elastic crane, identical to the one depicted in Fig. 2.4, but where the rod is not rigid, hence the length L becomes variable, and the associated potential energy stored in the rod is then given by:

$$V_{\text{rod}} = \frac{1}{2}K(L - L_0)^2 \quad (2.71)$$

where L_0 is the rest-length of the rod. We will now adopt the generalized coordinates

$$\mathbf{q} = \begin{bmatrix} x \\ \theta \\ L \end{bmatrix} \quad (2.72)$$

The kinetic energy of the system now needs to account for L being time-varying, but the computations performed previously do not change. I.e. we can compute:

$$\dot{\mathbf{p}}_m = \frac{\partial \mathbf{p}_m}{\partial \mathbf{q}} \dot{\mathbf{q}} = \begin{bmatrix} \dot{x} + L\dot{\theta} \cos\theta + \dot{L} \sin\theta \\ L\dot{\theta} \sin\theta - \dot{L} \cos\theta \end{bmatrix} \quad (2.73)$$

and

$$\|\dot{\mathbf{p}}_m\|^2 = \dot{\mathbf{p}}_m^\top \dot{\mathbf{p}}_m = \dots = L^2\dot{\theta}^2 + \dot{x}^2 + 2L\dot{\theta}\dot{x}\cos\theta + \dot{L}^2 + 2\dot{L}\dot{x}\sin\theta \quad (2.74)$$

The kinetic energy of the hanging crane is therefore given by:

$$T = \frac{1}{2}M\dot{x}^2 + \frac{1}{2}m\dot{\mathbf{p}}_m^\top \dot{\mathbf{p}}_m = \frac{1}{2}(m+M)\dot{x}^2 + \frac{1}{2}mL^2\dot{\theta}^2 + mL\dot{\theta}\dot{x}\cos\theta + \frac{1}{2}m\dot{L}^2 + m\dot{L}\dot{x}\sin\theta \quad (2.75)$$

Matrix $W(\mathbf{q})$ is not quite complex in this example. It reads as:

$$W(\mathbf{q}) = \begin{bmatrix} m+M & mL\cos\theta & m\sin\theta \\ mL\cos\theta & mL^2 & 0 \\ m\sin\theta & 0 & m \end{bmatrix} \quad (2.76)$$

It is the same as (2.60), but with an extra row and an extra column stemming from the new coordinate L . Our potential energy now includes the new term (2.71) corresponding to the elastic rod. It reads as:

$$V = -mgL\cos\theta + \frac{1}{2}K(L - L_0)^2 \quad (2.77)$$

We can then assemble the Lagrange function and compute the different terms of the Lagrange equation. We have:

$$\nabla_{\dot{\mathbf{q}}} \mathcal{L} = \begin{bmatrix} (M+m)\dot{x} + mL\dot{\theta}\cos\theta + m\dot{L}\sin\theta \\ mL^2\ddot{\theta} + mL\dot{x}\cos\theta \\ m\ddot{L} + m\dot{x}\sin\theta \end{bmatrix} = W(\mathbf{q})\dot{\mathbf{q}} \quad (2.78a)$$

$$\frac{d}{dt} \nabla_{\dot{\mathbf{q}}} \mathcal{L} = \begin{bmatrix} (M+m)\ddot{x} + mL\cos\theta\ddot{\theta} - mL\sin\theta\dot{\theta}^2 + 2m\dot{L}\cos\theta\dot{\theta} + m\ddot{L}\sin\theta \\ mL\cos\theta\ddot{x} + mL^2\ddot{\theta} - mL\dot{\theta}\dot{x}\sin\theta + m\dot{L}\dot{x}\cos\theta + 2mL\dot{L}\dot{\theta} \\ m\ddot{L} + m\ddot{x}\sin\theta + m\dot{\theta}\dot{x}\cos\theta \end{bmatrix} \quad (2.78b)$$

$$\nabla_{\mathbf{q}} \mathcal{L} = \begin{bmatrix} 0 \\ -m(Lg\sin\theta - \dot{L}\dot{x}\cos\theta + L\dot{\theta}\dot{x}\sin\theta) \\ mL\dot{\theta}^2 + m\dot{x}\cos\theta\dot{\theta} + K(L - L_0) + mg\cos\theta \end{bmatrix} \quad (2.78c)$$

The Lagrange equation then reads as:

$$\begin{bmatrix} (M+m)\ddot{x} + m\sin\theta\ddot{L} + mL\cos\theta\ddot{\theta} - mL\sin\theta\dot{\theta}^2 + 2m\dot{L}\cos\theta\dot{\theta} + \\ mL(L\ddot{\theta} + 2\dot{L}\dot{\theta} + \ddot{x}\cos\theta + g\sin\theta) \\ m\ddot{L} + m\sin\theta\ddot{x} + K(L - L_0) - mg\cos\theta - mL\dot{\theta}^2 \end{bmatrix} = 0 \quad (2.79)$$

As previously, the acceleration \ddot{q} enters linearly in (2.79). We can again observe that (2.79) can be written as:

$$W(q)\ddot{q} + \begin{bmatrix} m\dot{\theta}(2\dot{L}\cos\theta - L\dot{\theta}\sin\theta) \\ mL(2\dot{L}\dot{\theta} + g\sin\theta) \\ -mL\dot{\theta}^2 + K(L - L_0) - mg\cos\theta \end{bmatrix} = 0 \quad (2.80)$$

and the acceleration can be computed explicitly using:

$$\ddot{q} = -W(q)^{-1} \begin{bmatrix} m\dot{\theta}(2\dot{L}\cos\theta - L\dot{\theta}\sin\theta) \\ mL(2\dot{L}\dot{\theta} + g\sin\theta) \\ -mL\dot{\theta}^2 + K(L - L_0) - mg\cos\theta \end{bmatrix} \quad (2.81)$$

Unfortunately, (2.81) yields very long expressions and we skip writing it here.

4. Let us revisit the elastic crane example with a different choice of generalized coordinates. Instead of specifying the position of the hanging mass using (θ, L) (these are polar coordinates), it is equally valid to use cartesian coordinates. Let us describe the hanging mass using

$$\mathbf{p} = \begin{bmatrix} p_1 \\ p_2 \end{bmatrix} \in \mathbb{R}^2 \quad (2.82)$$

as illustrated in Fig. 2.5, and we chose to order our generalized coordinates q as:

$$q = \begin{bmatrix} x \\ p \end{bmatrix} \in \mathbb{R}^3 \quad (2.83)$$

The kinetic energy then reads as:

$$T = \frac{1}{2}M\dot{x}^2 + \frac{1}{2}m\dot{\mathbf{p}}^\top \dot{\mathbf{p}} \quad (2.84)$$

and the $W(q)$ matrix reads as:

$$W(q) = \begin{bmatrix} M & 0 & 0 \\ 0 & m & 0 \\ 0 & 0 & m \end{bmatrix} \quad (2.85)$$

and is constant. The potential energy in the rod reads as:

$$V_{\text{rod}} = \frac{1}{2}K(L - L_0)^2 \quad \text{where} \quad L = \left\| \mathbf{p} - \begin{bmatrix} x \\ 0 \end{bmatrix} \right\| \quad (2.86)$$

and the potential energy of the hanging mass is:

$$V_{\text{gravity}} = -mg p_2 \quad (2.87)$$

(note that the minus sign here is due to the reference frame chosen for the hanging mass, with the vertical basis vector oriented down, see Fig. 2.5). Hence the Lagrange function reads as:

$$\mathcal{L} = T - V = \frac{1}{2}M\dot{x}^2 + \frac{1}{2}m\dot{\mathbf{p}}^\top \dot{\mathbf{p}} - \frac{1}{2}K(L - L_0)^2 - mg p_2 \quad (2.88)$$

We can then compute:

$$\nabla_{\dot{\mathbf{q}}} \mathcal{L} = \begin{bmatrix} M\dot{x} \\ m\dot{\mathbf{p}} \end{bmatrix} \quad (2.89a)$$

$$\frac{d}{dt} \nabla_{\dot{\mathbf{q}}} \mathcal{L} = \begin{bmatrix} M\ddot{x} \\ m\ddot{\mathbf{p}} \end{bmatrix} \quad (2.89b)$$

$$\nabla_{\mathbf{q}} \mathcal{L} = \begin{bmatrix} 0 \\ 0 \\ mg \end{bmatrix} + \frac{K(L - L_0)}{L} \begin{bmatrix} p_1 - x \\ x - p_1 \\ -p_2 \end{bmatrix} \quad (2.89c)$$

And write the model as:

$$\begin{bmatrix} \ddot{x} \\ \ddot{\mathbf{p}} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ g \end{bmatrix} + \frac{K(L - L_0)}{L} \begin{bmatrix} M^{-1}(p_1 - x) \\ m^{-1}(x - p_1) \\ -m^{-1}p_2 \end{bmatrix} \quad (2.90)$$

5. Consider an elastic hanging chain as depicted in Fig. 2.7. The chain is made of N punctual masses each of mass m that can move in a 3-dimensional space, linked together by elastic rods of rigidity K , and linked to two fixed points by elastic rods as well. The “configuration” of the hanging chain can be described by listing the positions $\mathbf{p}_{1,\dots,N}$ of each mass in the vector $\mathbf{q} \in \mathbb{R}^{3N}$. We will label the fixed end-points as \mathbf{p}_0 and \mathbf{p}_{N+1} , and these two points will not be part of the generalized coordinates \mathbf{q} . For the sake of simplicity, we will consider that the elastic links have a rest-length $L_0 = 0$. The kinetic energy will read as:

$$T = \frac{1}{2}m\dot{\mathbf{q}}^\top \dot{\mathbf{q}} \quad (2.91)$$

hence the matrix $W(\mathbf{q})$ reads as:

$$W(\mathbf{q}) = mI \quad (2.92)$$

where $I \in \mathbb{R}^{3N \times 3N}$ is the identity matrix, and is therefore constant. The potential energy function reads as:

$$V = \underbrace{mg \sum_{k=1}^N [0 \ 0 \ 1] \mathbf{p}_k}_{\text{gravity}} + \underbrace{\frac{1}{2}K \sum_{k=0}^N \|\mathbf{p}_{k+1} - \mathbf{p}_k\|^2}_{\text{spring}} \quad (2.93)$$

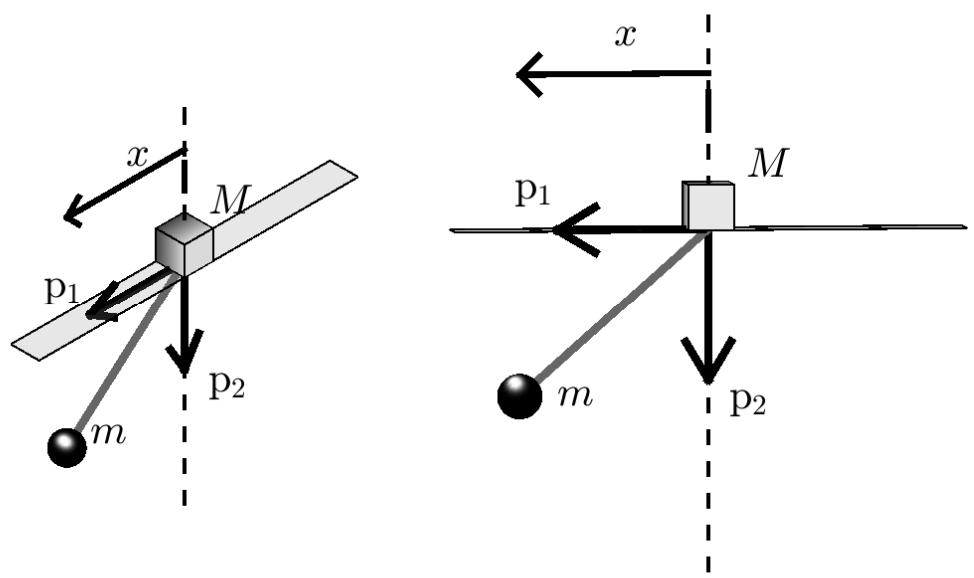


Figure 2.5: Simple hanging crane similar to Fig. 2.4. Here we use a cartesian coordinate system.

where \mathbf{p}_0 and \mathbf{p}_{N+1} are the fixed positions of the end points. The Lagrange function reads as:

$$\mathcal{L} = \frac{1}{2}m\dot{\mathbf{q}}^\top \dot{\mathbf{q}} - mg \sum_{k=1}^N [\begin{array}{ccc} 0 & 0 & 1 \end{array}] \mathbf{p}_k - \frac{1}{2}K \sum_{k=0}^N \|\mathbf{p}_{k+1} - \mathbf{p}_k\|^2 \quad (2.94)$$

We can then compute:

$$\nabla_{\dot{\mathbf{q}}} \mathcal{L} = m\dot{\mathbf{q}} \quad (2.95a)$$

$$\frac{d}{dt} \nabla_{\dot{\mathbf{q}}} \mathcal{L} = m\ddot{\mathbf{q}} \quad (2.95b)$$

$$\nabla_{\mathbf{q}} \mathcal{L} = K \begin{bmatrix} \mathbf{p}_0 - \mathbf{p}_1 \\ \mathbf{p}_1 - \mathbf{p}_2 \\ \vdots \\ \mathbf{p}_{N-1} - \mathbf{p}_N \end{bmatrix} + K \begin{bmatrix} \mathbf{p}_2 - \mathbf{p}_1 \\ \mathbf{p}_3 - \mathbf{p}_2 \\ \vdots \\ \mathbf{p}_{N+1} - \mathbf{p}_N \end{bmatrix} + m \begin{bmatrix} \mathbf{g} \\ \mathbf{g} \\ \vdots \\ \mathbf{g} \end{bmatrix}, \quad \mathbf{g} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ g \end{bmatrix} \quad (2.95c)$$

The Lagrange equation then delivers:

$$m\ddot{\mathbf{q}} + K \begin{bmatrix} \mathbf{p}_0 - \mathbf{p}_1 \\ \vdots \\ \mathbf{p}_{N-1} - \mathbf{p}_N \end{bmatrix} + K \begin{bmatrix} \mathbf{p}_2 - \mathbf{p}_1 \\ \vdots \\ \mathbf{p}_{N+1} - \mathbf{p}_N \end{bmatrix} + m \begin{bmatrix} \mathbf{g} \\ \vdots \\ \mathbf{g} \end{bmatrix} = 0 \quad (2.96)$$

which is trivial to put in an explicit form:

$$\ddot{\mathbf{q}} = -\frac{K}{m} \begin{bmatrix} \mathbf{p}_0 - \mathbf{p}_1 \\ \vdots \\ \mathbf{p}_{N-1} - \mathbf{p}_N \end{bmatrix} - \frac{K}{m} \begin{bmatrix} \mathbf{p}_2 - \mathbf{p}_1 \\ \vdots \\ \mathbf{p}_{N+1} - \mathbf{p}_N \end{bmatrix} - \begin{bmatrix} \mathbf{g} \\ \vdots \\ \mathbf{g} \end{bmatrix} \quad (2.97)$$

6. Let us consider now a 2-pendulum crane as illustrated in Fig. 2.6. We proceed with building the kinetic and potential energy functions for this system. To that end, it is best to start with describing the position of the two hanging masses (in a cartesian reference frame), as functions of the generalized coordinates

$$\mathbf{q} = \begin{bmatrix} x \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad (2.98)$$

The cart is at a position:

$$\mathbf{p}_0 = \begin{bmatrix} x \\ 0 \end{bmatrix} \quad (2.99)$$

The first hanging mass is at a cartesian position:

$$\mathbf{p}_1 = \mathbf{p}_0 + R(\theta_1) \begin{bmatrix} 0 \\ -L \end{bmatrix} \quad (2.100)$$

where

$$R(\theta_1) = \begin{bmatrix} \cos\theta_1 & -\sin\theta_1 \\ \sin\theta_1 & \cos\theta_1 \end{bmatrix} \quad (2.101)$$

and the second hanging mass is at a position: The first hanging mass is at a cartesian position:

$$\mathbf{p}_2 = \mathbf{p}_1 + R(\theta_2) \begin{bmatrix} 0 \\ -L \end{bmatrix} \quad (2.102)$$

We can then readily apply the chain rule (2.2), i.e. to obtain $\dot{\mathbf{p}}_{0,\dots,2}$. The kinetic energy function then reads as (the cart and the two haging masses are all of mass m):

$$T = \frac{1}{2}m \sum_{k=0}^2 \dot{\mathbf{p}}_k^\top \dot{\mathbf{p}}_k \quad (2.103)$$

but it is a fairly complex expression given by:

$$T = \frac{1}{2}\dot{\mathbf{q}}^\top W(\mathbf{q})\dot{\mathbf{q}} \quad (2.104)$$

where matrix $W(\mathbf{q})$ reads as:

$$W(\mathbf{q}) = m \begin{bmatrix} 3 & 2L\cos\theta_1 & L\cos\theta_2 \\ 2L\cos\theta_2 & 2L^2 & L\cos(\theta_1 - \theta_2) \\ L\cos\theta_2 & L\cos(\theta_1 - \theta_2) & L^2 \end{bmatrix} \quad (2.105)$$

The potential energy function is given by:

$$V = mg \sum_{k=1}^2 \mathbf{p}_{k,2} = -2mgL\cos\theta_1 - mgL\cos\theta_2 \quad (2.106)$$

We can then compute:

$$\nabla_{\dot{\mathbf{q}}} \mathcal{L} = W(\mathbf{q})\dot{\mathbf{q}} \quad (2.107a)$$

$$\frac{d}{dt} \nabla_{\dot{\mathbf{q}}} \mathcal{L} = W(\mathbf{q})\ddot{\mathbf{q}} + \frac{\partial}{\partial \mathbf{q}} (W(\mathbf{q})\dot{\mathbf{q}}) \dot{\mathbf{q}} \quad (2.107b)$$

$$\nabla_{\mathbf{q}} \mathcal{L} = -mL \begin{bmatrix} 0 \\ g\sin(\theta_1) + g\sin(\theta_1) + \dot{\theta}_1 \dot{x} \sin\theta_1 + \dot{x} \dot{\theta}_1 \sin\theta_1 + L\dot{\theta}_1 \dot{\theta}_2 \sin(\theta_1 - \theta_2) \\ g\sin\theta_2 + \dot{\theta}_2 \dot{x} \sin\theta_2 - L\dot{\theta}_1 \dot{\theta}_2 \sin(\theta_1 - \theta_2) \end{bmatrix} \quad (2.107c)$$

where

$$\frac{\partial}{\partial \mathbf{q}} (W(\mathbf{q})\dot{\mathbf{q}}) \dot{\mathbf{q}} = mL \begin{bmatrix} 2\sin\theta_1 \dot{\theta}_1^2 + \sin\theta_2 \dot{\theta}_2^2 \\ L\dot{\theta}_1 \dot{\theta}_2 \sin(\theta_1 - \theta_2) + 2\dot{x} \dot{\theta}_1 \sin\theta_1 - L\dot{\theta}_2^2 \sin(\theta_1 - \theta_2) \\ -L\dot{\theta}_1 \dot{\theta}_2 \sin(\theta_1 - \theta_2) + \dot{x} \dot{\theta}_2 \sin\theta_2 + L\dot{\theta}_1^2 \sin(\theta_1 - \theta_2) \end{bmatrix} \quad (2.108)$$

If one tries to assemble the resulting model in e.g. its explicit form, by computing (identically to (2.42)):

$$\ddot{\mathbf{q}} = W(\mathbf{q})^{-1} \left[\nabla_{\mathbf{q}} \mathcal{L} - \frac{\partial}{\partial \mathbf{q}} (W(\mathbf{q}) \dot{\mathbf{q}}) \dot{\mathbf{q}} \right] \quad (2.109)$$

However, the symbolic complexity resulting from computing (2.109) is high (and not reported here for this reason).

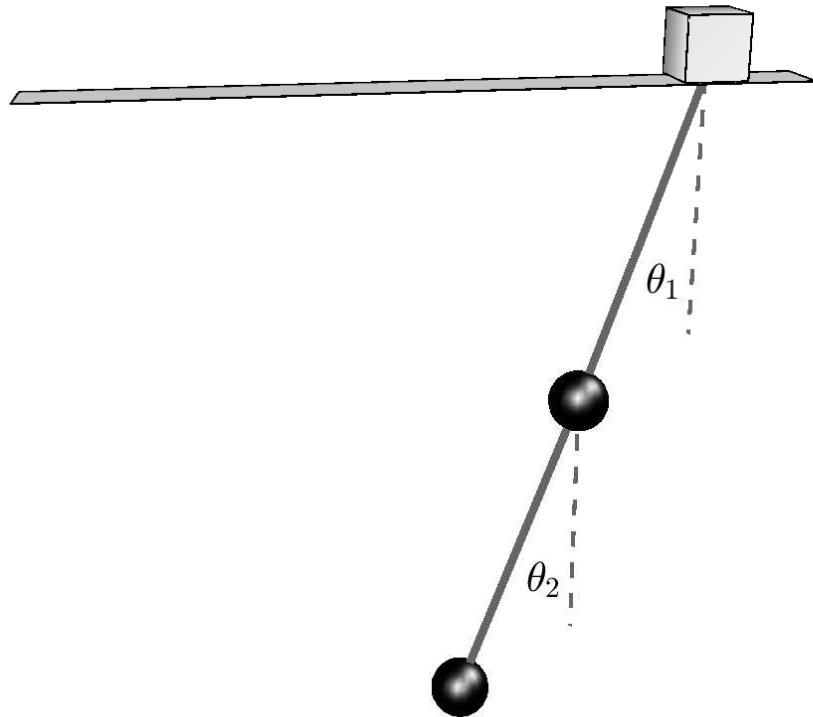


Figure 2.6: Simple hanging crane. The cart of mass m moves on a rail (position x , not shown here) and the two hanging masses also of mass m are linked to it via massless rods of length L . We describe the position of the hanging masses via the angle $\theta_{1,2}$ (relative to the verticals).

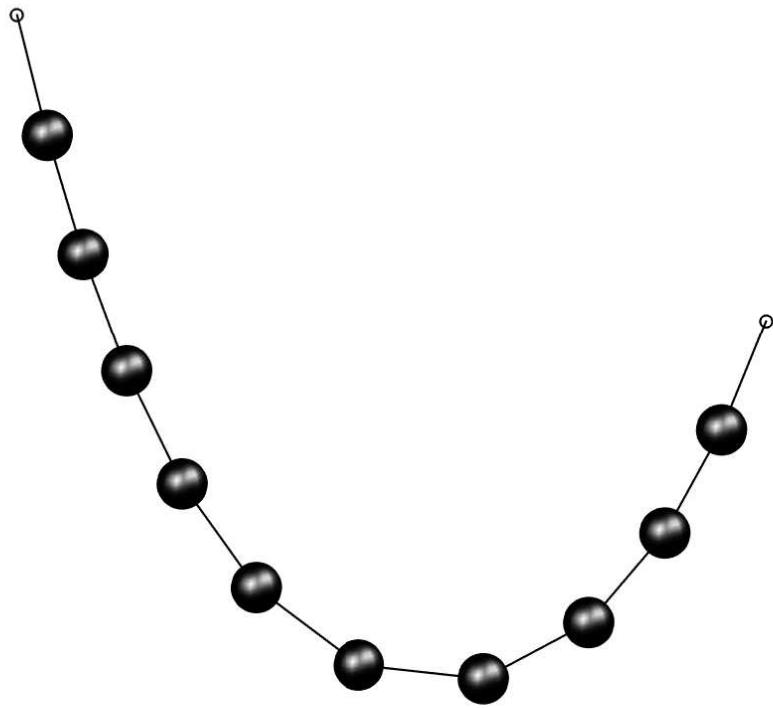


Figure 2.7: Illustration of the elastic hanging chain

□

2.4 EXTERNAL FORCES

The approach we have looked at so far does not include the possibility of having external forces and moments. We will close this gap now. One ought to observe that the Lagrange modelling approach is intrinsically an energy-based point of view of the mechanical system. As a result in the Lagrange approach, we will have to look at the external forces and moments in terms of the energy they deliver to or remove from the system. In other words, we will have to consider the work produced on the system by the external forces and moments. Work is related to motions occurring under forces and moments, and motions are described in the Lagrange approach as changes in the generalized coordinates q . The way external forces and moments are included in the Lagrange formulation in terms of **generalized forces** that we will label Q , which describe the amount of work produced on the system when moving the generalized coordinates. More formally, the generalized forces should satisfy:

$$\nabla_q E = Q \quad (2.110)$$

i.e. they describe the change of energy in the system when the generalized coordinates q are “moved a little bit”.

This concept is often presented in the literature via the following equation:

$$\delta W = \langle Q, \delta q \rangle \quad (2.111)$$

which is essentially saying the same as (2.110), i.e. a “small” motion δq , combined with the external forces and moments produces a “small” amount of work δW (a change of energy in the system), and the generalized forces Q relate the small motion to small amount of work produced. Once the generalized forces Q are known, they can be readily included in the Lagrange formalism using:

$$\frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{q}} - \frac{\partial \mathcal{L}}{\partial q} = F_q \quad (2.112)$$

A fairly systematic procedure can be derived from these concepts whenever a force of a torque is applied punctually somewhere in the system. Suppose that a force given by vector $F \in \mathbb{R}^n$ (where $n = 1, 2$ or 3 is the number of dimensions in which we are working) in a given fixed reference frame \mathcal{R} is applied at a specific point of the system, having a position $p \in \mathbb{R}^n$ in the same reference frame \mathcal{R} . We then observe that a small change in the generalized coordinates yields a small displacement of the position p given by the Jacobian $\frac{\partial p}{\partial q}$, and a small work:

$$\delta W = F^\top \frac{\partial p}{\partial q} \delta q \quad (2.113)$$

It follows that in this case the generalized force corresponding to F is given by:

$$Q = \frac{\partial p}{\partial q}^\top F = \nabla_q p F \quad (2.114)$$

This principle can be easily extended to several forces. When a set of forces $F_{1,\dots,m}$ is applied to a set of points $p_{1,\dots,m}$, then the generalized force is given by:

$$Q = \sum_{i=1}^m \frac{\partial p_i}{\partial q}^\top F_i \equiv \sum_{i=1}^m \nabla_q p_i F_i \quad (2.115)$$

In order to understand how to use these concepts in practice, let us see a few examples.

Example

- Consider the simple case of a point whose position is described in an orthonormal references frame, such that the generalized coordinates are $q \in \mathbb{R}^3$, and subject to a

force $\mathbf{F} \in \mathbb{R}^3$ described in the same reference frame. A small displacement $\mathbf{q} \rightarrow \mathbf{q} + \delta\mathbf{q}$ combined with the force \mathbf{F} yields a small amount of work:

$$\delta W = \langle \mathbf{F}, \delta\mathbf{q} \rangle = \mathbf{F}^\top \delta\mathbf{q} \quad (2.116)$$

such that in this specific case $\mathbf{Q} = \mathbf{F}$, i.e. the generalized force in this system is the force \mathbf{F} itself. One can verify that this is also given by (2.115). Indeed, since the position of the point is readily given by the generalized coordinates, i.e. $\mathbf{p} = \mathbf{q}$, it follows that an application of (2.115) yields (where $m = 1$):

$$\underbrace{\mathbf{Q} = \frac{\partial \mathbf{p}}{\partial \mathbf{q}}^\top}_{=I} \mathbf{F} = \mathbf{F} \quad (2.117)$$

- The very same principle applies to moments and rotations. Consider a axis of rotation whose position is described via the angle $\theta \in \mathbb{R}$, and subject to a moment $T \in \mathbb{R}$. The amount of work produces for a small motion $\delta\theta$ due to the moment T is then given by:

$$\delta W = \langle T, \delta\theta \rangle = T\delta\theta \quad (2.118)$$

such that the generalized force for this case is $\mathbf{Q} = T$. These two examples are fairly trivial. Let us consider next more complex ones.

- Consider the crane example illustrated in Fig. 2.6. Recall that the position of the two masses in a fixed reference frame is given by (see (2.100) and (2.102)):

$$\mathbf{p}_1 = \begin{bmatrix} x \\ 0 \end{bmatrix} + R(\theta_1) \begin{bmatrix} 0 \\ -L \end{bmatrix} \quad (2.119a)$$

$$\mathbf{p}_2 = \mathbf{p}_1 + R(\theta_2) \begin{bmatrix} 0 \\ -L \end{bmatrix} \quad (2.119b)$$

with the generalized coordinates

$$\mathbf{q} = \begin{bmatrix} x \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad (2.120)$$

We then observe that:

$$\frac{\partial \mathbf{p}_1}{\partial \mathbf{q}}^\top = \begin{bmatrix} 1 & 0 \\ L\cos\theta_1 & L\sin\theta_1 \\ 0 & 0 \end{bmatrix}, \quad \frac{\partial \mathbf{p}_2}{\partial \mathbf{q}}^\top = \begin{bmatrix} 1 & 0 \\ L\cos\theta_1 & L\sin\theta_1 \\ L\cos\theta_2 & L\sin\theta_2 \end{bmatrix} \quad (2.121)$$

If forces $\mathbf{F}_{1,2} \in \mathbb{R}^2$ are applied to masses m_1, m_2 respectively, then the generalized force then reads as:

$$\mathbf{Q} = \begin{bmatrix} 1 & 0 \\ L\cos\theta_1 & L\sin\theta_1 \\ 0 & 0 \end{bmatrix} \mathbf{F}_1 + \begin{bmatrix} 1 & 0 \\ L\cos\theta_1 & L\sin\theta_1 \\ L\cos\theta_2 & L\sin\theta_2 \end{bmatrix} \mathbf{F}_2 \quad (2.122)$$

2.5 CONSTRAINED LAGRANGE MECHANICS

We have been discussing so far cases where the generalized coordinates can “move independently” from each other, in the sense that any element of the set \mathbb{R}^{n_q} is admissible for the vector of generalized coordinates q . E.g.

- In the spring-mass system, any position x is “acceptable” for the system (even if for practical reasons the position of the mass may be restricted in a real system)
- In the hanging crane example, any value of θ and x are valid (even if in practice the rail may have a limited length and the hanging mass may not be physically allowed to get above the rail)

To clarify where we are going here, let us consider a simple example where the generalized coordinates are not free to “move independently”. Let us consider a “bowl” in 3D, described by the scalar equation $c(p) = 0$ (see Fig. 2.8 for an illustration), where $p \in \mathbb{R}^3$ are cartesian coordinates and

$$c(p) = p_3 - \frac{1}{4}(p_2^2 + p_1^2) \in \mathbb{R} \quad (2.123)$$

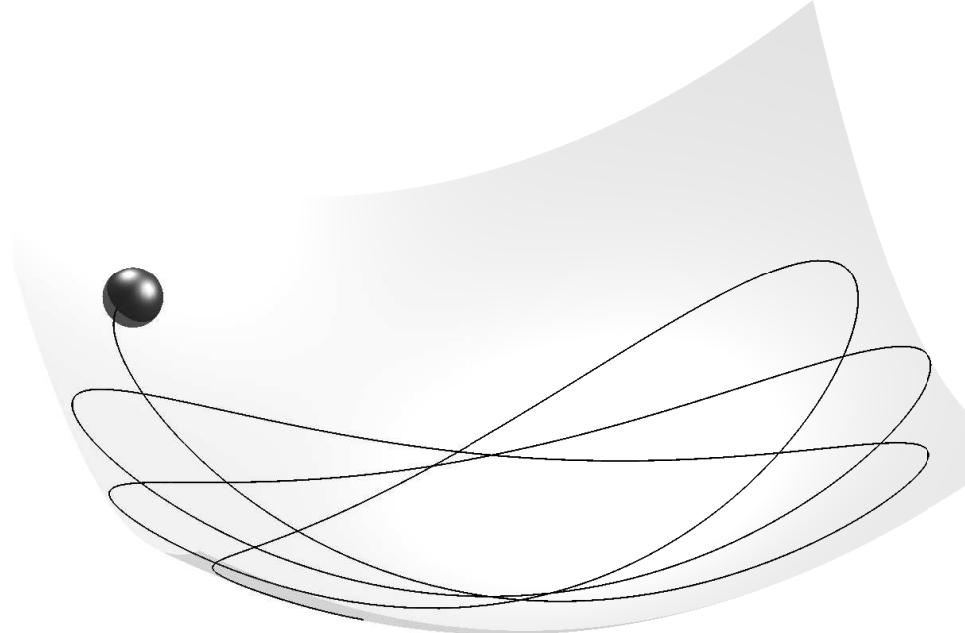


Figure 2.8: Illustration of the bowl example for (2.123).

Let us consider a mass m moving in 3D, but forced to slide on the surface of the bowl. The position of the mass can be described by the cartesian coordinates $p \in \mathbb{R}^3$, but since the mass moves at the surface of the bowl, the position p is not “free” to move everywhere, but

is “constrained” to move in the 2D “space” described by (2.123). Indeed, only positions that satisfy $c(p) = 0$ are admissible. Put differently, the generalized coordinates are not “free” to move independently.

More formally, we will consider mechanical systems of generalized coordinates $q \in \mathbb{R}^{n_q}$ where the generalized coordinates are constrained by $c(q) = 0$ with $c : \mathbb{R}^{n_q} \rightarrow \mathbb{R}^{n_c}$ and $n_c < n_q$. Mathematically, $c(q) = 0$ describes a *manifold* in \mathbb{R}^{n_q} , i.e. a (not necessarily flat) surface of dimension $n_q - n_c$. We ought to observe here that the number of dimension of the manifold in which the system moves, i.e. the difference $n_q - n_c$ is actually the number of degrees of freedom (DoF) that the system physically has. Being able to treat this type of problem systematically is extremely useful in a number of complex mechanical applications.

Fortunately, the Lagrange formalism we have discussed so far can be readily extended to constrained problems, with minimal changes. The construction of the kinetic and potential energy functions for the system can be done without any regard to the constraint function $c(q) = 0$. The constraints appear in the Lagrange function, which is then assembled as:

$$\mathcal{L}(q, \dot{q}, z) = T(q, \dot{q}) - V(q) - \langle z, c(q) \rangle \quad (2.124)$$

where z is a new set of variables, and usually labelled the *Lagrange multipliers* associated to the constraint function c . When function $c(q)$ has as image space the vector space \mathbb{R}^{n_c} (i.e. when it returns a “standard” vector), then $z \in \mathbb{R}^{n_c}$ and we can rewrite (2.124) simply as²:

$$\mathcal{L}(q, \dot{q}, z) = T(q, \dot{q}) - V(q) - z^\top c(q) \quad (2.125)$$

Beyond this point, the Lagrange formalism applies without any modification, i.e. the dynamics associated to a system of kinetic and potential energy T and V and constrained to evolve on the manifold $c(q)$ are described by:

$$\frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{q}} - \frac{\partial \mathcal{L}}{\partial q} = Q \quad (2.126a)$$

$$c(q) = 0 \quad (2.126b)$$

The mathematical justification for this construction are beyond the scope of this course, but we will try to build some intuitions via examples. Before deploying examples of application of the constrained Lagrange equation, we can do the same work as in (2.38) and (2.39). I.e. one can verify that the equalities:

$$\nabla_{\dot{q}} \mathcal{L} = W(q) \dot{q}. \quad (2.127)$$

and

$$\frac{d}{dt} \nabla_{\dot{q}} \mathcal{L} = \frac{d}{dt} (W(q) \dot{q}) = W(q) \ddot{q} + \frac{\partial}{\partial q} (W(q) \dot{q}) \dot{q} \quad (2.128)$$

²the generic form (2.124) allows one to consider constraint function in “exotic” vector spaces such as e.g. matrix spaces or Hilbert spaces. We will discuss only the simple case of “standard” vector spaces here.

still hold, and are actually not impacted by the presence of the constraints $c(q)$. The constraints modify the Lagrange equation via the term:

$$\nabla_q \mathcal{L} = \nabla_q T - \nabla_q V - \nabla_q c \cdot z \quad (2.129)$$

We ought to recall here remark 3. of Section 2.3: the term $\nabla_q \mathcal{L}$ in the Lagrange equation holds forces that are “intrinsic” to the system (potential, centrifugal). This remark still holds here, and we will see in the illustrative examples below that the term $\nabla_q c \cdot z$ is in fact akin to a force “keeping the system on the manifold c ”.

As observed previously in Section 2.3, we can assemble the Lagrange equation (2.126) in the form:

$$W(q) \ddot{q} + \frac{\partial}{\partial q} (W(q) \dot{q}) \dot{q} - \nabla_q T + \nabla_q V + \nabla_q c \cdot z = Q \quad (2.130a)$$

$$c(q) = 0 \quad (2.130b)$$

And the accelerations \ddot{q} can be explicitly expressed from (2.130a) as:

$$\ddot{q} = W(q)^{-1} \left[Q + \nabla_q T - \nabla_q V - \nabla_q c \cdot z - \frac{\partial}{\partial q} (W(q) \dot{q}) \dot{q} \right] \quad (2.131)$$

and as previously, for $W(q)$ constant, equation (2.44) reduces to:

$$\ddot{q} = W^{-1} [Q - \nabla_q V - \nabla_q c \cdot z] \quad (2.132)$$

We can observe again here that the Lagrange equation (and their explicit forms (2.131) and (2.132)) deliver the acceleration \ddot{q} as a function of q , \dot{q} , Q and of the Lagrange multipliers z . Hence, a simulation of the model could be produced for given initial conditions $q(0)$, $\dot{q}(0)$, given external forces $Q(t)$ (given at all time) and the Lagrange multipliers $z(t)$. Hence, in order to compute a simulation of the model equations, we need to calculate somehow the (time-varying) Lagrange multipliers z at every time instant of the simulation. We will investigate in the next section how to do that.

Let us now consider a list of examples in order to digest these theoretical concepts.

Examples

Let us go through a few examples of constrained Lagrange modelling and try to gather some intuitions on the meaning of the constrained Lagrange equation (2.126).

1. Let us start with our bowl example of Fig. 2.8 to detail the procedure behind the constrained Lagrange equation (2.126). The mass position is described in the cartesian coordinates $q \equiv p \in \mathbb{R}^3$, such that the kinetic and potential energy function read as:

$$T = \frac{1}{2} m \dot{p}^\top \dot{p}, \quad V = mg \begin{bmatrix} 0 & 0 & 1 \end{bmatrix} p \quad (2.133)$$

Note that we can assemble these functions by completely disregarding the existence of the constraint $c(q) = 0$. We then assemble the constrained Lagrange function as per (2.125):

$$\mathcal{L} = \frac{1}{2}m\dot{p}^\top \dot{p} - mg [\begin{array}{ccc} 0 & 0 & 1 \end{array}] p - z^\top \left(p_3 - \frac{1}{4}(p_2^2 + p_1^2) \right) \quad (2.134)$$

Note that here the constraint is scalar, such that $z \in \mathbb{R}$, and the scalar product $z^\top c(q)$ is a simple product here, i.e. we can rewrite the Lagrange function:

$$\mathcal{L} = \frac{1}{2}m\dot{p}^\top \dot{p} - mg [\begin{array}{ccc} 0 & 0 & 1 \end{array}] p - z \left(p_3 - \frac{1}{4}(p_2^2 + p_1^2) \right) \quad (2.135)$$

and deploy the Lagrange equation as usual:

$$\nabla_{\dot{q}} \mathcal{L} = m\dot{p} \quad (2.136a)$$

$$\frac{d}{dt} \nabla_{\dot{q}} \mathcal{L} = m\ddot{p} \quad (2.136b)$$

$$\nabla_q \mathcal{L} = -mg \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} - \nabla_q c \cdot z, \quad \text{where} \quad \nabla_q c = \begin{bmatrix} -\frac{1}{2}p_1 \\ -\frac{1}{2}p_2 \\ 1 \end{bmatrix} \quad (2.136c)$$

The dynamics of the mass in the bowl are then given by:

$$m\ddot{p} = -mg \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} - z \begin{bmatrix} -\frac{1}{2}p_1 \\ -\frac{1}{2}p_2 \\ 1 \end{bmatrix} \quad (2.137a)$$

$$0 = p_3 - \frac{1}{4}(p_2^2 + p_1^2) \quad (2.137b)$$

As detailed in the theory above, it is worth here briefly underlining that:

- equation (2.137a) delivers the system acceleration \ddot{p} for p and z known
 - equation (2.137b) is scalar, and describes the condition that the position p ought to satisfy in order to "be on" the manifold (i.e. the bowl here). However, one ought to observe that (2.137b) does not deliver the unknown z . As a matter of fact, z does not even appear in (2.137b). Hence, while (2.137) provides 4 equations for the 4 unknown variables \dot{q}, z, p , it does not provide them as such, as it fails to provide z .
2. Let us revisit the crane of Fig. 2.5 in cartesian coordinates, assuming that the link between the cart and the hanging mass is (infinitely) rigid. The kinetic energy function is built in the same way as in the elastic example above, i.e. we use the coordinates

$$p = \begin{bmatrix} p_1 \\ p_2 \end{bmatrix} \in \mathbb{R}^2 \quad (2.138)$$

for position of the hanging mass, and x for the position of the cart. We chose to order our generalized coordinates \mathbf{q} as:

$$\mathbf{q} = \begin{bmatrix} x \\ p \end{bmatrix} \quad (2.139)$$

The kinetic energy reads as (2.84), i.e.:

$$T = \frac{1}{2}M\dot{x}^2 + \frac{1}{2}m\dot{\mathbf{p}}^\top \dot{\mathbf{p}} \quad (2.140)$$

and the W matrix reads as (2.85), i.e.:

$$W = \begin{bmatrix} M & 0 & 0 \\ 0 & m & 0 \\ 0 & 0 & m \end{bmatrix} \quad (2.141)$$

and is constant. The potential energy now involves only the energy from gravity, i.e.

$$V = -mgp_2 \quad (2.142)$$

The constraint is then the distance between the cart and the hanging mass. The vector describing the link is given by:

$$\boldsymbol{\delta} = \mathbf{p} - \begin{bmatrix} x \\ 0 \end{bmatrix} \quad (2.143)$$

And the constraint imposed by the link in the system can e.g. be written as:

$$(\boldsymbol{\delta}^\top \boldsymbol{\delta})^{\frac{1}{2}} - L = 0 \quad (2.144)$$

For computational reasons, it is useful to “get rid” of the square root function in (2.144) and to scale the constraint by $\frac{1}{2}$, i.e. we write:

$$c(\mathbf{q}) = \frac{1}{2}(\boldsymbol{\delta}^\top \boldsymbol{\delta} - L^2) = 0 \quad (2.145)$$

which is equivalent to (2.144). The Lagrange function then reads as:

$$\mathcal{L} = \frac{1}{2}M\dot{x}^2 + \frac{1}{2}m\dot{\mathbf{p}}^\top \dot{\mathbf{p}} + mgp_2 - \frac{1}{2}z(\boldsymbol{\delta}^\top \boldsymbol{\delta} - L^2) \quad (2.146)$$

We can then deploy the Lagrange equation as usual:

$$\nabla_{\dot{\mathbf{q}}} \mathcal{L} = W\dot{\mathbf{q}} \quad (2.147a)$$

$$\frac{d}{dt} \nabla_{\dot{\mathbf{q}}} \mathcal{L} = W\ddot{\mathbf{q}} \quad (2.147b)$$

$$\nabla_{\mathbf{q}} \mathcal{L} = \begin{bmatrix} 0 \\ 0 \\ mg \end{bmatrix} - z \begin{bmatrix} x - p_1 \\ p_1 - x \\ p_2 \end{bmatrix} \quad (2.147c)$$

The model then reads as:

$$W\ddot{\mathbf{q}} - \begin{bmatrix} 0 \\ 0 \\ mg \end{bmatrix} + z \begin{bmatrix} x - p_1 \\ p_1 - x \\ p_2 \end{bmatrix} = 0 \quad (2.148a)$$

$$\frac{1}{2}(\boldsymbol{\delta}^\top \boldsymbol{\delta} - L^2) = 0 \quad (2.148b)$$

In a more explicit form, it reads as:

$$\ddot{\mathbf{q}} = W^{-1} \left(\begin{bmatrix} 0 \\ 0 \\ mg \end{bmatrix} - z \begin{bmatrix} x - p_1 \\ p_1 - x \\ p_2 \end{bmatrix} \right) \quad (2.149a)$$

$$0 = \frac{1}{2}(\boldsymbol{\delta}^\top \boldsymbol{\delta} - L^2) \quad (2.149b)$$

The observations already made above concerning solving (2.149) for z still hold here. It is interesting at this stage to compare (2.149) to its equivalent (2.70) developed in polar coordinates (i.e. using the angle θ to describe the position of the hanging mass). While (2.149) holds more equations than (2.70) (4 vs. 2), its symbolic complexity is lower. Indeed, while (2.70) comprises several trigonometric terms (\sin and \cos), and a division, (2.149) includes only bilinear terms (i.e. products of the variables).

3. We will revisit now the 2-pendulum crane illustrated in Fig. 2.6. Modelling this system using polar coordinates (i.e. the angles of the pendulums) was resulting in a model of very high symbolic complexity (see (2.107)-(2.108) and the following remarks). Let us see the outcome of approaching the modelling of the same system using cartesian coordinates and constrained Lagrange. We select the generalized coordinates:

$$\mathbf{q} = \begin{bmatrix} x \\ p_1 \\ p_2 \end{bmatrix}, \quad p_{1,2} \in \mathbb{R}^2 \quad (2.150)$$

The kinetic energy function reads as:

$$T = \frac{1}{2}m\dot{\mathbf{q}}^\top \dot{\mathbf{q}} \quad (2.151)$$

hence $W = mI$ (where I is the 3×3 identity matrix). The potential energy function reads as:

$$V = -mg [\ 0 \ 0 \ 1 \ 0 \ 1 \] \mathbf{q} \quad (2.152)$$

The constraints then read as:

$$c(\mathbf{q}) = \frac{1}{2} \begin{bmatrix} \boldsymbol{\delta}_1^\top \boldsymbol{\delta}_1 - L^2 \\ \boldsymbol{\delta}_2^\top \boldsymbol{\delta}_2 - L^2 \end{bmatrix} = 0 \quad (2.153)$$

where

$$\boldsymbol{\delta}_1 = \begin{bmatrix} x \\ 0 \end{bmatrix} - \mathbf{p}_1, \quad \boldsymbol{\delta}_2 = \mathbf{p}_2 - \mathbf{p}_1 \quad (2.154)$$

The Lagrange function reads as:

$$\mathcal{L} = \frac{1}{2} m \dot{\mathbf{q}}^\top \dot{\mathbf{q}} + mg \begin{bmatrix} 0 & 0 & 1 & 0 & 1 \end{bmatrix} \mathbf{q} - \frac{1}{2} \mathbf{z}^\top \begin{bmatrix} \boldsymbol{\delta}_1^\top \boldsymbol{\delta}_1 - L^2 \\ \boldsymbol{\delta}_2^\top \boldsymbol{\delta}_2 - L^2 \end{bmatrix} \quad (2.155)$$

where $\mathbf{z} \in \mathbb{R}^2$. We can then compute

$$\nabla_{\dot{\mathbf{q}}} \mathcal{L} = m \dot{\mathbf{q}} \quad (2.156a)$$

$$\frac{d}{dt} \nabla_{\dot{\mathbf{q}}} \mathcal{L} = m \ddot{\mathbf{q}} \quad (2.156b)$$

$$\nabla_{\mathbf{q}} \mathcal{L} = \begin{bmatrix} 0 \\ 0 \\ mg \\ 0 \\ mg \end{bmatrix} - \begin{bmatrix} x - [1 \ 0] \mathbf{p}_1 & 0 \\ -\boldsymbol{\delta}_1 & -\boldsymbol{\delta}_2 \\ 0 & \boldsymbol{\delta}_2 \end{bmatrix} \mathbf{z} \quad (2.156c)$$

The final model in a explicit form then reads as:

$$\ddot{\mathbf{q}} = \begin{bmatrix} 0 \\ 0 \\ g \\ 0 \\ g \end{bmatrix} - \frac{1}{m} \begin{bmatrix} x - [1 \ 0] \mathbf{p}_1 & 0 \\ -\boldsymbol{\delta}_1 & -\boldsymbol{\delta}_2 \\ 0 & \boldsymbol{\delta}_2 \end{bmatrix} \mathbf{z} \quad (2.157)$$

$$0 = \frac{1}{2} \begin{bmatrix} \boldsymbol{\delta}_1^\top \boldsymbol{\delta}_1 - L^2 \\ \boldsymbol{\delta}_2^\top \boldsymbol{\delta}_2 - L^2 \end{bmatrix} \quad (2.158)$$

This ought to be compared to (2.107)-(2.108), where the final model was not even provided because of its very high symbolic complexity. The reduction of symbolic complexity is resulting from the choice of cartesian coordinates, which yields a constant and diagonal W matrix. As a result, the complex terms $\frac{\partial}{\partial \mathbf{q}} (W(\mathbf{q}) \dot{\mathbf{q}}) \dot{\mathbf{q}}$ in

$$\ddot{\mathbf{q}} = W(\mathbf{q})^{-1} \left[Q + \nabla_{\mathbf{q}} T - \nabla_{\mathbf{q}} V - \nabla_{\mathbf{q}} \mathbf{c} \cdot \mathbf{z} - \frac{\partial}{\partial \mathbf{q}} (W(\mathbf{q}) \dot{\mathbf{q}}) \dot{\mathbf{q}} \right] \quad (2.159)$$

disappears, and the inverse $W(\mathbf{q})^{-1}$ is trivial.

2.5.1 HANDLING MODELS FROM CONSTRAINED LAGRANGE

In the previous section we have seen that we can write the accelerations $\ddot{\mathbf{q}}$

$$\ddot{\mathbf{q}} = W^{-1} [Q - \nabla_{\mathbf{q}} V - \nabla_{\mathbf{q}} \mathbf{c} \cdot \mathbf{z}] \quad (2.160)$$

As observed before, in order to compute a simulation from (2.160) we need the (time-varying) Lagrange multipliers z at every time instant of the simulation.

Unfortunately, a problem arises here. Indeed, since we use equation (2.126a) (or its equivalent form (2.130a)) to compute \ddot{q} , we are left with equation (2.166b) to determine the unknown z . However, equation (2.166b) *does not* deliver z , indeed it is not even a function of z , and thus cannot inform us - as such - about z . We will see in the following how to circumvent this problem, and we will put some further formalism around this issue in Section 5.

Running simulations for mechanical system requires one to be able to compute the system accelerations \ddot{q} as a function of the positions and velocities q, \dot{q} and of the external forces Q . In the unconstrained case, this is readily feasible as long as matrix $W(q)$ is full rank.

In the constrained case, we have raised the issue before that a model arising from constrained Lagrange and taking the form:

$$\frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{q}} - \frac{\partial \mathcal{L}}{\partial q} = Q \quad (2.161a)$$

$$c(q) = 0 \quad (2.161b)$$

does not deliver as such the accelerations of the system. Indeed, while equation (2.161a) delivers the acceleration \ddot{q} via its explicit version:

$$\ddot{q} = W(q)^{-1} \left[Q + \nabla_q T - \nabla_q V - \nabla_q c \cdot z - \frac{\partial}{\partial q} (W(q) \dot{q}) \dot{q} \right] \quad (2.162)$$

the accelerations \ddot{q} are function of the unknown z , which is not delivered by (2.161b), as $c(q)$ is not even a function of z . In this section we will see how to tackle this problem.

If one sets $z = 0$ in the dynamic equation (2.161a), then one can verify that the equation would be describing a "free" motion, i.e. what the system would do if the constraints $c(q) = 0$ were not present at all. E.g. in the "bowl" example above, the ball would free fall, and in the 2-pendulums crane example, the hanging masses would free fall and the cart would move along its rail as if not connected to anything. One can therefore construe $\nabla_q c \cdot z$ in (2.162) as a term that will enforce the constraints $c(q) = 0$ by manipulating \ddot{q} resulting from (2.162) via adequate selections of the variables z . As a matter of fact, the combination of terms

$$W(q)^{-1} \nabla_q c \cdot z \quad (2.163)$$

arising in (2.162) yield an acceleration that is added to the other "sources" of accelerations (e.g. stemming from forces, gravity, centrifugal effects, etc.). One can observe that $W(q)^{-1} \nabla_q c$ is a matrix of dimension $n_q \times n_c$, such that the term (2.163) yields accelerations in the subspace spanned by the columns of $W(q)^{-1} \nabla_q c$.

Let us then build some intuition behind what will happen here. We know that z ought to be chosen such that the accelerations \ddot{q} enforce the constraints $c(q) = 0$. However, the constraints specify conditions on the system position, not on its accelerations. We ought to understand, though, that the positions q are obviously influenced by the accelerations \ddot{q} (via 2 integrations), such that the accelerations influence $c(q)$ (via 2 integrations). In order to chose the z adequately, we need to make the impact of the accelerations on $c(q)$ appear explicitly.

Fortunately, this influence is fairly simple to unpack. Indeed, let us take two time derivatives of the constraints $c(q)$ (in order to "rewind" the two integrations from \ddot{q} to q). We should observe that if $c(q) = 0$ is enforced throughout the trajectory of the system (i.e. at every time instant), then:

$$\frac{d^k}{dt^k}c(q) = 0 \quad (2.164)$$

also hold at all time for any $k \geq 0$. We also observe that $\frac{d^2}{dt^2}c(q) = 0$ is a condition where the accelerations appear explicitly. In order to see that, we simply need to apply some chain rules:

$$\dot{c}(q, \dot{q}) = \frac{d}{dt}c(q) = \frac{\partial c}{\partial q}\dot{q} \quad (2.165a)$$

$$\ddot{c}(q, \dot{q}, \ddot{q}) = \frac{d^2}{dt^2}c(q) = \frac{\partial c}{\partial q}\ddot{q} + \frac{\partial}{\partial q}\left(\frac{\partial c}{\partial q}\dot{q}\right)\dot{q} \quad (2.165b)$$

We can then assemble the Lagrange equation (2.126a) and equation (2.165b) in the new model:

$$\frac{d}{dt}\frac{\partial \mathcal{L}}{\partial \dot{q}} - \frac{\partial \mathcal{L}}{\partial q} = Q \quad (2.166a)$$

$$\ddot{c}(q, \dot{q}, \ddot{q}) = 0 \quad (2.166b)$$

Or in the explicit form similar to (2.130):

$$W(q)\ddot{q} + \frac{\partial}{\partial q}(W(q)\dot{q})\dot{q} - \nabla_q T + \nabla_q V + \nabla_q c \cdot z = Q \quad (2.167a)$$

$$\frac{\partial c}{\partial q}\ddot{q} + \frac{\partial}{\partial q}\left(\frac{\partial c}{\partial q}\dot{q}\right)\dot{q} = 0 \quad (2.167b)$$

Recall that our original problem is that the constrained Lagrange equations (2.161) do not deliver the variables z , which are needed in order to compute the accelerations \ddot{q} . In order to solve the problem we have time-differentiated the constraints $c(q) = 0$, and obtained (2.167). However, It is not yet obvious that the modified constraints (2.167b) deliver z . As a matter of fact, they are still not a function of z . Let us investigate whether (2.167) is now capable of delivering the pair \ddot{q}, z . Here we can proceed via algebraic manipulations over

(2.167). However, it is much easier to address the question by observing that (2.167) is linear in \ddot{q} , z . Indeed one can rewrite (2.167) as:

$$\underbrace{\begin{bmatrix} W(q) & \nabla_q c \\ \nabla_q c^\top & 0 \end{bmatrix}}_{:=M(q)} \begin{bmatrix} \ddot{q} \\ z \end{bmatrix} = \begin{bmatrix} Q - \frac{\partial}{\partial q} (W(q)\dot{q})\dot{q} + \nabla_q T - \nabla_q V \\ -\frac{\partial}{\partial q} \left(\frac{\partial c}{\partial q} \dot{q} \right) \dot{q} \end{bmatrix} \quad (2.168)$$

Here we can readily see that (2.168) delivers \ddot{q} , z jointly if matrix $M(q)$ is full rank. One can also write the model in an explicit form:

$$\begin{bmatrix} \ddot{q} \\ z \end{bmatrix} = M(q)^{-1} \begin{bmatrix} Q - \frac{\partial}{\partial q} (W(q)\dot{q})\dot{q} + \nabla_q T - \nabla_q V \\ -\frac{\partial}{\partial q} \left(\frac{\partial c}{\partial q} \dot{q} \right) \dot{q} \end{bmatrix} \quad (2.169)$$

However, as observed previously for matrix $W(q)$ alone, inverting the symbolic matrix $M(q)$ can yield an extremely complex matrix $M(q)^{-1}$, even if $M(q)$ is fairly "simple". For this reason, it is often preferable to work with the model (2.168) in its implicit form, or to treat the matrix inversion $M(q)^{-1}$ numerically when deploying model (2.169).

Examples

Let us revisit some of the examples of Section 2.5 in the light of the model transformation described above.

1. The bowl example was ending with the model equations

$$m\ddot{p} = -mg \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} - z \begin{bmatrix} -\frac{1}{2}p_1 \\ -\frac{1}{2}p_2 \\ 1 \end{bmatrix} \quad (2.170a)$$

$$0 = p_3 - \frac{1}{4}(p_2^2 + p_1^2) \quad (2.170b)$$

where (2.170b) is $c(q) = 0$ written explicitly. In order to perform the model transformation detailed above on this model, we need to time-differentiate (2.170b) twice:

$$c(q) = p_3 - \frac{1}{4}(p_2^2 + p_1^2) \quad (2.171a)$$

$$\dot{c}(q) = \dot{p}_3 - \frac{1}{2}(p_2\dot{p}_2 + p_1\dot{p}_1) \quad (2.171b)$$

$$\ddot{c}(q) = \ddot{p}_3 - \frac{1}{2}(p_2\ddot{p}_2 + p_1\ddot{p}_1 + \dot{p}_2^2 + \dot{p}_1^2) \quad (2.171c)$$

These expressions can be obtained directly as above, or via using the construction (2.165). The transformed model then reads as:

$$m\ddot{p} = -mg \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} - z \begin{bmatrix} -\frac{1}{2}p_1 \\ -\frac{1}{2}p_2 \\ 1 \end{bmatrix} \quad (2.172a)$$

$$0 = \ddot{p}_3 - \frac{1}{2}(p_2\ddot{p}_2 + p_1\ddot{p}_1 + \dot{p}_2^2 + \dot{p}_1^2) \quad (2.172b)$$

It can be put in the form (2.168):

$$\underbrace{\begin{bmatrix} m & 0 & 0 & -\frac{1}{2}p_1 \\ 0 & m & 0 & -\frac{1}{2}p_2 \\ 0 & 0 & m & 1 \\ -\frac{1}{2}p_1 & -\frac{1}{2}p_2 & 1 & 0 \end{bmatrix}}_{:=M(q)} \begin{bmatrix} \ddot{p}_1 \\ \ddot{p}_2 \\ \ddot{p}_3 \\ z \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ mg \\ \dot{p}_2^2 + \dot{p}_1^2 \end{bmatrix} \quad (2.173)$$

The determinant of matrix $M(q)$ reads as:

$$\det(M(q)) = -\frac{m^2}{4}(p_1^2 + p_2^2 + 4) < 0 \quad (2.174)$$

and is non-zero for any position p , such that (2.173) is always well defined. The inverse of matrix $M(q)$ is fairly complex, such that writing (2.173) explicitly is tedious.

2. Let us look at the crane example, which had the original model

$$W\ddot{q} = \begin{bmatrix} 0 \\ 0 \\ mg \end{bmatrix} - z \begin{bmatrix} x - p_1 \\ p_1 - x \\ p_2 \end{bmatrix} \quad (2.175a)$$

$$0 = \frac{1}{2}(\boldsymbol{\delta}^\top \boldsymbol{\delta} - L^2) \quad (2.175b)$$

where $\boldsymbol{\delta} = p - \begin{bmatrix} x \\ 0 \end{bmatrix}$ and

$$W = \begin{bmatrix} M & 0 & 0 \\ 0 & m & 0 \\ 0 & 0 & m \end{bmatrix} \quad (2.176)$$

We can then proceed with time-differentiating the constraints:

$$c(q) = \frac{1}{2}(\boldsymbol{\delta}^\top \boldsymbol{\delta} - L^2) \quad (2.177a)$$

$$\dot{c}(q) = \boldsymbol{\delta}^\top \dot{\boldsymbol{\delta}} \quad (2.177b)$$

$$\ddot{c}(q) = \boldsymbol{\delta}^\top \ddot{\boldsymbol{\delta}} + \dot{\boldsymbol{\delta}}^\top \dot{\boldsymbol{\delta}} \quad (2.177c)$$

The model put in the form (2.168) then reads as:

$$\underbrace{\begin{bmatrix} W & x - p_1 & p_1 - x & p_2 \\ x - p_1 & p_1 - x & p_2 & 0 \end{bmatrix}}_{:=M(q)} \begin{bmatrix} \ddot{x} \\ \ddot{p} \\ z \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ mg \\ \dot{x}^2 + \dot{p}_1^2 + \dot{p}_2^2 - 2\dot{x}\dot{p}_1 \end{bmatrix} \quad (2.178)$$

We observe that the determinant of matrix $M(q)$ is non-zero as long as $\begin{bmatrix} x \\ 0 \end{bmatrix} \neq p$,

otherwise the last row and last column of $M(q)$ are zero, and the matrix becomes rank-deficient. This situation corresponds, physically, to the hanging mass coinciding with the cart (i.e. the length of the link is zero). The inverse of matrix $M(q)$ is very complex, such that writing (2.178) explicitly is not recommended.

□

2.6 CONSISTENCY CONDITIONS

Before finishing this Section, we need to approach an important consequence of the model transformation detailed above in terms of model simulation. The original model that is describing the physical system reads as:

$$\frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{q}} - \frac{\partial \mathcal{L}}{\partial q} = Q \quad (2.179a)$$

$$c(q) = 0 \quad (2.179b)$$

while the transformed model reads as:

$$\frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{q}} - \frac{\partial \mathcal{L}}{\partial q} = Q \quad (2.180a)$$

$$\ddot{c}(q, \dot{q}, \ddot{q}) = 0 \quad (2.180b)$$

An important question to ask here is whether the transformed model (2.180) is equivalent to the original model (2.179). Since the trajectories arising from the original model (2.179) satisfy $c = 0$ at all time, they also satisfies (2.180b) at all time. Hence, the trajectories of the original model are also trajectories of the transformed model (2.180). Unfortunately, the converse is not always true. Indeed, the trajectories of the transformed model (2.179) satisfy $\ddot{c} = 0$ at all time, which does not entail that they satisfy $c = 0$ at all time. Actually, if the trajectories satisfy $\ddot{c} = 0$ at all time, then one can verify that the constraints c must obey the time evolution:

$$c(q(t)) = c(q(0)) + t \cdot \dot{c}(q(0), \dot{q}(0)) \quad (2.181)$$

such that $c(q(t)) = 0$ does not necessarily hold. Fortunately, (2.181) delivers the conditions that are required in order for the trajectories of the transformed model (2.180) to be trajectories of the original model (2.179). Indeed, (2.181) readily tells us that $c(q) = 0$ holds at all time if:

$$C(q(0), \dot{q}(0)) := \begin{bmatrix} c(q(0)) \\ \dot{c}(q(0), \dot{q}(0)) \end{bmatrix} = 0 \quad (2.182a)$$

hence, the transformed model delivers meaningful trajectories if the initial conditions $q(0), \dot{q}(0)$ satisfy $C(q(0), \dot{q}(0)) = 0$. Condition (2.182a) is referred to as the **consistency conditions**

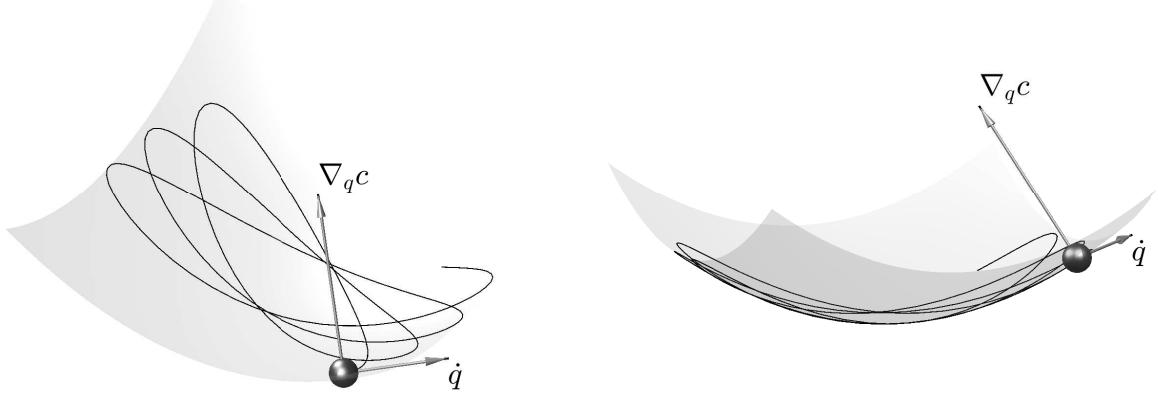


Figure 2.9: Consistency conditions for the bowl example. The condition $\nabla_q c^\top \dot{q} = 0$ simply requires that $\dot{q}(0)$ is tangent to the bowl.

associated to the model transformation. They are required in order to impose initial conditions in the model that are physically meaningful.

In order to build some intuition behind these observation using the bowl example above. The consistency conditions for the bowl read as:

$$C(q(0), \dot{q}(0)) := \begin{bmatrix} p_3 - \frac{1}{4}p_1^2 - \frac{1}{4}p_2^2 \\ \dot{p}_3 - \frac{1}{2}\dot{p}_1 p_1 - \frac{1}{2}\dot{p}_2 p_2 \end{bmatrix} \quad (2.183)$$

While these expressions do not carry much intuition, they have a very strong geometric interpretation. The first condition in C simply states that the initial condition $q(0)$ must satisfy the constraints $c(q)$, and is a quite natural requirement. The second condition states that:

$$\dot{c}(q(0), \dot{q}(0)) = \frac{\partial c}{\partial q} \dot{q} = \nabla_q c^\top \dot{q} = 0 \quad (2.184)$$

i.e. it requires that the scalar product between the gradient $\nabla_q c$ and the velocities \dot{q} is zero. In other words, it requires the velocities \dot{q} to be orthogonal to the gradient $\nabla_q c$. This requirement has a very intuitive interpretation in the bowl example, see Fig. 2.9. Indeed, the gradient $\nabla_q c$ is describing a normal to the surface described by the equation $c(q(0)) = 0$. In our bowl example, this surface is the bowl itself (more complex systems may not deliver such simple picture), and the condition $\nabla_q c(q(0))^\top \dot{q}(0) = 0$ essentially requires that the initial velocities $\dot{q}(0)$ are tangent to the bowl surface. This requirement is physically needed in order for the ball to slide on the surface of the bowl.

2.7 CONSTRAINTS DRIFT

We have seen in the previous section that the initial conditions of the transformed model must satisfy some consistency conditions in order for the transformed model to deliver

meaningful trajectories. The consistency conditions can be derived from the observation that the constraints $c(q)$ obey the time evolution (2.181), i.e.:

$$c(q(t)) = c(q(0)) + t \cdot \dot{c}(q(0), \dot{q}(0)) \quad (2.185)$$

on the trajectories $q(t)$ computed from the transformed model. However, this time evolution holds if and only if the condition $\ddot{c} = 0$ is imposed. While this condition is imposed by the transformed model (2.166) (and all the variants we have looked at), when simulating the system, the condition is imperfectly imposed for numerical reasons. As a result, the time evolution (2.185) is "noisy", and it tends to accumulate the numerical errors over time. This effect is illustrated in Fig. 2.10, which is the outcome of simulating the transformed model (2.172) with consistent initial conditions, but a low integration accuracy. As a result, the condition $\ddot{c} = 0$ is not enforced very accurately, and tends to create a drift in the constraints.

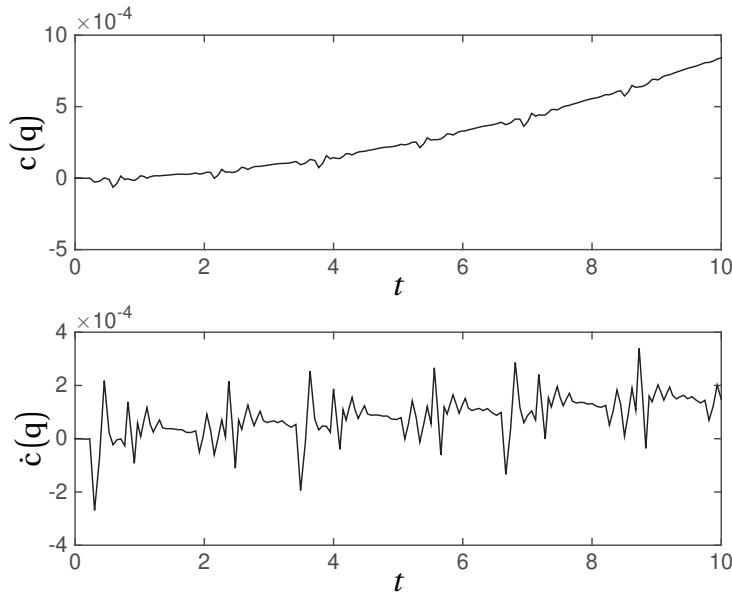


Figure 2.10: Illustration of the constraints drift for the bowl example.

This issue can be treated via **Baumgarte stabilization**, whereby the transformed model does not impose $\ddot{c} = 0$, but rather dynamics on c that stabilizes them to zero. This can be e.g. done by imposing:

$$\ddot{c} + 2\alpha\dot{c} + \alpha^2 c = 0 \quad (2.186)$$

for some $\alpha > 0$ (such that the dynamics (2.186) are stable, with two real poles at $-\alpha$). The transformed model then reads as:

$$\frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{q}} - \frac{\partial \mathcal{L}}{\partial q} = Q \quad (2.187a)$$

$$\ddot{c} + 2\alpha\dot{c} + \alpha^2 c = 0 \quad (2.187b)$$

3 NEWTON METHOD

The Newton method is a central tool for simulation, optimization, and estimation. It is therefore no surprise that we will encounter it here. Generally speaking, the Newton method aims at solving a set of equations, that we can e.g. write as:

$$\boldsymbol{\varphi}(x, y) = 0 \quad (3.1)$$

in $x \in \mathbb{R}^{n_x}$ for a given $y \in \mathbb{R}^{n_y}$. Clearly, in order for the problem of finding x to be well-posed, one needs (3.1) to hold “as many equations as unknown variables”, i.e. function $\boldsymbol{\varphi}$ is

$$\boldsymbol{\varphi} : \mathbb{R}^{n_x} \times \mathbb{R}^{n_y} \mapsto \mathbb{R}^{n_x} \quad (3.2)$$

We will see later that this condition is not sufficient. When the system of equations described by (3.1) is nonlinear, finding a solution x can typically not be done explicitly (i.e. we typically cannot provide a set of symbolic expressions that describe x as a function of y). This does not mean that x is not perfectly well described as a function of y by (3.1). Formally, one says that x is an **implicit** function of y defined by (3.1).

The Newton method then allows one to compute x as a function of y numerically. In this chapter we will introduce the basics of the Newton method.

3.1 BASIC IDEA OF NEWTON METHOD

The Newton method replaces the problem of solving the nonlinear set of equations (3.1) by solving a succession of surrogate linear approximations. These linear approximations are built by successive linearizations of (3.1), i.e. instead of solving (3.1), we solve:

$$\boldsymbol{\varphi}(x_+, y) \approx \boldsymbol{\varphi}(x, y) + \frac{\partial \boldsymbol{\varphi}(x, y)}{\partial x} (x_+ - x) = 0 \quad (3.3)$$

for x_+ , based on a given guess x . Note that (3.3) is a first-order Taylor approximation of $\boldsymbol{\varphi}(x, y)$, i.e.

$$\boldsymbol{\varphi}(x_+, y) = \boldsymbol{\varphi}(x, y) + \frac{\partial \boldsymbol{\varphi}(x, y)}{\partial x} (x_+ - x) + \mathcal{O}(\|x_+ - x\|^2) \quad (3.4)$$

and is an affine form in x_+ , i.e. x_+ can be obtained from (3.3) explicitly from:

$$x_+ = x - \frac{\partial \boldsymbol{\varphi}(x, y)}{\partial x}^{-1} \boldsymbol{\varphi}(x, y) \quad (3.5)$$

It is common to label

$$\Delta x = - \frac{\partial \boldsymbol{\varphi}(x, y)}{\partial x}^{-1} \boldsymbol{\varphi}(x, y) \quad (3.6)$$

as a (full) *Newton step* that corrects the guess x to the update x_+ . See Fig 3.1 for an illustration. Using this construction, the Newton method is then *iterative*, in the sense that the (3.5) is repeated until the solution to $\varphi(x, y) = 0$ is reached (or close enough). More specifically, the Newton method uses the algorithm:

Algorithm: Full-step Newton method

Input: Variable y , initial guess x , and tolerance tol

while $\|\varphi(x, y)\|_\infty \geq \text{tol}$ **do**

 Compute

$$\varphi(x, y) \quad \text{and} \quad \frac{\partial \varphi(x, y)}{\partial x} \quad (3.7)$$

 Compute the Newton step

$$\frac{\partial \varphi(x, y)}{\partial x} \Delta x + \varphi(x, y) = 0 \quad (3.8)$$

 Take the Newton step

$$x \leftarrow x + \Delta x \quad (3.9)$$

return x

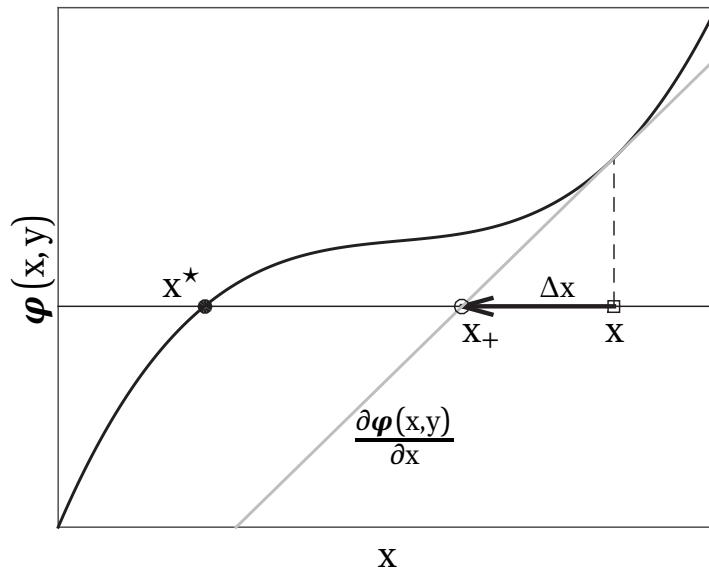


Figure 3.1: Illustration of the Newton principle on a scalar equation (i.e. $n_x = 1$ and y is not used here). The update x_+ is obtained by finding the root of the linear approximation (grey line) of $\varphi(x, y)$ built at a candidate point x . The update x_+ is closer to the solution x^* of $\varphi(x, y) = 0$ than the candidate x .

A few crucial remarks ought to be made here:

- If it converges, the Newton method converges to x^* a solution of $\varphi(x, y) = 0$
- Each step of the Newton method requires evaluating the function $\varphi(x, y)$ and its Jacobian $\frac{\partial \varphi(x)}{\partial x}$, and solving the linear system (3.8).
- The Newton method requires that the square Jacobian matrix $\frac{\partial \varphi(x, y)}{\partial x}$ is full rank (i.e. invertible), in order for the linear system (3.8) to be well-posed.
- If function $\varphi(x, y)$ is linear in x (and well posed), then the Newton method finds the solution x^* in one step. It is then fully equivalent to solving the linear system.

3.2 CONVERGENCE OF THE NEWTON METHOD

A crucial question for the Newton method is whether the Newton algorithm converges or not. If it converges, then it has to converge to a solution³ x^* of $\varphi(x, y) = 0$. Indeed, as long as $\varphi(x, y) \neq 0$, the Newton algorithm performs corrections $\Delta x \neq 0$, and only stops correcting the initial guess when $\varphi(x, y) \approx 0$ (depending on the tolerance). Unfortunately, the convergence of the Newton method is not trivial. We will discuss some basic results in this section.

3.2.1 CONVERGENCE RATE

Let us first assume that the Newton iteration presented above converges, and let us investigate its convergence rate, i.e how quickly $\|x - x^*\|$ decreases. This question is addressed as follows:

Theorem 3. *if the full step Newton iteration converges, then it converges at a quadratic rate, i.e.*

$$\|x_+ - x^*\| \leq C \cdot \|x - x^*\|^2 \quad (3.10)$$

for some constant $C > 0$.

Before delivering a formal proof of the quadratic convergence rate, one ought to understand that this result is fairly intuitive. Indeed, one can observe that at every iteration, the error that the Newton step is making is of order 2. Hence, every iteration “removes” the first-order inaccuracy between the solution candidate and the solution to $\varphi(x, y) = 0$. It follows that at every step, the Newton method retains an error of order $\|x - x^*\|^2$. This observation alone give an intuition of why (3.10) ought to be true. We provide next a proof of Theorem 3. This proof is not part of the examination though.

³Note that there may be more than one solution! Hence the Newton method could converge to different points, depending on where it is started.

Proof. for the sake of simplicity, let us dismiss the fixed argument y , and let us write the Newton step as

$$\Delta x = -M^{-1}\varphi(x) \quad (3.11)$$

We first observe that since $\varphi(x^*) = 0$, the equality

$$x_+ - x^* = x - x^* - M^{-1}(\varphi(x) - \varphi(x^*)) \quad (3.12)$$

holds. We then use classical results from analysis to observe that:

$$\varphi(x) - \varphi(x^*) = \int_0^1 \frac{\partial \varphi(x + \tau(x^* - x))}{\partial x} (x - x^*) d\tau \quad (3.13)$$

Using (3.12) and (3.13), we can write:

$$x_+ - x^* = x - x^* - M^{-1} \left(\int_0^1 \frac{\partial \varphi(x + \tau(x^* - x))}{\partial x} (x - x^*) d\tau \right) \quad (3.14)$$

Or equivalently

$$x_+ - x^* = M^{-1} \left(M - \int_0^1 \frac{\partial \varphi(x + \tau(x^* - x))}{\partial x} d\tau \right) (x - x^*) \quad (3.15)$$

We can modify (3.15) as:

$$x_+ - x^* = M^{-1} \left(M - \frac{\partial \varphi(x)}{\partial x} - \int_0^1 \left[\frac{\partial \varphi(x + \tau(x^* - x))}{\partial x} - \frac{\partial \varphi(x)}{\partial x} \right] d\tau \right) (x - x^*) \quad (3.16)$$

We can then bound (using triangular inequalities):

$$\begin{aligned} \|x_+ - x^*\| &\leq \left\| M^{-1} \left(M - \frac{\partial \varphi(x)}{\partial x} \right) \right\| \|x - x^*\| \\ &\quad + \left\| M^{-1} \int_0^1 \left[\frac{\partial \varphi(x + \tau(x^* - x))}{\partial x} - \frac{\partial \varphi(x)}{\partial x} \right] d\tau \right\| \|x - x^*\| \end{aligned} \quad (3.17)$$

Since $M = \frac{\partial \varphi(x)}{\partial x}$, the first term is zero. Moreover, the term under the integral sign in the second term can be bounded (on any closed set around x^*) by

$$\left\| M^{-1} \left(\frac{\partial \varphi(x)}{\partial x} - \frac{\partial \varphi(x^*)}{\partial x} \right) \right\| \leq c \|x - x^*\| \quad (3.18)$$

for some constant $c > 0$ for φ smooth. It follows that

$$\begin{aligned} \left\| M^{-1} \int_0^1 \left[\frac{\partial \varphi(x + \tau(x^* - x))}{\partial x} - \frac{\partial \varphi(x)}{\partial x} \right] d\tau \right\| &\leq \int_0^1 \left\| M^{-1} \left[\frac{\partial \varphi(x + \tau(x^* - x))}{\partial x} - \frac{\partial \varphi(x)}{\partial x} \right] \right\| d\tau \\ &\leq \int_0^1 c \|x - x^*\| \tau d\tau = \frac{c}{2} \|x - x^*\| \end{aligned} \quad (3.19)$$

such that

$$\|x_+ - x^*\| \leq \frac{c}{2} \|x - x^*\|^2 \quad (3.20)$$

□

Let us illustrate this convergence rate for the example provided in Fig. 3.3. A few remarks can be useful here:

- A quadratic convergence rate is a very strong convergence. It means that at every iteration the number of accurate digits is doubled, i.e. over the iterations, the error $\|x - x^*\|$ follows a decay of e.g. the form

$$\begin{aligned}\|x - x^*\| &= 10^{-1} \\ \|x - x^*\| &= 10^{-2} \\ \|x - x^*\| &= 10^{-4} \\ \|x - x^*\| &= 10^{-8} \\ \|x - x^*\| &= 10^{-16}\end{aligned}$$

Hence a few iterations are typically sufficient to reach machine precision ($\epsilon = 1e^{-16}$)

- While (3.18) appears very technical, it can be understood as a “measure” of how nonlinear the function $\varphi(x)$ is. Indeed, if $\varphi(x)$ is linear in x , i.e.

$$\varphi(x) = Ax + b \quad (3.21)$$

(we dismiss the argument y here) for a matrix A and a vector b , then

$$\frac{\partial \varphi(x)}{\partial x} = A, \quad \forall x \quad (3.22)$$

and the left-hand side of (3.18) is zero, such that $c = 0$. Conversely, if $\varphi(x)$ is strongly nonlinear, then its Jacobian $\frac{\partial \varphi(x)}{\partial x}$ varies a lot and the difference

$$\frac{\partial \varphi(x + \tau(x^* - x))}{\partial x} - \frac{\partial \varphi(x)}{\partial x} \quad (3.23)$$

can be very large, yielding a large constant c .

- The proof above also informs us on the condition required for the Newton iteration to converge. Indeed, in order for the bound (3.10) to guarantee the decay of $\|x - x^*\|$, one ought to require that the initial guess x provided to the algorithm is such that:

$$\|x - x^*\| \leq 2c^{-1} \quad (3.24)$$

It is interesting here to understand (3.24) properly. Indeed, it tells us that

1. in order for the full-step Newton iteration to converge, it ought to be provided with an initial guess that is **close enough** to a solution x^* . This observation is illustrated in Fig. 3.2.
2. “how far” the initial guess provided to the full-step Newton iteration can be from a solution x^* is larger if c is small, i.e. if $\varphi(x)$ is close to linear. Conversely, if $\varphi(x)$ is very nonlinear, then c is large and the initial guess provided to the full-step Newton iteration must be close to a solution x^* in order to guarantee the convergence of the iteration

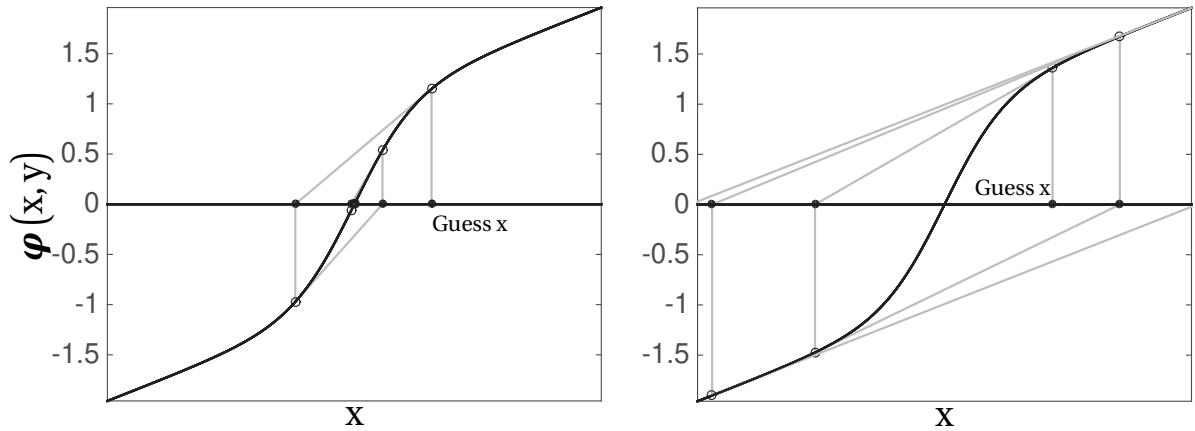


Figure 3.2: Newton iteration on a nonlinear, scalar function $\varphi(x, y)$ (five steps are displayed here). On the left graph the iteration converges to the solution x^* of $\varphi(x^*, y) = 0$, while on the right graph the iteration diverges. Convergence is obtained for the initial guess provided to the full-step Newton iteration being close enough to a solution x^* , while an initial guess further away (possibly) results from an unstable iteration.

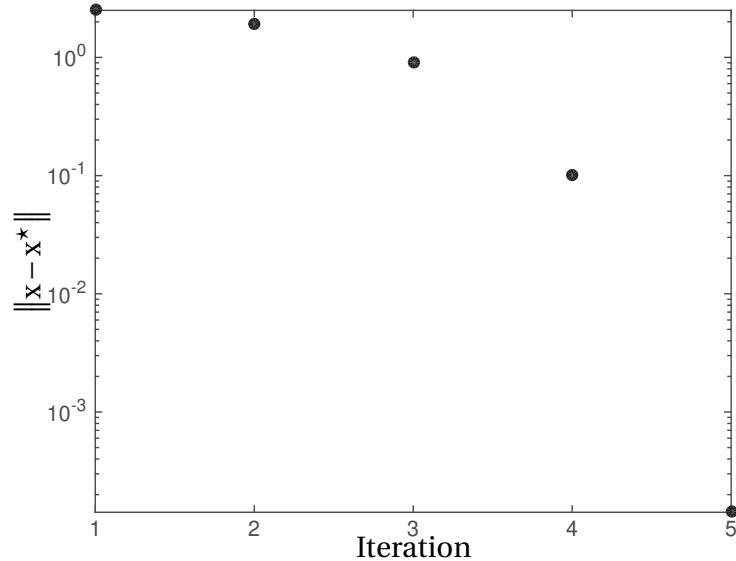


Figure 3.3: Converging full-step Newton iteration. The quadratic convergence rate yields a decrease of the solution error $\|x - x^*\|$ that “collapses” downward in a semi-log plot.

3.2.2 REDUCED NEWTON STEPS

The full-step Newton algorithm as presented above can diverge if the initial guess provided to the Newton iteration is “too far” from a solution x^* of $\varphi(x, y) = 0$. Fortunately, this shortcoming can be addressed using reduced Newton steps. In this section we present the justification and deployment of the approach,

The motivation behind reduced Newton steps can be stated as follows.

Theorem 4. *if it exists, the Newton step Δx solution of the linear system*

$$\frac{\partial \varphi(x, y)}{\partial x} \Delta x + \varphi(x, y) = 0 \quad (3.25)$$

is a descent direction for function $\varphi(x, y)$, i.e.

$$\|\varphi(x + t \cdot \Delta x, y)\| < \|\varphi(x, y)\| \quad (3.26)$$

for some $t \in]0, 1]$.

Note that the result above is valid for any choice of norm $\|\cdot\|$ (this follows directly from the equivalence between norms), but it is easiest to prove using the 2-norm.

Proof. inequality (3.26) is equivalent to⁴

$$\|\varphi(x + t \cdot \Delta x, y)\|^2 - \|\varphi(x, y)\|^2 < 0 \quad (3.27)$$

The first-order Taylor expansion of (3.27) yields:

$$\frac{d}{dt} \|\varphi(x + t \cdot \Delta x, y)\|^2 \Big|_{t=0} \cdot t + \mathcal{O}(t^2) < 0 \quad (3.28)$$

such that (3.27) holds for $t > 0$ small enough if

$$\frac{d}{dt} \|\varphi(x + t \cdot \Delta x, y)\|^2 \Big|_{t=0} < 0 \quad (3.29)$$

We observe that

$$\begin{aligned} \frac{d}{dt} \|\varphi(x + t \cdot \Delta x, y)\|^2 \Big|_{t=0} &= \frac{d}{dt} (\varphi(x + t \cdot \Delta x, y)^\top \varphi(x + t \cdot \Delta x, y)) \Big|_{t=0} \\ &= \varphi(x, y)^\top \frac{d}{dt} \varphi(x + t \cdot \Delta x, y) \Big|_{t=0} \end{aligned} \quad (3.30)$$

$$= \varphi(x, y)^\top \frac{\partial \varphi(x, y)}{\partial x} \Delta x \quad (3.31)$$

$$= -\varphi(x, y)^\top \varphi(x, y) < 0 \quad (3.32)$$

if $\|\varphi(x, y)\| > 0$. □

⁴squaring the 2-norm is helpful here, as it makes it smooth

An important consequence of this result is that while the Newton method presented above may diverge, a careful selection of *reduced Newton steps*, i.e. modifications of x in the direction of the Newton step Δx but possibly scaled down must necessarily converge, as long as the Newton steps Δx exist. This motivates the modification of the full-step Newton algorithm into:

Algorithm: Newton method with reduced steps

Input: Variable y , initial guess x , and tolerance tol

while $\|\varphi(x, y)\|_\infty \geq \text{tol}$ **do**

 Compute

$$\varphi(x, y) \quad \text{and} \quad \frac{\partial \varphi(x, y)}{\partial x} \quad (3.33)$$

 Compute the Newton step

$$\frac{\partial \varphi(x, y)}{\partial x} \Delta x + \varphi(x, y) = 0 \quad (3.34)$$

 Select step size $t \in]0, 1]$

 Take Newton step

$$x \leftarrow x + t \Delta x \quad (3.35)$$

return $x \approx x^*$

Finding the adequate step size $t \in]0, 1]$ is typically performed via a *line-search* strategy, based on testing first a full step ($t = 1$), and then reducing it if the condition (3.26) is not met. Clever methods exist in order to decide whether a step-size t is "good enough", or too short or too long. Other methods also exist in order to adjust the Newton step Δx and guarantee the convergence (see e.g. *trust-region techniques*), but they are more difficult to describe and implement.

A Newton iteration equipped with the possibility of taking reduced steps is guaranteed to converge to a solution x^* of $\varphi(x, y) = 0$ as long as the Newton steps Δx exist throughout the iterations. It shall be underlined here that when reduced steps (i.e. using $t < 1$) are necessary, then the quadratic convergence rate detailed in Theorem 3 is lost. The Newton iteration with reduced steps has the following behavior:

- when x is "far" from a solution x^* , then reduced steps ($t < 1$) are necessary, and the algorithm converges "slowly". The resulting convergence rate can be very poor, even though it is often close to linear in practice
- after a certain amount of iteration, x becomes close enough to x^* and full steps ($t = 1$) are acceptable. The convergence then becomes quadratic. Once x is close to x^* , full

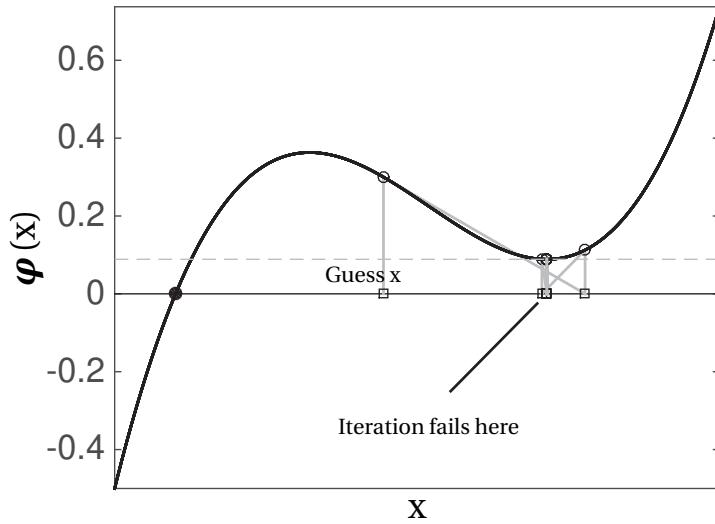


Figure 3.4: Newton iteration with reduced steps on a nonlinear, scalar function $\varphi(x)$ (five steps are displayed here). Here the iteration does not diverge, but it fails at a point where $\frac{\partial \varphi(x,y)}{\partial x} = 0$. At this point, the linear system (3.8) does not have a well-defined solution and the Newton step Δx ceases to exist.

steps can be taken all the way to x^* and the convergence is very fast.

The Newton algorithm with reduced steps converges globally to a solution x^* unless the linear system (3.34) becomes ill-posed. In order to develop some intuition on this potential (and frequent) issue, let us consider the example illustrated in Fig. 3.4. A simple inspection reveals that the Newton iteration may converge to a point at which $\frac{\partial \varphi(x,y)}{\partial x} = 0$, where the Newton step ceases to exist. In such a case, Theorem 4 ceases to apply, and the Newton iteration fails. One can also readily see that if the Newton iteration was started closer to the solution x^* (black dot in the graph), then it would converge to x^* .

3.3 IMPLICIT FUNCTION THEOREM

We have presented the Newton method as a numerical approach to solve the equations

$$\varphi(x,y) = 0 \quad (3.36)$$

in terms of x . In many cases, the equations are dependent on some other argument, labelled y here, which in practice can gather e.g. parameters or data appearing in the equations. If one changes the y , it is natural to expect that the solution x to $\varphi(x,y) = 0$ is affected. Formally, we ought to think of x as being implicitly a function of y , i.e. we should consider x as a function $x(y)$ that we usually cannot write explicitly, but which is implicitly defined by

$$\varphi(x(y),y) = 0 \quad (3.37)$$

Some questions naturally arise then, such as e.g. is $x(y)$ well-defined, is it differentiable with respect to y and what is its derivative. These questions are object of the Implicit Func-

tion Theorem (IFT), which is one of the cornerstones in the field of numerical analysis.

Theorem 5. (*IFT, simplified version*): let function $\varphi(x, y)$ be smooth, and consider a point \bar{x}, \bar{y} such that $\varphi(\bar{x}, \bar{y}) = 0$. Suppose that the Jacobian

$$\frac{\partial \varphi(x, y)}{\partial x} \Big|_{x=\bar{x}, y=\bar{y}} \quad (3.38)$$

is full rank. Then there exists an open set Y around the point \bar{y} in which there exist a unique, smooth function $x(y)$ satisfying:

$$\varphi(x(y), y) = 0, \quad \forall y \in Y \quad (3.39)$$

Moreover, the Jacobian of function $x(y)$ is given by

$$\frac{\partial x(y)}{\partial y} = -\frac{\partial \varphi(x, y)}{\partial x}^{-1} \frac{\partial \varphi(x, y)}{\partial y} \Big|_{x=x(y), y} \quad (3.40)$$

The proof of the IFT is fairly involved and we will not do it here. However, equality (3.40) is trivial to prove.

Proof. (of (3.40)) Because of (3.39), the equality:

$$\frac{d}{dy} \varphi(x(y), y) = 0, \quad \text{holds} \quad \forall y \in Y \quad (3.41)$$

Using the chain rule, we observe that:

$$\frac{\partial \varphi(x, y)}{\partial x} \frac{\partial x(y)}{\partial y} + \frac{\partial \varphi(x, y)}{\partial y} \Big|_{x=x(y), y} = 0 \quad (3.42)$$

which yields (3.40). \square

The importance of the IFT in numerical analysis can hardly be understated, but let us briefly review what it provides us:

- even though function $x(y)$ cannot be written explicitly and can only be evaluated numerically (using a Newton iteration), its derivative is trivial to compute, using (3.40) at the point x found by the Newton iteration for given point y .
- the IFT guarantees that if $\varphi(x, y) = 0$ is well posed (i.e. the Jacobian (3.38) is full rank) at a point y , then it has a neighborhood where it is also well posed, i.e. we can “move” y (locally), and x “moves” accordingly, and also locally.

- One ought to observe that the assumption of the Jacobian $\frac{\partial \varphi(\bar{x}, \bar{y})}{\partial x}$ being full rank is also required for the Newton iteration to be well-behaved (i.e. for the linear system (3.34) to be well-posed). In other words, if the Newton iteration is well-behaved, then the IFT holds. Conversely, the assumption of the IFT must hold in order for a system of equation $\varphi(x, y) = 0$ to be solvable for x .

We will use the IFT in several places in the following.

3.4 JACCOBIAN APPROXIMATION

In some applications of the Newton iteration, the evaluation of the Jacobian $\frac{\partial \varphi(x)}{\partial x}$ is very expensive. It can then be useful to consider using an approximation that is less expensive to evaluate. Let us label this approximation:

$$M \approx \frac{\partial \varphi(x)}{\partial x} \quad (3.43)$$

The resulting *Newton-type* step reads as:

$$\Delta x = -M^{-1} \varphi(x) \quad (3.44)$$

The use of an approximate Jacobian has an impact on the convergence of the Newton method. Let us briefly review how it is changed. In order to do that, we will revisit Theorem 3 and 4 with the Jacobian approximation.

Theorem 6. *the convergence of the full-step Newton method with an approximate Jacobian follows:*

$$\|x_+ - x^*\| \leq \left(\kappa + \frac{c}{2} \|x - x^*\| \right) \|x - x^*\| \quad (3.45)$$

for some constants $c, \kappa > 0$.

Proof. We can re-use (3.17), repeated here for simplicity:

$$\begin{aligned} \|x_+ - x^*\| &\leq \left\| M^{-1} \left(M - \frac{\partial \varphi(x)}{\partial x} \right) \right\| \|x - x^*\| \\ &\quad + \left\| M^{-1} \int_0^1 \left[\frac{\partial \varphi(x + \tau(x^* - x))}{\partial x} - \frac{\partial \varphi(x)}{\partial x} \right] d\tau \right\| \|x - x^*\| \end{aligned} \quad (3.46)$$

and we observe that if $M \neq \frac{\partial \varphi(x)}{\partial x}$ the first term in the right-hand side of the inequality does not disappear. We can then bound (on any closed set around x^*):

$$\left\| M^{-1} \left(M - \frac{\partial \varphi(x)}{\partial x} \right) \right\| = \left\| I - M^{-1} \frac{\partial \varphi(x)}{\partial x} \right\| \leq \kappa, \quad \forall x \quad (3.47)$$

while the second term in the right-hand side of (3.46) can be bounded as in Theorem 3. Bound (3.45) follows. \square

It is useful here to make some remarks concerning Theorem 6.

- The bound (3.45) predicts a linear convergence rate, due to the term $\kappa \|x - x^*\|$ in (3.45). A linear convergence rate is significantly slower than a quadratic convergence rate, as it does not go through a “collapse” to very small values. See Fig. 3.5 for an illustration.

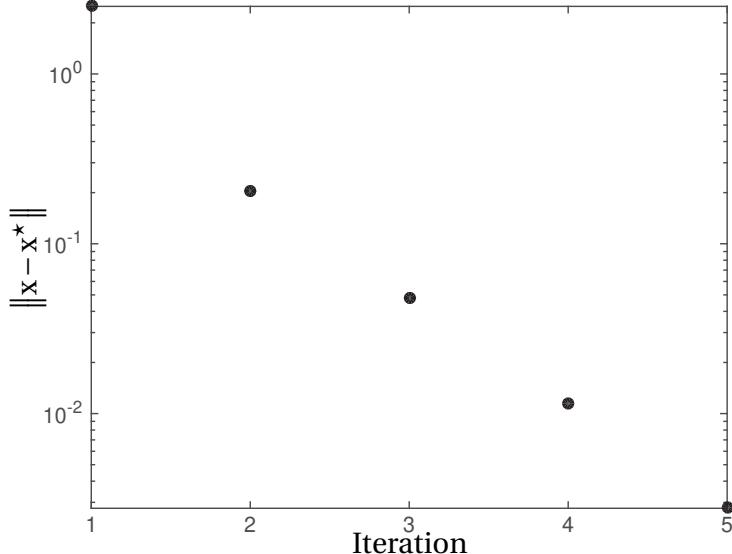


Figure 3.5: Stable full-step Newton iteration with approximate Jacobian. The linear convergence rate yields a decrease of the solution error $\|x - x^*\|$ that follows a “line” downward in a semi-log plot. This is to be contrasted with Fig. 3.3.

- Bound (3.47) informs us on the impact of the error between M and the true Jacobian $\frac{\partial \varphi(x)}{\partial x}$. Indeed, for $M = \frac{\partial \varphi(x)}{\partial x}$, $\kappa = 0$ can be selected, and one recovers a quadratic convergence (see Th. 3). If M and $\frac{\partial \varphi(x)}{\partial x}$ are “very different”, then κ has to be large. One can observe from (3.45) that $\kappa < 1$ must hold in order for the convergence of the full-step quasi-Newton iteration to be guaranteed (even for x being arbitrarily close to x^*).

3.5 NEWTON METHODS FOR UNCONSTRAINED OPTIMIZATION

An important application of the Newton method is for solving minimization problems. In this chapter, we will consider the case of unconstrained optimization, i.e. problems of the form:

$$\min_x \Phi(x, y) \quad (3.48)$$

for a scalar function $\Phi : \mathbb{R}^{n_x} \times \mathbb{R}^{n_y} \rightarrow \mathbb{R}$. If function Φ is smooth with respect to x , the solution x^* to (3.48) is given by solving:

$$\nabla_x \Phi(x, y) = 0 \quad (3.49)$$

for x . The Newton methods detailed above then apply directly to solving (3.49). More specifically, a Newton iteration for solving (3.49) reads as:

Algorithm: Newton method for optimization

Input: Variable y , initial guess x , and tolerance tol

while $\|\nabla_x \Phi(x, y)\|_\infty \geq \text{tol}$ **do**

Compute

$$\nabla_x \Phi(x, y) \quad \text{and} \quad \nabla_x^2 \Phi(x, y) \quad (3.50)$$

Compute the Newton step

$$\nabla_x^2 \Phi(x, y) \Delta x + \nabla_x \Phi(x, y) = 0 \quad (3.51)$$

Select step size $t \in]0, 1]$

Take Newton step

$$x \leftarrow x + t \Delta x \quad (3.52)$$

return $x \approx x^*$

Let us review the results presented in the previous sections in the specific context of optimization.

- The convergence of Newton for Optimization is quadratic in a neighborhood of the solution x^* and “slow” otherwise. If the Hessian $\nabla_x^2 \Phi(x, y)$ is approximated by M , then the method has a linear convergence rate (i.e. Theorem 6 applies).
- The iteration can fail at a point where (3.51) is ill-defined, i.e. at a point where the Hessian of Φ is rank-deficient.

In the context of optimization problems, an additional remark can be made concerning the choice of Hessian approximation, summarized hereafter.

Theorem 7. *if $M > 0$ (positive-definite), then the quasi-Newton step*

$$\Delta x = -M^{-1} \nabla_x \Phi(x, y) \quad (3.53)$$

is a descent direction for the cost function $\Phi(x, y)$, i.e.

$$\Phi(x + t \Delta x, y) < \Phi(x, y) \quad (3.54)$$

for some $t > 0$.

Proof. for the sake of simplicity, let us write the proof dismissing the argument y . We observe that

$$\Phi(x + t \Delta x) - \Phi(x) = \frac{\partial \Phi(x)}{\partial x} \Delta x t + \mathcal{O}(t^2) \quad (3.55)$$

$$= -t \nabla_x \Phi(x, y)^\top M^{-1} \nabla_x \Phi(x, y) + \mathcal{O}(t^2) < 0 \quad (3.56)$$

if $M > 0$ and for $t > 0$ sufficiently small. \square

Hence it is sufficient to choose a positive-definite Hessian approximation to guarantee that the Newton iteration converges. However, the results of Theorem 6 still apply, hence a poor Hessian approximation results in a strong degradation of the convergence rate.

3.5.1 GAUSS-NEWTON HESSIAN APPROXIMATION

Let us consider optimization problems of the form

$$\min_{\mathbf{x}} \frac{1}{2} \|\boldsymbol{\phi}(\mathbf{x}) - \mathbf{y}\|^2 \quad (3.57)$$

for some function $\boldsymbol{\phi} : \mathbb{R}^{n_x} \times \mathbb{R}^{n_\phi}$.

Optimization problems of the form (3.57) are very common in system identification, see Section 4. The gradient and Hessian of the cost function $\Phi(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\boldsymbol{\phi}(\mathbf{x}) - \mathbf{y}\|^2$, then read as:

$$\nabla_{\mathbf{x}} \Phi(\mathbf{x}, \mathbf{y}) = \nabla_{\mathbf{x}} \boldsymbol{\phi}(\mathbf{x}) (\boldsymbol{\phi}(\mathbf{x}) - \mathbf{y}) \quad (3.58)$$

$$\nabla_{\mathbf{x}}^2 \Phi(\mathbf{x}, \mathbf{y}) = \nabla_{\mathbf{x}} \boldsymbol{\phi}(\mathbf{x}) \nabla_{\mathbf{x}} \boldsymbol{\phi}(\mathbf{x})^\top + \left[\nabla_{x_i, x_j} \boldsymbol{\phi}(\mathbf{x}) (\boldsymbol{\phi}(\mathbf{x}) - \mathbf{y}) \right]_{i,j} \quad (3.59)$$

where $[.]_{i,j}$ denotes a matrix where the elements i, j are detailed between the brackets.

For many fitting problems, the evaluation of the second term in the Hessian $\nabla_{\mathbf{x}}^2 \Phi(\mathbf{x}, \mathbf{y})$ is expensive, such that the Gauss-Newton Hessian approximation:

$$\nabla_{\mathbf{x}}^2 \Phi(\mathbf{x}, \mathbf{y}) \approx M_{GN} = \nabla_{\mathbf{x}} \boldsymbol{\phi}(\mathbf{x}) \nabla_{\mathbf{x}} \boldsymbol{\phi}(\mathbf{x})^\top \quad (3.60)$$

is often preferred. The Gauss-Newton Hessian approximation (3.60) is valid if:

- the function $\boldsymbol{\phi}(\mathbf{x})$ is not very nonlinear, such that its second-order derivatives $\nabla_{x_i, x_j} \boldsymbol{\phi}(\mathbf{x})$ are small, or
- the optimal solution to the fitting problem (3.57) yields $\boldsymbol{\phi}(\mathbf{x}^*) \approx \mathbf{y}$, such that $\boldsymbol{\phi}(\mathbf{x}) - \mathbf{y}$ is small.

Both cases justify the dismissal of the second term $\left[\nabla_{x_i, x_j} \boldsymbol{\phi}(\mathbf{x}) (\boldsymbol{\phi}(\mathbf{x}) - \mathbf{y}) \right]_{i,j}$ in $\nabla_{\mathbf{x}}^2 \Phi(\mathbf{x}, \mathbf{y})$. It is useful to observe that the Gauss-Newton Hessian approximation is by construction semi-positive definite, as $M_{GN} \geq 0$. In order to ensure its strict positiveness, it is common to add some *regularization*, i.e. to use a Hessian approximation of the form:

$$M = \nabla_{\mathbf{x}} \boldsymbol{\phi}(\mathbf{x}) \nabla_{\mathbf{x}} \boldsymbol{\phi}(\mathbf{x})^\top + \alpha I \quad (3.61)$$

for some $\alpha > 0$ reasonably small.

3.5.2 CONVEX OPTIMIZATION

Before concluding this section, we ought to briefly talk about convexity in optimization problems. For a (sufficiently) smooth cost function $\Phi(x, y)$, the solution x^* must satisfy the condition

$$\nabla_x \Phi(x^*, y) = 0 \quad (3.62)$$

However, a point satisfying condition (3.62) is not necessarily the solution to the underlying optimization problem (3.48). This problem is illustrated in Fig. 3.6.

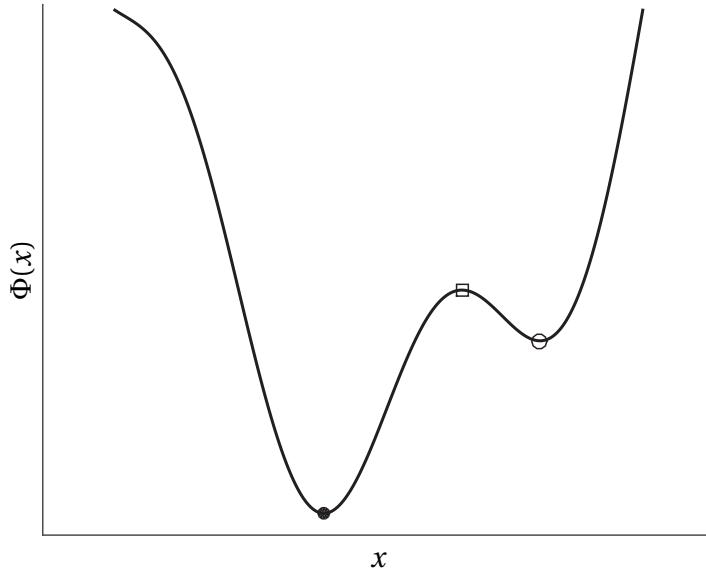


Figure 3.6: Illustration of the difficulty of finding the minimum of a function $\Phi(x)$ via solving (3.62). Several points satisfy (3.62) here. The black dot is the true minimum, but (3.62) also yields a so-called local minimum (circle), which is lower than its neighbours but not lower than the black dot. Condition (3.62) also returns (possibly local) maxima (square).

Mathematically speaking we would say that condition (3.62) is necessary but not sufficient to deliver the solution to (3.48). An exception to this lack of equivalence between condition (3.62) and solving (3.48) is when the cost function $\Phi(x, y)$ is **convex** in x , i.e. if

$$\nabla_x^2 \Phi(x, y) > 0, \quad \forall x \quad (3.63)$$

In this case, the solution to (3.62) is unique and delivers the solution to (3.48). Moreover, the Newton iteration with reduced steps is guaranteed to converge to the solution of (3.48).

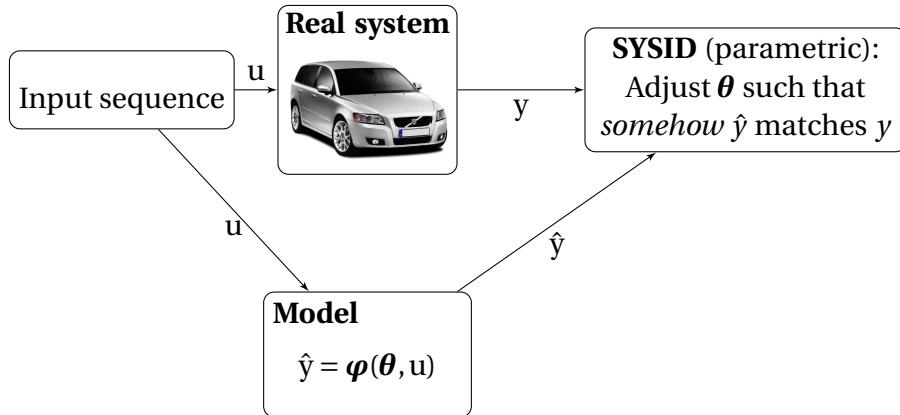
3.6 SUMMARY

The results concerning the Newton method are a bit convoluted. We provide hereafter a brief summary condensing the main points:

- Exact reduced Newton steps Δx improves φ for sufficiently small step sizes $t \in]0, 1]$
- Inexact reduced Newton steps Δx improve φ for a sufficiently small step size $t \in]0, 1]$ if M is sufficiently close to $\frac{\partial \varphi}{\partial x}$. In the context of optimization, $M > 0$ and sufficiently small steps $t \in]0, 1]$ reduce the cost function Φ .
- Exact full ($t = 1$) Newton steps converge quadratically if close enough to the solution
- Inexact full ($t = 1$) Newton steps converge linearly if close enough to the solution and if the Jacobian approximation is "sufficiently good"
- The Newton iteration fails if $\frac{\partial \varphi}{\partial x}$ becomes singular
- Newton methods with reduced steps converge in two phases: damped (slow) phase when reduced steps ($t < 1$) are needed, quadratic/ linear when full steps are possible.

4 SYSTEM IDENTIFICATION (SYSID)

In this Chapter, we will consider the problem of estimating the parameters of a model of a system based on data collected on the system. The principle is illustrated below. An input sequence u is applied to the real system, generating an output sequence y . A model (that can take various forms) holds a set of parameters θ , and generates an output sequence \hat{y} from the input sequence u . The goal of SYSID is to adjust the parameters θ such that the model output \hat{y} matches in some sense the output of the real system y .



Note that we will work in discrete time for SYSID. Indeed, since the measurements are always collected in a discrete fashion on the real system, it makes sense to build a theory where the model output is also discrete.

The system identification procedure we will consider here will then hold the following objects:

- A mathematical model of the system, aimed at explaining an “input-output” relationship:

$$\hat{y} = \phi(\theta, u) \quad (4.1)$$

where u is our input and \hat{y} our model output, and θ a set of model parameters.

- A set of data (or measurements), y physically taken from the system itself.
- A noise realization e explaining the mismatch between the model and the measurements, which is often (but not always!) considered as additive, i.e.:

$$y = \hat{y} + e = \phi(\theta, u) + e \quad (4.2)$$

The question that will be at the heart of this chapter is: for a given model ϕ , a given input u and a given data set y , what is the best model parameter vector θ ?

A classic approach to answer that question is to use a least-squares minimization, i.e. to compute the parameters by solving:

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \|y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u)\|^2 \quad (4.3)$$

Note that the least-squares approach typically tries to minimize the noise e explaining the mismatch between the model and the real system. Indeed, one can trivially see that the choice of additive noise (4.2):

$$e = y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u) \quad (4.4)$$

While the least-squares approach is very commonly used in order to solve the system identification problem, it is very important to understand what it means, why and when it works, and what alternatives can be used when it fails.

Note that in the context of dynamic systems, the input-output of the system are typically input and output sequences, i.e.

$$y = \begin{bmatrix} \hat{y}_0 \\ \vdots \\ \hat{y}_{N-1} \end{bmatrix}, \quad u = \begin{bmatrix} u_0 \\ \vdots \\ u_{N-1} \end{bmatrix} \quad (4.5)$$

and a (causal) model is then in the form:

$$\hat{y} = \begin{bmatrix} \hat{y}_0 \\ \vdots \\ \hat{y}_{N-1} \end{bmatrix} = \boldsymbol{\varphi}(\boldsymbol{\theta}, u) = \begin{bmatrix} \varphi_0(\boldsymbol{\theta}, u_0) \\ \varphi_1(\boldsymbol{\theta}, u_0, u_1) \\ \vdots \\ \varphi_{N-1}(\boldsymbol{\theta}, u_0, \dots, u_{N-1}) \end{bmatrix} \quad (4.6)$$

Before building some theory, let us consider a simple example of the least-squares approach applied to a basic "FIR" (Finite-Impulse Response) model.

Examples

Let us consider a first example of parameter estimation for dynamic systems. Consider a model of the form:

$$\hat{y}_k = \sum_{i=0}^n \boldsymbol{\theta}_i u_{k-i} \quad (4.7)$$

having $n + 1$ parameters. We can also rewrite it in the compact form:

$$\hat{y} = \begin{bmatrix} y_0 \\ \vdots \\ y_{N-1} \end{bmatrix} = \underbrace{\begin{bmatrix} u_0 & \dots & u_{-n} \\ u_1 & \dots & u_{1-n} \\ \vdots & & \vdots \\ u_{N-1} & \dots & u_{N-1-n} \end{bmatrix}}_{:=U} \boldsymbol{\theta} = U\boldsymbol{\theta} \quad (4.8)$$

Hence the model is linear in the parameters $\boldsymbol{\theta}$, and reads as:

$$\boldsymbol{\varphi}(\boldsymbol{\theta}, \mathbf{u}) = \mathbf{U}\boldsymbol{\theta} \quad (4.9)$$

The least-square estimation problem then reads as:

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \|\mathbf{y} - \boldsymbol{\varphi}(\boldsymbol{\theta}, \mathbf{u})\|^2 = \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \|\mathbf{y} - \mathbf{U}\boldsymbol{\theta}\|^2 \quad (4.10)$$

We will see in a bit that the solution to (4.10) is:

$$\hat{\boldsymbol{\theta}} = (\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top \mathbf{y} \quad (4.11)$$

We will discuss in detail the meaning of this construction later in this Chapter. Let us illustrate next this example for an FIR model with $n + 1 = 20$. The next Figure shows the outcome $\hat{\boldsymbol{\theta}}$ resulting from (4.11) for different inputs (left-hand graphs), yielding different output of the real system and model (centred graphs), and set of parameters $\hat{\boldsymbol{\theta}}$ (right-hand side graphs)

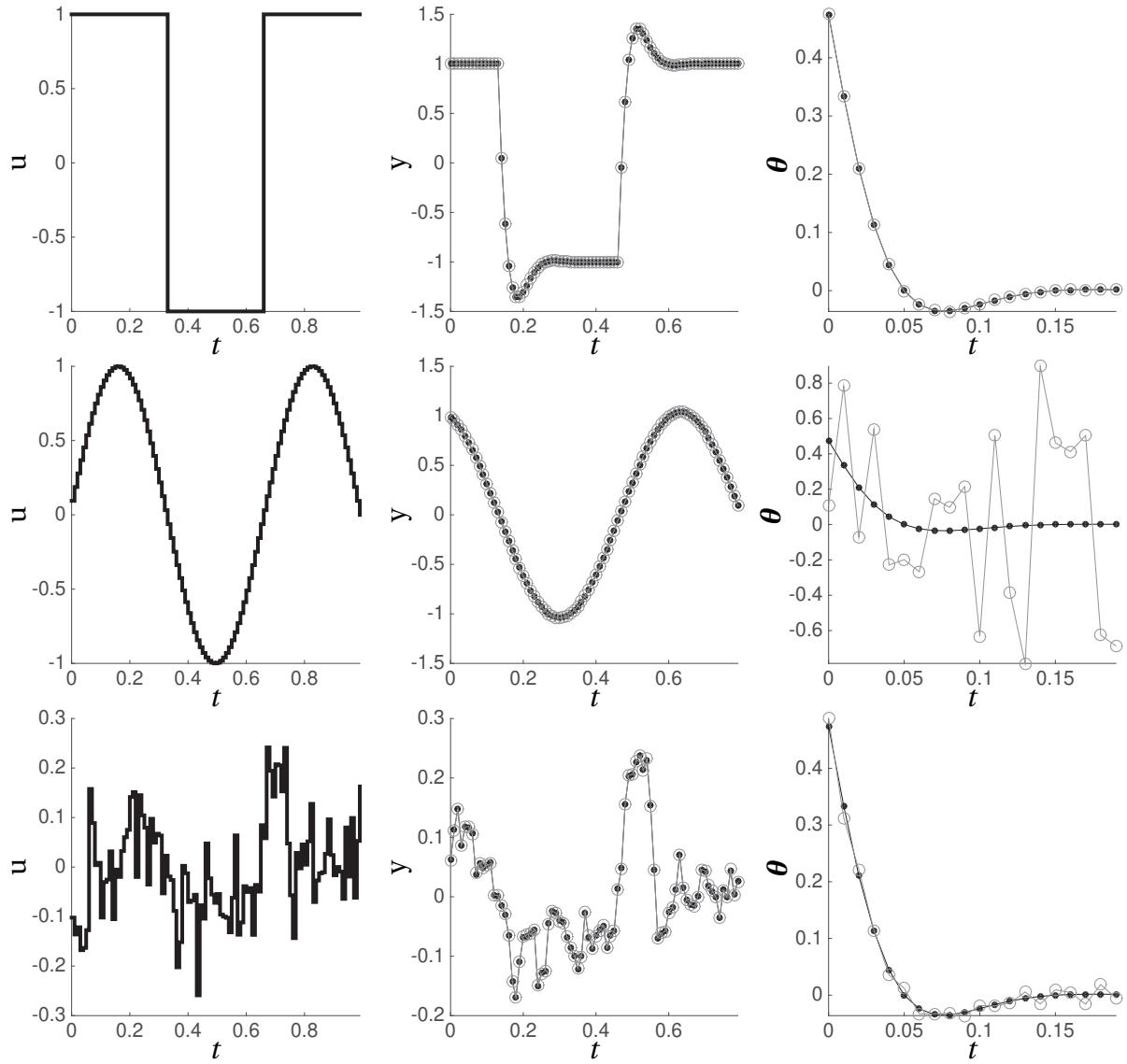


Figure 4.1: Illustration of the least-squares approach for a model of the form (4.7) for $n + 1 = 20$. The real system here is also of the form (4.7) with $n + 1 = 20$, but a slight measurement noise is introduced. The input sequences u are displayed on the left-hand side graphs. The black signals represent the outputs y (centred graphs) and parameters (right-hand side graphs) of the real system, while the grey signals represent the outputs \hat{y} (centred graphs) and estimated parameters (right-hand side graphs) of the model, where the parameters are estimated using (4.11).

One can observe that the choice of input has a strong impact on the quality of the parameter estimation $\hat{\theta}$, and that a poor parameter estimation can be obtained while the model output fits the system output fairly well. We will discuss these observations further in these notes.

□

Next we will build a theory behind the system identification problem. In order to do so, we need to start from the Bayesian principle.

4.1 BAYESIAN REASONING & BAYESIAN ESTIMATION

In order to make sense of the theory coming below, let us consider the following game. Consider we have an opaque bag containing N balls, possibly black (B) and white (W). We draw a ball, note its color, *put it back in the bag*, and repeat. Suppose we made n draws and saw only white balls. What can we say about the hypothesis (H) that “all balls in the bag are white”?

We can start with making a few intuitive observations:

- H (the hypothesis that all balls are white) is true or false, but with limited observations we cannot conclude with certainty that it is true or false
- Regardless of n (even with $n \gg N$), we can still not be certain about H. We may just be “unlucky” and keep drawing only white balls even though there are some black balls in the bag.
- Nonetheless, it should be clear that as n becomes larger and larger (lots of only white balls drawn), we should feel more and more confident that all balls in the bag are white.

How to approach these observations in a formal way, and build a framework where we can discuss the validity of a given hypothesis using limited information? This problem is addressed via Bayesian reasoning. Bayesian reasoning is the mathematical tool that allows medical doctors to conclude that a new drug is effective using trials. It is the tool that allows biologists to assess the state of our environment with limited observations. It is the tool that allows physicists to use experiments to validate theories. It is a tool that allows the detection of e.g. fraudulent trading activities from unusual patterns in stock exchanges. In this section we will have a brief look at the formalism used in Bayesian reasoning.

Bayesian reasoning replaces the true/false, or equivalently “0/1”, status of a hypothesis by a scale that can take all values in $[0, 1]$. One classic interpretation of this scale is to look at it as the *probability that the hypothesis is true*, where 1 means that the hypothesis is certainly true, and 0 means that it is certainly wrong. Note that these probabilities will always be conditional, conditioned on the existing observations that ought to inform us on the validity of the hypothesis. In the mathematical language, we are talking about a *conditional probability* here, which is usually noted as:

$$\mathbb{P} [H | O] \tag{4.12}$$

and means “Probability that the hypothesis H is true provided that we gathered the observation O”. More specifically, in our B/W balls game, we would have:

$$\mathbb{P}[\underbrace{\text{All balls in the bag are white}}_{H} \mid \underbrace{n \text{ white balls drawn}}_O] \quad (4.13)$$

Experimental sciences consider a hypothesis as validated if $\mathbb{P}[H \mid O]$ is sufficiently high (typically at least 0.95, but this value is purely arbitrary!). In the context of estimation, we use Bayesian reasoning in a slightly different way. We will rather consider making a choice between different competing hypothesis (say H_1, \dots, H_N) and make that choice by simply favoring the hypothesis that has the highest probability of being true. In the game described above, the hypothesis H_k would be:

$$H_k : \text{there are } k \text{ white balls in the bag} \quad (4.14)$$

and one would assess the probability of the different H_k to be true provided the observations, and select the one having the highest probability.

Let us now apply Bayesian reasoning in the context of parameter estimation. Consider the model

$$\hat{y} = \boldsymbol{\varphi}(\boldsymbol{\theta}, u) \quad (4.15)$$

mapping the inputs u into the model outputs \hat{y} . We want to find the parameter $\boldsymbol{\theta}$ that best explains the measurements (observations) y collected on the real system.

We can then consider y as the observations, and the different possible values for $\boldsymbol{\theta}$ as competing hypothesis, that we assess via the conditional probability:

$$\mathbb{P}[\boldsymbol{\theta} \mid y] \quad (4.16)$$

meaning “probability that a *given* parameter $\boldsymbol{\theta}$ is true *provided* the measurements y ”. One ought to observe that $\mathbb{P}[\boldsymbol{\theta} \mid y]$ is in fact a function of $\boldsymbol{\theta}$ and y , which returns a scalar value in the interval $[0, 1]$. Selecting the parameters $\boldsymbol{\theta}$ that explain best the observations y then boils down to solving:

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \mathbb{P}[\boldsymbol{\theta} \mid y] \quad (4.17)$$

i.e. we ought to select the parameters $\hat{\boldsymbol{\theta}}$ as the parameters *most likely to be correct (true)*.

Unfortunately, evaluating (4.16) is impossible for most cases of interest. Fortunately, Bayes theorem saves the day. Indeed, as we will soon see, computing $\mathbb{P}[y \mid \boldsymbol{\theta}]$ is in most cases a lot easier than (4.16), and the Bayes theorem allows us to perform this “reversal” of the hypothesis and observations. An application of the Bayes theorem to (4.16) yields:

$$\mathbb{P}[\boldsymbol{\theta} \mid y] = \frac{\mathbb{P}[y \mid \boldsymbol{\theta}] \mathbb{P}[\boldsymbol{\theta}]}{\int \mathbb{P}[y \mid \boldsymbol{\theta}] \mathbb{P}[\boldsymbol{\theta}] d\boldsymbol{\theta}} \quad (4.18)$$

A slight difficulty with (4.18) is to provide the *prior* distribution $\mathbb{P}[\boldsymbol{\theta}]$, which is basically “the probability that a parameter $\boldsymbol{\theta}$ is correct *before any observation has been made*”. This distribution is referred to as the “prior” in the Bayesian context as in “prior knowledge”. In the framework we are developing here, though, we ought to adopt an *indifferent* view of the competing values of the parameter $\boldsymbol{\theta}$, i.e. we ought to consider that since we have no idea what that parameter is a priori (before we have collected observations), all values have equal probability. More formally, we ought to consider that:

$$\mathbb{P}[\boldsymbol{\theta}] = c \quad (4.19)$$

for some constant⁵ c .

It follows that (4.18) simplifies to:

$$\mathbb{P}[\boldsymbol{\theta} | \mathbf{y}] = \frac{\mathbb{P}[\mathbf{y} | \boldsymbol{\theta}]}{\int \mathbb{P}[\mathbf{y} | \boldsymbol{\theta}] d\boldsymbol{\theta}} \propto \mathbb{P}[\mathbf{y} | \boldsymbol{\theta}] \quad (4.20)$$

(here we use \propto for “proportional to”). We then observe that

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \mathbb{P}[\boldsymbol{\theta} | \mathbf{y}] = \arg \max_{\boldsymbol{\theta}} \mathbb{P}[\mathbf{y} | \boldsymbol{\theta}] \quad (4.21)$$

i.e. *selecting the parameters most likely to be true* as done in (4.17) is equivalent to *selecting the parameters most likely to yield the measurements \mathbf{y}* as done in (4.21), right-hand side. The latter is referred to as Maximum-Likelihood Estimation (MLE). We discuss this approach further next, for the specific problem of estimation model parameters.

4.2 MAXIMUM LIKELIHOOD ESTIMATION

Let us start with looking at the likelihood function. In words, the concept described by this function is: what is the probability of observing a set of measurements \mathbf{y} for a given set of parameters $\boldsymbol{\theta}$? As detailed in (4.2), we often assume (but not always) that the mismatch between the measurements \mathbf{y} and the model $\hat{\mathbf{y}} = \boldsymbol{\varphi}(\boldsymbol{\theta}, \mathbf{u})$ can be explained by some additive noise \mathbf{e} . Asking what is the probability of observing the given measurements \mathbf{y} for a given parameter value $\boldsymbol{\theta}$ then boils down to asking what is the probability that the additive noise takes the value:

$$\mathbf{e} = \mathbf{y} - \boldsymbol{\varphi}(\boldsymbol{\theta}, \mathbf{u}) \quad (4.22)$$

for a given $\boldsymbol{\theta}$ i.e. the Likelihood function read as:

$$L(\boldsymbol{\theta}) = \mathbb{P}[\mathbf{y} = \mathbf{e} + \boldsymbol{\varphi}(\boldsymbol{\theta}, \mathbf{u})] \quad (4.23)$$

and can be understood as follows:

⁵an attentive reader may notice here that (4.19) is actually not a distribution. Indeed, the choice (4.19) yields $\int \mathbb{P}[\boldsymbol{\theta}] d\boldsymbol{\theta} \neq 1$. In fact, the integral is unbounded. Fortunately, this is not a problem in the context of the Bayes theorem, and (4.19) is sometimes referred to as a pseudo-distribution.

- The hypothesis used in (4.23) is $y = e + \varphi(\theta, u)$. It describes the mismatch between the measurements y and the model output $\varphi(\theta, u)$ “explained” by the noise e . It is important here to understand that since y, u are known and θ is provided (as an argument in the Likelihood function), the only random variable here is e , hence the notion of probability applies to e .
- We can easily make the latter observation explicit by rewriting (4.23) as:

$$L(\theta) = \mathbb{P} [e = y - \varphi(\theta, u)] \quad (4.24)$$

Here we form the hypothesis that the noise e took a specified value $y - \varphi(\theta, u)$ (specified since y, u are known, and θ is provided).

An important observation ought to be made here. In order to be able to specify the likelihood function (4.24) further, we need to decide on a “noise model”, i.e. we need to specify what is the probability distribution of the random variable e , such that we can translate (4.23) into specific expressions.

In order to digest these concepts, let us make some simple examples.

Examples

1. Let us start with a trivial, but hopefully helpful case. Consider a single measurement, $\hat{y} \in \mathbb{R}$ used to estimate a single parameter $\theta \in \mathbb{R}$, based on the underlying model:

$$\hat{y} = \theta \cdot u \quad (4.25)$$

where $u \in \mathbb{R}$ is a given (single) input. Estimating the parameter θ in this case is fairly obvious, as one ought to compute $\hat{\theta} = \frac{\hat{y}}{u}$. Let us nevertheless use this example in order to unpack the notion of likelihood function (and the concepts that we will introduce later on). Since we expect the measurement y to be corrupted by measurement error, we ought to view our estimation problem as:

$$y = \theta \cdot u + e \quad (4.26)$$

where e explains the difference between the measurement and the model. The likelihood function then reads as:

$$L(\theta) = \mathbb{P} [e = y - \theta \cdot u] \quad (4.27)$$

In order to further specify the likelihood function, we need to pin a distribution to the random variable e . If we e.g. use a normal centred distribution, i.e. $e \sim \mathcal{N}(0, \sigma)$, then we can write:

$$\mathbb{P} [e = x] = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}x^2} \quad (4.28)$$

and our likelihood function reads as:

$$L(\theta) = \mathbb{P} [e = y - \theta \cdot u] = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(y-\theta \cdot u)^2} \quad (4.29)$$

2. Consider the basic problem of taking N measurements of a scalar quantity to estimate (e.g. the temperature in a room) in order to reduce the risk of having an erroneous estimation. We label θ the quantity to estimate via direct measurements, and we write the trivial model

$$\hat{y} = \underbrace{\begin{bmatrix} \theta \\ \vdots \\ \theta \end{bmatrix}}_{:=\boldsymbol{\theta}} \quad (4.30)$$

representing the fact that we take direct measurements of that quantity. Note that we do not have "inputs" in this example. The error vector e is then given by:

$$y = \hat{y} + e \quad (4.31)$$

The likelihood function then reads as:

$$L(\boldsymbol{\theta}) = \mathbb{P}[e = y - \boldsymbol{\theta}] \quad (4.32)$$

In order to further specify the likelihood function, we need a probability distribution for the random variable e . If we select a noise model (i.e. a probability distribution for the random variable e) using e.g. the very common centred *multi-variate* normal distribution:

$$e \sim \mathcal{N}(0, \Sigma) \quad (4.33)$$

for which the probability distribution reads as:

$$\mathbb{P}[e = x] = \frac{1}{\det(2\pi\Sigma)^{\frac{1}{2}}} e^{-\frac{1}{2}x^\top \Sigma^{-1} x} \quad (4.34)$$

Note that if the random variable e is uncorrelated, i.e. e_i is independent of e_j for $i \neq j$, then matrix Σ is diagonal, and (4.34) becomes a product of scalar normal distributions. It then reads as:

$$\mathbb{P}[e = x] = \prod_{k=0}^{N-1} \frac{1}{\sqrt{2\pi\Sigma_k}} e^{-\frac{1}{2\Sigma_k}x_k^2} \quad (4.35)$$

where Σ_k are the elements in the diagonal of matrix Σ , and correspond to the variance of the individual elements e_k , i.e. $\Sigma_k = \sigma_k^2$.

Using these information, we can write the likelihood function (4.32) more explicitly, i.e.

$$L(\boldsymbol{\theta}) = \mathbb{P}[e = y - \boldsymbol{\theta}] = \frac{1}{\det(2\pi\Sigma)^{\frac{1}{2}}} e^{-\frac{1}{2}(y - \boldsymbol{\theta})^\top \Sigma^{-1} (y - \boldsymbol{\theta})} \quad (4.36)$$

and if we adopt an uncorrelated noise model of variance σ_k^2 , (4.36) translates to

$$L(\boldsymbol{\theta}) = \mathbb{P}[e = y - \boldsymbol{\theta}] = \prod_{k=0}^{N-1} \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2\sigma_k^2}(y_k - \theta)^2} \quad (4.37)$$

3. Consider the model:

$$\boldsymbol{\varphi}(\boldsymbol{\theta}, \mathbf{u}) = \boldsymbol{\theta}_1 \mathbf{u} + \boldsymbol{\theta}_0 \quad (4.38)$$

which is *affine in input*, and suppose we obtain the data y . The likelihood function then reads as:

$$L(\boldsymbol{\theta}) = \mathbb{P}[e = y - \boldsymbol{\theta}_1 \mathbf{u} - \boldsymbol{\theta}_0] \quad (4.39)$$

If we select a noise model using e.g. the centred normal distribution again:

$$\mathbf{e} \sim \mathcal{N}(0, \Sigma) \quad (4.40)$$

we can write the likelihood function (4.39) as:

$$L(\boldsymbol{\theta}) = \mathbb{P}[e = y - \boldsymbol{\theta}_1 \mathbf{u} - \boldsymbol{\theta}_0] = \frac{1}{\det(2\pi\Sigma)^{\frac{1}{2}}} e^{-\frac{1}{2}(y - \boldsymbol{\theta}_1 \mathbf{u} - \boldsymbol{\theta}_0)^\top \Sigma^{-1} (y - \boldsymbol{\theta}_1 \mathbf{u} - \boldsymbol{\theta}_0)} \quad (4.41)$$

and if we adopt a normally distributed uncorrelated noise model of variance σ_k^2 , (4.41) translates to:

$$\mathbb{P}[e = y - \boldsymbol{\theta}_1 \mathbf{u} - \boldsymbol{\theta}_0] = \prod_{k=0}^{N-1} \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2\sigma_k^2}(y_k - \boldsymbol{\theta}_1 u_k - \boldsymbol{\theta}_0)^2} \quad (4.42)$$

4. In order to conclude this list of simple examples, let us revisit the first example in this list, but assuming now that the random variable e does not follow a normal centred distribution, but e.g. a uniform distribution in the interval $[-1, 1]$. The distribution of e then reads as:

$$\mathbb{P}[e = x] = \begin{cases} \frac{1}{2} & \text{if } x \in [-1, 1] \\ 0 & \text{otherwise} \end{cases} \quad (4.43)$$

Our likelihood function then reads as:

$$L(\boldsymbol{\theta}) = \mathbb{P}[e = y - \boldsymbol{\theta} \cdot \mathbf{u}] = \begin{cases} \frac{1}{2} & \text{if } y - \boldsymbol{\theta} \cdot \mathbf{u} \in [-1, 1] \\ 0 & \text{otherwise} \end{cases} \quad (4.44)$$

which we could rewrite as:

$$\mathbb{P}[e = y - \boldsymbol{\theta} \cdot \mathbf{u}] = \begin{cases} \frac{1}{2} & \text{if } -1 \leq y - \boldsymbol{\theta} \cdot \mathbf{u} \leq 1 \\ 0 & \text{otherwise} \end{cases} \quad (4.45)$$

□

Equipped with the concept of likelihood function, one can provide a very meaningful solution to the problem of estimating model parameters from data: the most sensible parameter to select is the one that maximizes the likelihood function. In other words, the maximum likelihood estimation approach selects the parameter as the one for which the model

has the highest probability of delivering the observed data. Mathematically, we express the maximum likelihood estimator (MLE) as:

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \mathbb{P} [e = y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u)] \quad (4.46)$$

For a given model $\boldsymbol{\varphi}$, measurements y , and input u , the MLE optimization problem (4.46) returns as estimation the value of $\boldsymbol{\theta}$ that achieves the maximum of the likelihood function. We can now revisit our examples above, applying the MLE approach.

Examples

1. Our first problem had the likelihood function:

$$L(\theta) = \mathbb{P} [e = y - \theta \cdot u] = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(y-\theta\cdot u)^2} \quad (4.47)$$

which is essentially a Gaussian function. It takes its maximum when the argument of the exponential is zero, i.e. when $y - \theta \cdot u = 0$. It follows that, unsurprisingly enough, the MLE of θ in this case is simply $\hat{\theta} = \frac{y}{u}$.

2. Our second example had the likelihood function:

$$L(\boldsymbol{\theta}) = \mathbb{P} [e = y - \boldsymbol{\theta}] = \prod_{k=0}^{N-1} \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2\sigma_k^2}(y_k-\theta)^2} \quad (4.48)$$

In order to apply the MLE, we need to maximize (4.48) over $\boldsymbol{\theta}$. This can be achieved with a bit of calculus. Indeed, one can verify that:

$$\nabla_{\boldsymbol{\theta}} \mathbb{P} [e = y - \boldsymbol{\theta}] = \sum_{i=0}^{N-1} \prod_{k=0}^{N-1} \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2\sigma_k^2}(y_k-\theta)^2} \frac{1}{\sigma_i^2} (y_i - \theta) \quad (4.49)$$

$$= \mathbb{P} [e = y - \boldsymbol{\theta}] \cdot \sum_{i=0}^{N-1} \frac{1}{\sigma_i^2} (y_i - \theta) \quad (4.50)$$

and $\nabla_{\boldsymbol{\theta}} \mathbb{P} [e = y - \boldsymbol{\theta}] = 0$ is achieved by:

$$\sum_{i=0}^{N-1} \frac{1}{\sigma_i^2} (y_i - \theta) = 0 \quad (4.51)$$

and yields:

$$\hat{\theta} = \frac{1}{\sum_{i=0}^{N-1} \frac{1}{\sigma_i^2}} \sum_{i=0}^{N-1} \frac{1}{\sigma_i^2} y_i \quad (4.52)$$

One can observe that (4.52) is a weighted average of the measurements \hat{y}_i , where the measurements having a low variance σ_i^2 receive a proportionally higher weight. If all measurements have the same variance, then (4.52) boils down to a simple average:

$$\hat{\theta} = \frac{1}{N} \sum_{i=0}^{N-1} y_i \quad (4.53)$$

3. In this example with an affine model, we obtained the likelihood function:

$$L(\boldsymbol{\theta}) = \mathbb{P} [e = y - \boldsymbol{\theta}_1 u - \boldsymbol{\theta}_0] = \prod_{k=0}^{N-1} \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2\sigma_k^2}(y_k - \boldsymbol{\theta}_1 u_k - \boldsymbol{\theta}_0)^2} \quad (4.54)$$

We can apply the same principles as in the previous example, and get:

$$\nabla_{\boldsymbol{\theta}} \mathbb{P} [e = y - \boldsymbol{\theta}] = \mathbb{P} [e = y - \boldsymbol{\theta}] \cdot \sum_{i=0}^{N-1} \frac{1}{\sigma_i^2} (y_i - \boldsymbol{\theta}_1 u_i - \boldsymbol{\theta}_0) \begin{bmatrix} 1 \\ u_i \end{bmatrix} \quad (4.55)$$

The MLE $\boldsymbol{\theta}$ is then given by:

$$\sum_{i=0}^{N-1} \frac{1}{\sigma_i^2} (y_i - \boldsymbol{\theta}_1 u_i - \boldsymbol{\theta}_0) \begin{bmatrix} 1 \\ u_i \end{bmatrix} = 0 \quad (4.56)$$

While this equation is linear in $\boldsymbol{\theta}$ and can therefore be solved fairly easily, we will leave the final solution aside, and make more sense of what we are doing here in the next section.

4. Our last example had the likelihood function:

$$L(\theta) = \mathbb{P} [e = y - \theta \cdot u] = \begin{cases} \frac{1}{2} & \text{if } -1 \leq y - \theta \cdot u \leq 1 \\ 0 & \text{otherwise} \end{cases} \quad (4.57)$$

Unfortunately, this likelihood function is discontinuous, such that trying to find its maximum by solving $\nabla_{\theta} \mathbb{P} [e = y - \theta \cdot u] = 0$ is of no help. However, function (4.57) takes only two possible values, the largest of which is $\frac{1}{2}$. Hence, any choice of parameter θ that achieves the value $\frac{1}{2}$ is a valid MLE. More specifically, we can describe the MLE of θ as:

$$\hat{\theta} = \{ \theta \mid -1 \leq y - \theta \cdot u \leq 1 \} \quad (4.58)$$

if the set is not empty. If the set is empty, then (4.57) takes only 0 as a value, and any θ is a valid MLE (though not a very meaningful one).

□

4.3 LEAST-SQUARES ESTIMATION

We will now develop the least-squares estimation (and its possible variants) from the MLE principle. The big picture is that the MLE can be converted into a *fitting problem* for many types of distributions of the error e . In particular, for an additive noise having a normal centred distribution, the MLE is equivalent to a least-square estimation (LSE). Let us start with this specific transformation.

Consider the MLE (4.46) for a given model $\boldsymbol{\varphi}$, a given input u and measurements y :

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \mathbb{P} [e = y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u)] \quad (4.59)$$

and for a normal centred distribution of e , i.e.

$$\mathbb{P} [e = x] = \frac{1}{\det(2\pi\Sigma)^{\frac{1}{2}}} e^{-\frac{1}{2}x^\top \Sigma^{-1} x} \quad (4.60)$$

Then the MLE reads as⁶:

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} e^{-\frac{1}{2}(y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u))^\top \Sigma^{-1} (y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u))} \quad (4.61)$$

We will now proceed with a sequence of transformation on (4.61). First consider taking the log on the likelihood function in (4.61). Because the log function is monotonous, a function and its log achieve their maximum at the same point, hence:

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \log \left[e^{-\frac{1}{2}(y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u))^\top \Sigma^{-1} (y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u))} \right] \quad (4.62)$$

holds, or equivalently:

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} -\frac{1}{2} (y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u))^\top \Sigma^{-1} (y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u)) \quad (4.63)$$

Then one can observe that the maximum of a function f is achieved at the same point as the minimum of $-f$, i.e.

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{2} (y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u))^\top \Sigma^{-1} (y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u)) \quad (4.64)$$

holds. We then observe that (4.64) is a least-squares problem that is often written as:

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \|y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u)\|_{\Sigma^{-1}}^2 \quad (4.65)$$

i.e. the MLE (4.61) delivers the exact same parameter estimation $\hat{\boldsymbol{\theta}}$ as the LSE (4.65). Let us state this observation as clearly as possible hereafter.

The MLE in the presence of an additive error term e that is centred normal distributed is equivalent to a Least-Squares Estimation.

⁶Note that we can dismiss the constant $\frac{1}{\det(2\pi\Sigma)^{\frac{1}{2}}}$ here, as it does not change the value of $\boldsymbol{\theta}$ that achieves the maximum of the likelihood function

Example

- Let us consider a first example of parameter estimation for dynamic systems. Consider a model of the form:

$$\hat{y}_k = \sum_{i=0}^n \boldsymbol{\theta}_i u_{k-i} \quad (4.66)$$

having $n + 1$ parameters. We can also rewrite it in the compact form:

$$\hat{\mathbf{y}} = \begin{bmatrix} y_0 \\ \vdots \\ y_{N-1} \end{bmatrix} = \underbrace{\begin{bmatrix} u_0 & \dots & u_{-n} \\ u_1 & \dots & u_{1-n} \\ \vdots & & \vdots \\ u_{N-1} & \dots & u_{N-1-n} \end{bmatrix}}_{:=U} \boldsymbol{\theta} = U\boldsymbol{\theta} \quad (4.67)$$

The measurements \mathbf{y} are then given by:

$$\mathbf{y} = U\boldsymbol{\theta} + \mathbf{e} \quad (4.68)$$

Consider that the error term \mathbf{e} is normally distributed and centred, such that the MLE associated to estimating the model parameters $\boldsymbol{\theta}$ reads as:

$$\mathbb{P} [\mathbf{e} = \mathbf{y} - U\boldsymbol{\theta}] = \frac{1}{\det(2\pi\Sigma)^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{y}-U\boldsymbol{\theta})^\top \Sigma^{-1} (\mathbf{y}-U\boldsymbol{\theta})} \quad (4.69)$$

and the equivalent LSE problem reads as:

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \|\mathbf{y} - U\boldsymbol{\theta}\|_{\Sigma^{-1}}^2 \quad (4.70)$$

- The transformation (??) to (4.65) can be generalized to many distributions. Let us illustrate this principle on e.g. an uncorrelated "fair" normal centered distribution, i.e.⁷

$$\mathbb{P} [\mathbf{e} = \mathbf{x}] \sim \prod_{k=0}^{N-1} e^{-\frac{x_k^2}{2\sigma^2(1+\alpha|x_k|)}} \quad (4.71)$$

⁷Note that σ^2 is not the variance anymore. This distribution describes a random variable that "looks like" it is normal for e small, but that makes the occurrence of large e more likely than under a normal distribution. This type of distribution is often referred to as a "fat-tail" distribution, because it looks like a Gaussian distribution but with "fatter" tails.

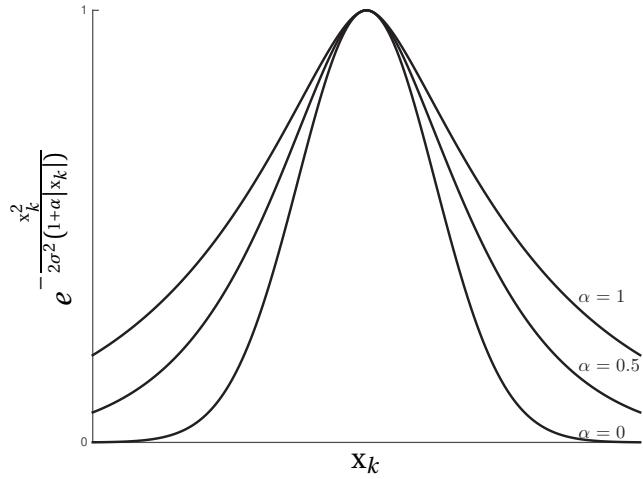


Figure 4.2: Illustration of the fair distribution for a few values of the constant α .

The MLE problem for model (4.68) then reads as:

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \prod_{k=0}^{N-1} e^{-\frac{(y_k - U_k \boldsymbol{\theta})^2}{2\sigma^2(1+\alpha|y_k - U_k \boldsymbol{\theta}|)}} \quad (4.72)$$

and its transformation via the log function reads as:

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma^2} \sum_{k=0}^{N-1} \frac{(y_k - U_k \boldsymbol{\theta})^2}{1 + \alpha |y_k - U_k \boldsymbol{\theta}|} \quad (4.73)$$

We can further rewrite it as:

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{2\sigma^2} \| \mathbf{y} - \mathbf{U} \boldsymbol{\theta} \|_{\text{Fair}} \quad (4.74)$$

where we use the notation

$$\| \mathbf{x} \|_{\text{Fair}} = \sqrt{\sum_{k=0}^{N-1} \frac{x_k^2}{1 + \alpha |x_k|}} \quad (4.75)$$

Note that we are somewhat abusing the standard notation $\| \cdot \|$ here, as (4.75) does not define a norm (see definition of a mathematical norm in Sec. 1.3).

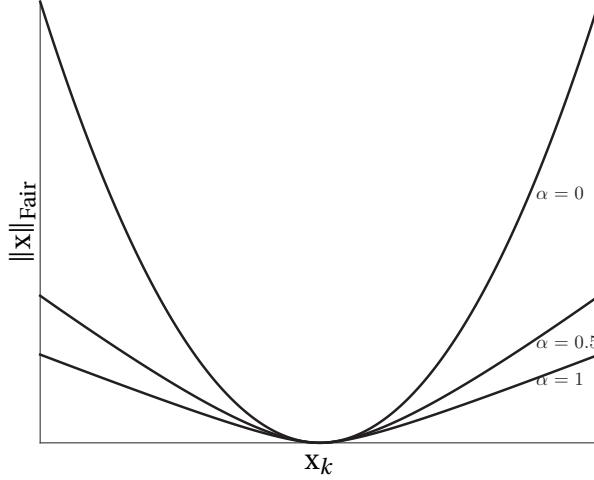


Figure 4.3: Illustration of the cost function (4.75) associated to the fair distribution for a few values of the constant α .

4.4 SOLUTION TO THE LEAST-SQUARES PROBLEM

In this section, we will review the solution of least-squares problem (4.65), i.e. problems of the general form:

$$\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{2} \|y - \boldsymbol{\varphi}(\boldsymbol{\theta}, u)\|_W^2 \quad (4.76)$$

We need to distinguish here two cases: linear least-squares problems (LLS) where the model $\boldsymbol{\varphi}(\boldsymbol{\theta}, u)$ is linear in $\boldsymbol{\theta}$ and nonlinear least-squares problems (NLS) where the model $\boldsymbol{\varphi}(\boldsymbol{\theta}, u)$ is nonlinear in $\boldsymbol{\theta}$. Note that the model function can be nonlinear in u and be linear in $\boldsymbol{\theta}$ nonetheless. While both LLS and NLS use the same mathematical principles, the solution approaches to obtain the solution are radically different. Indeed, LLS problems have an explicit solution, while NLS problems usually require the deployment of a Newton-type method (see Sec. 3).

Linear least-squares matching e.g. the form obtained in (4.70) can be generally written as:

$$\hat{\boldsymbol{\theta}} = \min_{\boldsymbol{\theta}} \frac{1}{2} \|y - A\boldsymbol{\theta}\|_W^2 \quad (4.77)$$

which equivalently read as:

$$\hat{\boldsymbol{\theta}} = \min_{\boldsymbol{\theta}} \frac{1}{2} (y - A\boldsymbol{\theta})^\top W (y - A\boldsymbol{\theta}) \quad (4.78)$$

The solution approach to solving (4.77) relies on finding the roots of the gradient of the cost function underlying (4.77), i.e. labelling $J(\boldsymbol{\theta}, y) = \frac{1}{2} \|y - A\boldsymbol{\theta}\|_W^2$, the parameter estimation $\hat{\boldsymbol{\theta}}$ is given by solving:

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}, y) = \nabla_{\boldsymbol{\theta}} \frac{1}{2} \|y - A\boldsymbol{\theta}\|_W^2 = 0 \quad (4.79)$$

The equation stemming from (4.79) then read as:

$$A^\top W (A\theta - y) = 0 \quad (4.80)$$

and delivers the standard solution to LLS problems, i.e.:

$$\hat{\theta} = (A^\top W A)^{-1} A^\top W y \quad (4.81)$$

There is a handful of crucial observations to be made here:

- The parameter estimation $\hat{\theta}$, if it exists, is a linear function of the data/measurements y .
- The parameter estimation $\hat{\theta}$, if it exists, is unique (this is to be contrasted with NLS problems)
- The parameter estimation $\hat{\theta}$ exists if and only if the square matrix $A^\top W A$ is full rank, i.e. invertible.

When the model $\varphi(\theta, u)$ is nonlinear in θ , then the parameter estimation is still obtained by finding the root of the gradient of the underlying cost function, i.e.

$$\nabla_{\theta} J(\theta, y) = \nabla_{\theta} \frac{1}{2} \|y - \varphi(\theta, u)\|_W^2 = 0 \quad (4.82)$$

but the developments above do not apply further. Equation (4.82) also reads as:

$$\nabla_{\theta} \varphi(\theta, u) W (y - \varphi(\theta, u)) = 0 \quad (4.83)$$

and is nonlinear in θ . It is therefore generally impossible to solve explicitly, and one typically deploys a numerical method in order to find θ . The methods detailed on Section 3.5 can be used.

4.5 COVARIANCE, BIAS & CONSISTENCY

There are some crucial properties underlying parameter estimation that we ought to discuss here. These properties essentially discuss the extent to which a chosen estimation technique delivers a meaningful estimation $\hat{\theta}$. Discussing these properties is the object of this section. First, let us consider a generic parameter estimation technique:

$$\hat{\theta} = \Psi(u, y, N) \quad (4.84)$$

taking the inputs u and the measurements y , and turning them into a parameter estimation. The integer $N \in \mathbb{N}$ represents the number of data point used in building $\hat{\theta}$ (e.g. N in (4.67)). We observe that e.g. a Least-Squares estimation approach such as (4.65) can be written in the form (4.84), but also all the examples presented above.

It is important to observe here that the data y ought to be considered as the realization of a random variable, e.g. of the noise e if we consider that $y = \varphi(\boldsymbol{\theta}, u) + e$. It follows that, $\hat{\boldsymbol{\theta}}$ becomes also a realization of a random variable, as it is a function of y through the estimator Ψ in (4.84). These observations, in practice, imply that if one runs the estimation procedure a number of times, each time using a set of newly collected data y , the outcome of that process is that $\hat{\boldsymbol{\theta}}$ will be different each time, and therefore build over many experiments a distribution itself. It is then useful to discuss what the distribution of $\hat{\boldsymbol{\theta}}$ looks like. Is it wide or narrow? Does it average to the true parameter $\boldsymbol{\theta}_*$? How does the number of data point N impact the distribution of $\hat{\boldsymbol{\theta}}$? These are the questions we discuss next. We need to define first a basic terminology on estimators. Let us label $\boldsymbol{\theta}_*$ the “true” model parameter⁸ that we are trying to estimate through $\hat{\boldsymbol{\theta}}$.

1. An estimator Ψ is labelled **unbiased** if

$$\mathbb{E}[\hat{\boldsymbol{\theta}}] = \boldsymbol{\theta}_* \quad (4.85)$$

where $\mathbb{E}[\cdot]$ labels the expected value operator (i.e. roughly speaking the average over a very large number of experiments giving different $\hat{\boldsymbol{\theta}}$). In words, if one uses an unbiased estimator, running the estimation process a large number of times to collect many $\hat{\boldsymbol{\theta}}$ and averaging them yields the correct parameter $\boldsymbol{\theta}_*$.

2. An estimator is **consistent** if⁹

$$\lim_{N \rightarrow \infty} \hat{\boldsymbol{\theta}} = \boldsymbol{\theta}_* \quad (4.86)$$

In words, if one collects a very large number of data points y (i.e. $N \rightarrow \infty$), the parameter estimation $\hat{\boldsymbol{\theta}}$ will match the correct value.

3. For an **efficient** estimator, the estimates $\hat{\boldsymbol{\theta}}$ become normally distributed for N large enough, and the variance meets the *Cramer-Rao lower bound*. The latter bound comes from information theory and is essentially “the best one can do for a given set of data”. An estimator is therefore efficient if it makes the “best possible use” of the data. It also implies that for $N \rightarrow \infty$ (and some conditions on the problem and data), the normal distribution of the estimates becomes “arbitrarily narrow” (i.e. it tends to a Dirac distribution), such that $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}_*$ almost surely.

Two remarks ought to be useful here.

⁸it is debatable that such a true parameter actually exists. Most models are indeed not capturing the reality perfectly such that there is no such thing as a “true” parameter. It is, however, very useful to reason based on the assumption that there is a true parameter when discussing estimation techniques. Indeed, if such a true parameter could be found, it would arguably be desirable from an estimation technique that it does find this parameter in one way or another.

⁹to be more formal, we ought to say that $\hat{\boldsymbol{\theta}} \rightarrow \boldsymbol{\theta}_*$ *almost surely* as $N \rightarrow \infty$, where “almost surely” means “with probability 1”. The distinction between this statement and the limit (4.86) is subtle and of secondary importance here.

- Efficiency implies consistency. This observation is explained at the end of point 3. Indeed if the estimates distribution tends to a Dirac function, then $\hat{\boldsymbol{\theta}}$ will be arbitrarily close to $\boldsymbol{\theta}_*$ with probability 1.
- Efficiency or consistency do not necessarily imply unbiasedness, i.e. an estimator can be efficient (and) or consistent and still biased. Biasedness has to do with the behavior of estimators for limited (possibly very few) data point, and whether the estimates average to the correct parameter value for *a large amount of experiments*, as opposed to a large amount of data.

We are now ready to discuss the properties of Maximum-Likelihood Estimation (MLE).

- An MLE estimator is **in general biased**. It is unbiased in some special cases, though
 - An MLE estimator is always **efficient** and therefore **consistent**

Let us now turn to more practical considerations. Consider the solution

$$\hat{\boldsymbol{\theta}} = (A^\top W A)^{-1} A^\top W \mathbf{y} \quad (4.87)$$

given in (4.81) to the least-squares problem (4.77), i.e.

$$\hat{\boldsymbol{\theta}} = \min_{\boldsymbol{\theta}} \frac{1}{2} \|A\boldsymbol{\theta} - \mathbf{y}\|_W^2 \quad (4.88)$$

Here we can readily observe how $\hat{\boldsymbol{\theta}}$ is a (linear) function of the data \mathbf{y} . Let us assume that there is a true parameter $\boldsymbol{\theta}_*$ and that the measurements are given by the model disturbed by normally distributed, centered, additive noise, i.e.

$$\mathbf{y} = A\boldsymbol{\theta}_* + \mathbf{e} \quad (4.89)$$

It follows that

$$\hat{\boldsymbol{\theta}} = (A^\top W A)^{-1} A^\top W (A\boldsymbol{\theta}_* + \mathbf{e}) = \boldsymbol{\theta}_* + (A^\top W A)^{-1} A^\top W \mathbf{e} \quad (4.90)$$

We can observe that if the noise \mathbf{e} is centred (i.e. $\mathbb{E}[\mathbf{e}] = 0$), then **regardless of the distribution**:

$$\mathbb{E}[\hat{\boldsymbol{\theta}}] = \boldsymbol{\theta}_* \quad (4.91)$$

holds, i.e. the linear least-squares estimation is unbiased. Hence if the noise is normal distributed, then LSE is a special case of MLE that achieves unbiasedness. We also observe that since $\hat{\boldsymbol{\theta}}$ is an affine transformation of \mathbf{e} , then if \mathbf{e} is normally distributed then so is $\hat{\boldsymbol{\theta}}$. We can assess the covariance of the estimates, i.e.

$$\mathbb{E}[(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_*)(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_*)^\top] = (A^\top W A)^{-1} A^\top W \Sigma_e W A (A^\top W A)^{-1} \quad (4.92)$$

where we use $\Sigma_e = \mathbb{E}[\mathbf{e}\mathbf{e}^\top]$. Finally, in order for our least squares problem to be an MLE (see (4.70)), we need e to be normal distributed, and we need to choose $W = \Sigma_e^{-1}$ (see (4.70)). We then get:

$$\Sigma_{\hat{\theta}} = \mathbb{E}[(\hat{\theta} - \theta_*)(\hat{\theta} - \theta_*)^\top] = (A^\top \Sigma_e^{-1} A)^{-1} \quad (4.93)$$

We therefore observe that for Σ_e large, $\Sigma_{\hat{\theta}}$ is also large, i.e. in words, a lot of measurement noise yield a large co-variance of the parameters. Moreover, we observe that if matrix $A^\top \Sigma_e^{-1} A$ is close to being rank deficient (i.e. ill-conditioned or “close to being not invertible”), then $\Sigma_{\hat{\theta}}$ becomes large (at least in some directions). This explains the observations made in Figure 4.1. Indeed, in the Least-Squares problem (4.11), the matrix $U = A$ and $W = I$. Matrix $A^\top \Sigma_e^{-1} A = U^\top U$ is then built based on the input sequence u , such that the co-variance of the parameter estimation depends highly on the choice of input sequence. The sinusoidal input (middle row) in Figure 4.1 yields a matrix $U^\top U$ that is close to singular and yields a very poor parameter estimation. The other signals yield better conditioned matrices and yield better estimations.

Let us wrap up our observations.

- A Linear Least Square (LLS) problem is equivalent to a MLE if the noise is normally distributed, centred, and for the adequate choice of weights (i.e. $W = \Sigma_e^{-1}$).
- A LLS is unbiased if the noise is additive and centered (regardless of its distribution)
- The covariance of the estimates $\hat{\theta}$, i.e. $\Sigma_{\hat{\theta}}$, depends on matrix A and the noise covariance Σ_e (regardless of its distribution)

4.6 ESTIMATION FOR LINEAR DYNAMIC SYSTEMS

Let us now apply what we have established so far in the context of estimation for dynamic systems. We will work here with discrete-time linear dynamic systems. System identification for continuous systems is a bit more intricate and we will leave it out of these course notes. In general, linear dynamic systems can be represented via the equation:

$$\hat{y}_k = \sum_{i=0}^{N_b} b_i u_{k-i} + \sum_{i=1}^{N_a} a_i \hat{y}_{k-i} \quad (4.94)$$

describing the model output y at time k as a function of the past (and possibly present) inputs u_k, \dots, u_{k-N_b} and past outputs $\hat{y}_{k-1}, \dots, \hat{y}_{k-N_a}$. The model parameters here are the vectors a and b , i.e. we can e.g. write:

$$\boldsymbol{\theta} = \begin{bmatrix} a \\ b \end{bmatrix} \quad (4.95)$$

We can then investigate how to deploy the estimation techniques investigated above to this specific type of model.

4.6.1 ESTIMATION OF FINITE IMPULSE RESPONSE MODELS

To start simple, let us first consider a simple special case of (4.94), namely Finite Impulse Response models, which have $a = 0$, i.e.

$$\hat{y}_k = \sum_{i=0}^{N_b} b_i u_{k-i} \quad (4.96)$$

We have in fact already considered this case in an example above (see (4.66)) for the case $\hat{y}_k \in \mathbb{R}$. Let us expand on this case. Recall that a sequence of model output can be put in a matrix form:

$$\hat{\mathbf{y}} = \begin{bmatrix} y_0 \\ \vdots \\ y_{N-1} \end{bmatrix} = \underbrace{\begin{bmatrix} u_0 & \dots & u_{-N_b} \\ u_1 & \dots & u_{1-n} \\ \vdots & & \vdots \\ u_{N-1} & \dots & u_{N-1-N_b} \end{bmatrix}}_{:=U} \boldsymbol{\theta} = U\boldsymbol{\theta} \quad (4.97)$$

where we label $b = \boldsymbol{\theta}$. If the mismatch between the measurements \mathbf{y} and the model output $\hat{\mathbf{y}}$ are to be explained by additive, centered, normally distributed noise \mathbf{e} , i.e.

$$\mathbf{y} = \hat{\mathbf{y}} + \mathbf{e} \quad (4.98)$$

then MLE of the model parameters $\hat{\boldsymbol{\theta}}$ is a linear least squares problem (see (4.70)):

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \|\mathbf{y} - U\boldsymbol{\theta}\|_{\Sigma^{-1}}^2 \quad (4.99)$$

The estimate $\hat{\boldsymbol{\theta}}$ is then efficient and consistent (delivers true parameter for N large enough), and is also unbiased because it is given by a linear least squares problem. The covariance of $\hat{\boldsymbol{\theta}}$ is given by the formula (4.93) and reads as

$$\Sigma_{\hat{\boldsymbol{\theta}}} = \mathbb{E} \left[(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_*) (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_*)^\top \right] = (U^\top \Sigma_e^{-1} U)^{-1} \quad (4.100)$$

In the specific case that the noise \mathbf{e} is white (i.e. noise samples e_i and $e_{j \neq i}$ are uncorrelated), and has a uniform covariance (i.e. all noise sample e_i have the same variance σ^2), then the estimate covariance is given by:

$$\Sigma_{\hat{\boldsymbol{\theta}}} = \sigma^2 (U^\top U)^{-1} \quad (4.101)$$

We observe that if we have a "large" amount of data, i.e. if $N \gg N_b$, then matrix $U^\top U$ looks as follows:

$$\boxed{U^\top} \quad \boxed{U} = \boxed{U^\top U}$$

i.e. it is made of the product of a “long” lying matrix by a “tall” matrix, and has dimensions $N_b \times N_b$, which creates some “richness” in matrix $U^\top U$. The opposite, i.e. $N < N_b$ (insufficient data) makes matrix $U^\top U$ rank deficient, and the estimates $\hat{\theta}$ are undefined (infinite covariance!).

4.6.2 AUTO-REGRESSIVE MODELS - SIMULATION ERROR

Let us now turn to more general models in the form (4.94), which we label as auto-regressive as the past model outputs participate in building the next model output. Let us consider a series of model outputs $\hat{y}_{0,\dots,N}$ and investigate how they relate to the model parameters. We have seen that in the FIR case, this relationship can be put in the linear form $\hat{y} = U\theta$. Is it still the case for more generic models of the form (4.94)? It is in fact easy to verify that the model output of generic models are in general not a linear function of the parameters θ . Indeed consider the simple example:

$$\hat{y}_k = b_0 u_k + a_1 \hat{y}_{k-1} \quad (4.102)$$

We can then compute:

$$\hat{y}_1 = b_0 u_1 + a_1 \hat{y}_0 \quad (4.103a)$$

$$\hat{y}_2 = b_0 u_2 + a_1 (b_0 u_1 + a_1 \hat{y}_0) = b_0 u_2 + a_1 b_0 u_1 + a_1^2 \hat{y}_0 \quad (4.103b)$$

$$\vdots \quad (4.103c)$$

We can deduce from these developments that more generally:

$$\hat{y}_k = a_1^k \hat{y}_0 + b_0 \sum_{i=1}^k a_1^{k-i} u_i \quad (4.104)$$

We observe that the model output is a linear function of the model inputs u and initial conditions \hat{y}_0 (this is expected as the model is linear), but it is a nonlinear function of the model parameters b_0, a_1 as they are combined nonlinearly (powers and multiplications). We can conclude that it is impossible to write the model output using a linear expression of the form (4.97).

Nonetheless, we observe that (4.104) can be put in the generic model form (4.15) i.e.:

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_N \end{bmatrix} = \begin{bmatrix} a_1 \hat{y}_0 + b_0 u_1 \\ \vdots \\ a_1^N \hat{y}_0 + b_0 \sum_{i=1}^N a_1^{N-i} u_i \end{bmatrix} = \varphi(\theta, u, \hat{y}_0) \quad (4.105)$$

such that the developments of section 4.3 hold, and in particular the MLE-based parameter estimation for (4.105) is given by (4.65), i.e. it is a **Nonlinear Least-Squares problem** of the form:

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \|y - \varphi(\theta, u, \hat{y}_0)\|_{\Sigma^{-1}}^2 \quad (4.106)$$

The estimates $\hat{\theta}$ must then be obtained (in most cases) numerically via a Newton method solving (4.83). It may be useful here to underline that **even though the model is linear, the resulting least-squares problem is nonlinear**. In this case, the resulting estimates $\hat{\theta}$ hold the properties of MLE, namely they are efficient and consistent, but **they are generally biased**.

4.6.3 AUTO-REGRESSIVE MODELS - PREDICTION ERROR

It is time now to reflect on what we are trying to do when estimating the parameters of the AR model (4.94) using collected data y . The model outputs generated by the AR model (4.94) stem from the input sequence u and the initial conditions \hat{y}_0 , such as is e.g. explicitly shown in example (4.105). The model parameters are then adjusted so that the sequence of model outputs follows as closely as possible the data y . This approach, put in other words, seeks at finding model parameters a, b such that the model provides a good *simulation* of the real system. While this approach makes sense in a number of cases, it is not the only view one can adopt to adjusting the model parameters.

Another view consists in finding the model parameters that deliver a small *prediction error*. In that context, one would want to build a model that - fed with the most recent measurements y - delivers a good estimation of what the real system will do next. More specifically, one then considers the model:

$$\hat{y}_k = \sum_{i=0}^{N_b} b_i u_{k-i} + \sum_{i=1}^{N_a} a_i y_{k-i} \quad (4.107)$$

instead of (4.94). The difference between (4.107) and (4.94) is that (4.107) builds \hat{y}_k based on past *measurements* y_{k-i} as opposed to past *model outputs* \hat{y}_{k-i} as in (4.94). The Prediction Error model (4.107) will then focus on building model outputs \hat{y}_k that are good prediction of the future system output based on the data measured so far y_{k-i} . The mismatch:

$$\hat{y}_k - y_k \quad (4.108)$$

between the model output \hat{y}_k and the real system output or measurement y_k is called a prediction error, and adjusting the model parameters a, b according to

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{2} \|\hat{y} - y\|_{\Sigma^{-1}}^2 \quad (4.109)$$

where

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_N \end{bmatrix} \quad (4.110)$$

is given by (4.107). This approach is called a **Prediction Error Method (PEM)**.

We ought to come back to the question whether a linear or nonlinear least-squares problem results from this alternative approach. To that end, we need to understand if \hat{y} in (4.107) can be written as a linear (or affine) function of the model parameters a, b . Answering that question is fairly straightforward. Indeed, we observe that:

$$\hat{y} = \begin{bmatrix} \hat{y}_0 \\ \vdots \\ \hat{y}_N \end{bmatrix} = \underbrace{\begin{bmatrix} y_{-1} & \dots & y_{-N_a} & u_0 & \dots & u_{-N_b} \\ y_0 & \dots & y_{1-N_a} & u_1 & \dots & u_{1-N_b} \\ \vdots & & & & & \\ y_{N-1} & \dots & y_{N-N_a} & u_N & \dots & u_{N-N_b} \end{bmatrix}}_M \underbrace{\begin{bmatrix} a_1 \\ \vdots \\ a_{N_a} \\ b_0 \\ \vdots \\ b_{N_b} \end{bmatrix}}_{\boldsymbol{\theta}} = M\boldsymbol{\theta} \quad (4.111)$$

hence the model output sequence \hat{y} can indeed be written as a linear function of the model parameters $\boldsymbol{\theta}$. The Least-Squares problem (4.109) underlying the PEM method is therefore linear and reads as:

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \|M\boldsymbol{\theta} - y\|_{\Sigma^{-1}}^2 \quad (4.112)$$

where matrix M is in fact also a function of the measurements y .

4.6.4 PREDICTION ERROR & NOISE MODEL

We now need to turn to a not really trivial question, which is to understand how we assume the noise enters into the problem when using the PEM approach. This question appears certainly unclear at this stage, so let us try to unpack it step by step. Let us get back to the FIR model first, where (see (4.96))

$$\hat{y}_k = \sum_{i=0}^{N_b} b_i u_{k-i} \quad (4.113)$$

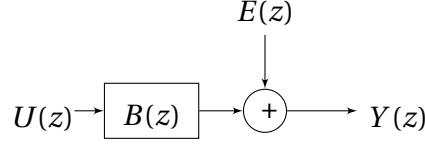
and

$$y_k = e_k + \sum_{i=0}^{N_b} b_i u_{k-i} \quad (4.114)$$

We observe first that in the FIR case the “simulation” approach of section 4.6.2 and the PEM of section 4.6.3 are in fact identical. Indeed, there is not “past data” that can be fed into the model here since the model output depends purely on the input u . If we take the Z -transform of (4.114), we obtain:

$$Y(z) = E(z) + \underbrace{\sum_{i=0}^{N_b} b_i z^{-i}}_{B(z)} U(z) \quad (4.115)$$

such that the FIR setup assumes that the signals are treated via the block-diagram



This picture ought appear fairly unsurprising. The noise e is additive in equations (4.114) and (4.115) and so is it in our block diagram. However, we will apply the same process to the AR model (4.107) used in the PEM now, and a different picture will emerge.

In the PEM context, recall that the model reads as (see (4.107)):

$$\hat{y}_k = \sum_{i=0}^{N_b} b_i u_{k-i} + \sum_{i=1}^{N_a} a_i y_{k-i} \quad (4.116)$$

such that

$$y_k = e_k + \sum_{i=0}^{N_b} b_i u_{k-i} + \sum_{i=1}^{N_a} a_i y_{k-i} \quad (4.117)$$

The Z -transform of (4.117) reads as:

$$Y(z) = E(z) + \underbrace{\sum_{i=0}^{N_b} b_i z^{-i} U(z)}_{=B(z)} + \sum_{i=1}^{N_a} a_i z^{-i} Y(z) \quad (4.118)$$

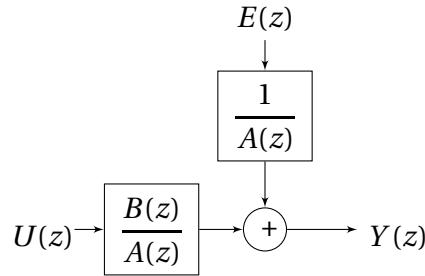
i.e.

$$A(z)Y(z) = E(z) + B(z)U(z) \quad (4.119)$$

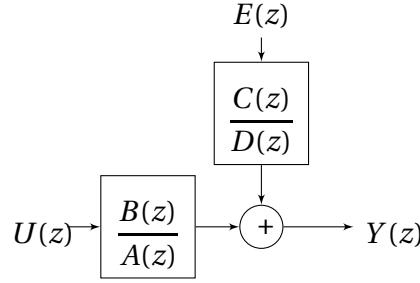
where $A(z) = 1 - \sum_{i=1}^{N_a} a_i z^{-i}$. Hence we finally have that:

$$Y(z) = \frac{B(z)}{A(z)}U(z) + \frac{1}{A(z)}E(z) \quad (4.120)$$

In a block-diagram representation, this relationship would look like:



And we see that the noise is no longer simply additive in the model dynamics, but is in fact “filtered” via the model denominator $A(z)$. An intuitive way of construing this observation is to simply observe that the noise in the PEM dynamics (4.117) does not in fact enter simply added to the model output, but also enters via the past measurements $y_{k-1}, \dots, y_{k-N_a}$ which are subject to that noise, and are filtered by the “A” part of the model. The block-diagram above shows that we are in fact assuming in the PEM context that the noise is filtered by $1/A(z)$. It begs then the question: can we build more generic models of the “filter” affecting the noise in the PEM? More specifically, can we build a model for PEM that underlies the generic block-diagram



i.e. we would like to find a PEM-like formulation that corresponds to the Z -transform:

$$Y(z) = \frac{B(z)}{A(z)}U(z) + \frac{C(z)}{D(z)}E(z) \quad (4.121)$$

Unfortunately, the following developments are a bit tricky (and understanding them is not required for the exam), but we present them next for the interested readers. We first rewrite (4.121) as

$$Y(z) = G(z)U(z) + \underbrace{H(z)E(z)}_{:=V(z)} \quad (4.122)$$

In the time domain, we can write $V(z)$ as:

$$v_k = \sum_{j=0}^{\infty} h_j e_{k-j} = e_k + \sum_{j=1}^{\infty} h_j e_{k-j} \quad (4.123)$$

where h is the impulse response corresponding to $H(z)$ and where we assume $h_0 = 1$ (without loss of generality). At this stage, in order to deliver a prediction \hat{y}_k , we would need a prediction \hat{v}_k of v_k . It makes sense to write \hat{v}_k as:

$$\hat{v}_k = \mathbb{E}[v_k] = \mathbb{E}[e_k] + \sum_{j=1}^{\infty} h_j e_{k-j} = \sum_{j=1}^{\infty} h_j e_{k-j} \quad (4.124)$$

hence

$$\hat{v}_k = \left(\sum_{j=0}^{\infty} h_j e_{k-j} \right) - e_k \quad (4.125)$$

If we convert the latter expression in the Z -domain, we see that

$$\hat{V}(z) = (H(z) - 1) E(z) \quad (4.126)$$

As the noise realisations $e_{k-1}, \dots, e_{-\infty}$ are indirectly known through the past measurements $y_{k-1}, \dots, y_{-\infty}$ using:

$$E(z) = \frac{1}{H(z)} V(z) \quad (4.127)$$

Hence

$$\hat{V}(z) = \frac{H(z) - 1}{H(z)} V(z) = (1 - H(z)^{-1}) V(z) \quad (4.128)$$

holds. From there on, we just have some algebraic manipulations to do in the Z domain. Similarly to (4.122), we observe that

$$\begin{aligned} \hat{Y}(z) &= G(z)U(z) + \hat{V}(z) = G(z)U(z) + (1 - H(z)^{-1}) V(z) \\ &= G(z)U(z) + (1 - H(z)^{-1})(Y(z) - G(z)U(z)) \end{aligned} \quad (4.129)$$

Hence we finally get:

$$\hat{Y}(z) = H(z)^{-1} G(z)U(z) + (1 - H(z)^{-1}) Y(z) \quad (4.130)$$

or equivalently:

$$H(z)\hat{Y}(z) = G(z)U(z) + (H(z) - 1) Y(z) \quad (4.131)$$

Using $H(z) = \frac{C(z)}{D(z)}$ and $G(z) = \frac{B(z)}{A(z)}$, we can also write:

$$\frac{C(z)}{D(z)} \underbrace{\hat{Y}(z)}_{\text{prediction}} = \frac{B(z)}{A(z)} U(z) + \left(\frac{C(z)}{D(z)} - 1 \right) \underbrace{Y(z)}_{\text{measurements}} \quad (4.132)$$

Let us consider two simple examples of application of this formula:

- If we want noise model $\frac{C(z)}{D(z)} = 1$ (output noise), we need to use

$$\underbrace{\hat{Y}(z)}_{\text{prediction}} = \frac{B(z)}{A(z)} U(z) \quad (4.133)$$

i.e.

$$A(z) \underbrace{\hat{Y}(z)}_{\text{prediction}} = B(z)U(z) \quad (4.134)$$

This is a pure simulation model (not using past measurements)

- If we want noise model $\frac{C(z)}{D(z)} = \frac{B(z)}{A(z)}$ (input noise), we need to use

$$\frac{B(z)}{A(z)} \underbrace{\hat{Y}(z)}_{\text{prediction}} = \frac{B(z)}{A(z)} U(z) + \left(\frac{B(z)}{A(z)} - 1 \right) \underbrace{Y(z)}_{\text{measurements}}$$

i.e.

$$\text{i.e. } B(z) \underbrace{\hat{Y}(z)}_{\text{prediction}} = B(z)U(z) + (B(z) - A(z)) \underbrace{Y(z)}_{\text{measurements}}$$

5 DIFFERENTIAL ALGEBRAIC EQUATIONS (DAEs)

We have seen so far formally ordinary differential equations (ODEs), which describe the time evolution of a vector of variables via relationships of the form:

$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t), t) \quad (5.1)$$

where $\mathbf{x} \in \mathbb{R}^{n_x}$ is a vector of *differential states* and $\mathbf{u} \in \mathbb{R}^{n_u}$ a vector of “external variables”, often referred to as inputs (any variable that is not defined by (5.1), but rather “externally” can be classified as an “input”). For the sake of simplicity, the time dependence (t) of these variables is usually omitted.

In this chapter, we will approach a new form of differential equations, labelled *Differential-Algebraic Equations* (DAEs), which are used a lot in the modelling of complex and large-scale systems. DAEs have features and properties that are unlike ODEs. We will investigate them here.

5.1 WHAT ARE DAEs?

Simply put, DAEs are a set of equations that do not define directly the entire state. This happens e.g. if we have a mix of differential equations and purely algebraic equations (equations where the time-differentiated variables do not appear) or if the equations mix *differential states*, i.e. variables that appear as time-differentiated in the equations, and *algebraic states*, i.e. variables (excluding inputs and parameters) whose time derivative does not appear at all. Before introducing the formal definition of a DAE, let us make some simple examples.

Example

- Consider the differential equation

$$\dot{\mathbf{x}}_1 = 2\mathbf{x}_1 + \mathbf{x}_2 \quad (5.2a)$$

$$0 = 3\mathbf{x}_1 - \mathbf{x}_2 \quad (5.2b)$$

Here we clearly observe that the state variable \mathbf{x}_2 does not appear in its time-differentiated form $\dot{\mathbf{x}}_2$, such that the equations do not define $\dot{\mathbf{x}}_2$, but rather \mathbf{x}_2 . Put differently, \mathbf{x}_2 is not a differential state but an algebraic state. Hence (5.2) is a DAE.

- Consider the differential equation

$$\dot{\mathbf{x}}_1 + 2\mathbf{x}_1 + \dot{\mathbf{x}}_2 + \mathbf{x}_2 = 0 \quad (5.3a)$$

$$2\dot{\mathbf{x}}_1 + \mathbf{x}_1 + 2\dot{\mathbf{x}}_2 + 2\mathbf{x}_2 = 0 \quad (5.3b)$$

here both variables \dot{x}_1 and \dot{x}_2 appear time-differentiated. However, one can also observe that replacing (5.3b) by $-2 \cdot (5.3a) + (5.3b)$ yields:

$$\dot{x}_1 + 2x_1 + \dot{x}_2 + x_2 = 0 \quad (5.4a)$$

$$-3x_1 = 0 \quad (5.4b)$$

Here we observe that (5.4) does not define the differential states \dot{x}_1 and \dot{x}_2 (in the sense that one cannot solve (5.4) to obtain \dot{x}_1, \dot{x}_2). However, (5.4) still provides a well-defined trajectory for x_1 and x_2 . Indeed, (5.4b) specifies that $x_1 = 0$, such that (5.4a) reads as the simple ODE:

$$\dot{x}_2 + x_2 = 0 \quad (5.5)$$

which can be solved for x_2 .

□

In order for the notion of DAE to be rigorous, we need a more formal definition.

Definition 1. consider a differential equation having the state vector $s \in \mathbb{R}^{n_s}$, and defined by the equation

$$F(\dot{x}, x, u, t) = 0 \quad (5.6)$$

The differential equation (5.6) is a DAE if

$$\det\left(\frac{\partial F}{\partial \dot{x}}\right) = 0 \quad (5.7)$$

i.e. if the Jacobian $\frac{\partial F}{\partial \dot{x}}$ is rank-deficient.

To clarify this definition, let us test it on (5.2) and (5.3).

Example

- Differential equation (5.2) can be written as

$$F(\dot{x}, x) = \begin{bmatrix} \dot{x}_1 - 2x_1 - x_2 \\ 3x_1 - x_2 \end{bmatrix} = 0 \quad (5.8a)$$

such that

$$\frac{\partial F}{\partial \dot{x}} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \quad (5.9)$$

is rank-deficient because the second column (row) of the Jacobian of F is zero. One can generalize this observation: algebraic states yield columns of zeros in the Jacobian $\frac{\partial F}{\partial \dot{x}}$ and make it rank-deficient.

- Differential equation (5.3) can be written as

$$F(\dot{x}, x) = \begin{bmatrix} \dot{x}_1 + 2x_1 + \dot{x}_2 + x_2 \\ 2\dot{x}_1 + x_1 + 2\dot{x}_2 + 2x_2 \end{bmatrix} = 0 \quad (5.10)$$

such that

$$\frac{\partial F}{\partial \dot{x}} = \begin{bmatrix} 1 & 1 \\ 2 & 2 \end{bmatrix} \quad (5.11)$$

is rank-deficient because its determinant is zero. Here the rank-deficient determinant captures the fact that (5.3) does not define the differential state properly.

□

In practice, most often DAEs are arising because some of the states are purely algebraic (i.e. they do not appear time-differentiated), as in (5.2). In order to stress the difference between differential and algebraic states, it is common to use the notation x for the differential states, and z the algebraic states. E.g. (5.2) would be written as:

$$\dot{x} = 2x + z \quad (5.12a)$$

$$0 = 3x - z \quad (5.12b)$$

Definition 1 allows one to assess whether a differential equation is a DAE or an ODE in ambiguous cases. However, one can trivially write differential equations that cannot be categorized so simply. Let us consider two simple examples.

Example Let us investigate two “freak” differential equations that can switch between being ODEs and DAEs

- Consider the simple scalar differential equation

$$u\dot{x} + x = 0 \quad (5.13)$$

here we observe that $F = u\dot{x} + x$ and

$$\frac{\partial F}{\partial \dot{x}} = u \quad (5.14)$$

Hence for $u \neq 0$ (5.13) is an ODE and reads as:

$$\dot{x} = -\frac{x}{u} \quad (5.15)$$

while for $u = 0$ (5.13) is an “DAE” (in fact it is purely algebraic) and reads as:

$$x = 0 \quad (5.16)$$

Hence (5.13) can switch between being an ODE and a DAE depending on the input u .

- Consider the differential equation

$$\dot{x}_1 + x_1 - u = 0 \quad (5.17a)$$

$$(x_1 - x_2) \dot{x}_2 + x_1 - x_2 = 0 \quad (5.17b)$$

One can verify that

$$\det\left(\frac{\partial F}{\partial \dot{x}}\right) = x_1 - x_2 \quad (5.18)$$

Hence for initial conditions $x(0)$ satisfying $x_1(0) = x_2(0)$, (5.17) starts as a DAE and obeys the dynamics:

$$\dot{x}_1 + x_1 - u = 0 \quad (5.19a)$$

$$x_1 - x_2 = 0 \quad (5.19b)$$

hence it enforces $x_1(t) = x_2(t)$ for all time, and (5.17) remains a DAE. However, if the initial conditions satisfy $x_1(0) \neq x_2(0)$, then (5.17) starts as an ODE and remains an ODE throughout its trajectories.

These examples are meant to draw the attention to the fact that the notion of DAE can be fairly convoluted. In most practical cases, however, DAEs arise because the differential equations hold variables that do not appear time-differentiated. In this context, the problem of DAEs having exotic behaviors as in the examples above does not arise.

□

5.2 DIFFERENT FORMS OF DAEs

DAEs come in different forms, which are useful to recognize, as these different forms can have different intrinsic properties. We will briefly look at these forms here.

- **Fully-implicit DAEs** are in the form (5.6) of definition 1, i.e. they read as

$$F(\dot{x}, x, u) = 0 \quad (5.20)$$

where

$$\det\left[\frac{\partial F}{\partial \dot{x}}\right] = 0 \quad (5.21)$$

If the rank-deficiency arises because some the states x do not appear as time-differentiated in F (hence creating columns of zeros in the Jacobian $\frac{\partial F}{\partial \dot{x}}$), then they are commonly labelled as "z" states, and (5.20) is rewritten as:

$$F(\dot{x}, x, z, u) = 0 \quad (5.22)$$

Note that here condition (5.21) becomes

$$\det\left(\begin{bmatrix} \frac{\partial F}{\partial \dot{x}} & \frac{\partial F}{\partial \dot{z}} \end{bmatrix}\right) = \det\left(\begin{bmatrix} \frac{\partial F}{\partial \dot{x}} & 0 \end{bmatrix}\right) = 0 \quad (5.23)$$

hence (5.22) is by construction a DAE as \dot{z} is not an argument in F .

- **Semi-explicit DAEs** split explicitly the differential and algebraic equations, and can generally be written in the form:

$$\dot{x} = f(x, z, u) \quad (5.24a)$$

$$0 = g(x, z, u) \quad (5.24b)$$

A semi-explicit DAE can be trivially written as a fully-implicit DAE by defining:

$$F(\dot{x}, x, z, u) = \begin{bmatrix} \dot{x} - f(x, z, u) \\ g(x, z, u) \end{bmatrix} = 0 \quad (5.25)$$

Conversely, a fully implicit can be trivially written as a semi-explicit DAE by introducing some “helper variables” (labelled v here), i.e.

$$\dot{x} = v \quad (5.26a)$$

$$0 = F(v, x, z, u) \quad (5.26b)$$

We observe that (5.24a) is a DAE by construction. Indeed, one can apply definition 1 on (5.25) (using the same construction as in (5.23)) and observe that:

$$\det \left(\begin{bmatrix} \frac{\partial F}{\partial x} & \frac{\partial F}{\partial z} \end{bmatrix} \right) = \det \left(\begin{bmatrix} I & 0 \\ 0 & 0 \end{bmatrix} \right) = 0 \quad (5.27)$$

- **Linear DAEs** are DAEs where the underlying functions are linear. They are often put in the form:

$$E\dot{x} = Ax + Bu \quad (5.28)$$

where definition 1 requires that E is rank-deficient in order for (5.28) to be a DAE. In a semi-explicit form, and underlying the distinction between algebraic and differential states, a linear DAE would read as:

$$\dot{x} = Ax + Bu + Cz \quad (5.29a)$$

$$0 = Dx + Eu + Fz \quad (5.29b)$$

For matrix F full rank, we observe that the algebraic variables z can be eliminated using (5.29b) and replaced into (5.29a) such that the DAE can be reduced to an ODE with the addition of some “output variables” z :

$$\dot{x} = (A - CF^{-1}D)x + (B - CF^{-1}E)u \quad (5.30a)$$

$$z = -F^{-1}(Dx + Eu) \quad (5.30b)$$

Similarly to ODEs, nonlinear DAEs can be linearized into linear DAEs. E.g. a fully-implicit DAE of the form (5.22) can be linearized to

$$\frac{\partial F}{\partial \dot{x}} \Delta \dot{x} + \frac{\partial F}{\partial x} \Delta x + \frac{\partial F}{\partial z} \Delta z + \frac{\partial F}{\partial u} \Delta u = 0 \quad (5.31)$$

and a semi-explicit DAE of the form can be linearized to

$$\Delta \dot{x} = \frac{\partial f}{\partial x} \Delta x + \frac{\partial f}{\partial u} \Delta u + \frac{\partial f}{\partial z} \Delta z \quad (5.32a)$$

$$0 = \frac{\partial g}{\partial x} \Delta x + \frac{\partial g}{\partial u} \Delta u + \frac{\partial g}{\partial z} \Delta z \quad (5.32b)$$

and therefore takes the form (5.29). It is interesting to note here that the elimination of the algebraic variables Δz in (5.32) is possible if the Jacobian $\frac{\partial g}{\partial z}$ is full rank throughout the trajectory. We will get back to this notion soon.

5.3 TIKHONOV THEOREM

Most physical systems comprise some dynamics that are very fast in comparison with the main dynamics of the system. During the modelling of the system, these very fast dynamics are often disregarded. E.g. when modelling mechanical systems, many components of the system are considered as infinitely rigid, and their dynamics dismissed. For electrical systems, the inductance and capacitance of the cables are often neglected. Experienced engineers typically know which dynamics matter and ought to be included in their model, and which ones can be safely neglected. While the elimination of very fast dynamics is most often performed based on the intuition and experience of the engineer, it has strong foundations in mathematics. The elimination of fast dynamics is one of the reasons why DAEs arise in modelling. In this section we will see how and why eliminating fast dynamics from a model is justified.

Theorem 8. Consider the ordinary differential equation (ODE):

$$\dot{x} = f(x, z) \quad (5.33a)$$

$$\epsilon \dot{z} = g(x, z) \quad (5.33b)$$

where $0 < \epsilon \ll 1$ is very small. Let us label $x_\epsilon(t), z_\epsilon(t)$ the solution to (5.33). Suppose:

- the dynamics $\dot{z} = g(x, z)$ are stable $\forall x$
- matrix $\frac{\partial g}{\partial z}$ is full rank (i.e. invertible) everywhere

then

$$\lim_{\epsilon \rightarrow 0} x_\epsilon(t), z_\epsilon(t) = x_0(t), z_0(t) \quad (5.34)$$

where $x_0(t), z_0(t)$ is the solution of

$$\dot{x} = f(x, z) \quad (5.35a)$$

$$0 = g(x, z) \quad (5.35b)$$

The limit (5.34) is to be understood in the sense of "almost everywhere" formally defined in the context of measure theory. From an intuitive point of view, the trajectories match everywhere but on a union of time intervals that are "infinitely short". An illustration of the Tikhonov theorem is provided in the following example.

Examples consider the linear system

$$\dot{x}_1 = -x_1 + x_2 \quad (5.36a)$$

$$\epsilon \dot{x}_2 = x_1 - ax_2 + u \quad (5.36b)$$

The hypothesis of Theorem 8 are satisfied for (5.36) if $a > 0$, such that $\frac{\partial g}{\partial z} = -a$ is full rank and (5.36b) is stable. The limit case $\epsilon = 0$ reads as

$$\dot{x}_1 = -x_1 + x_2 \quad (5.37a)$$

$$0 = x_1 - ax_2 + u \quad (5.37b)$$

We then distinguish two cases

- if $a \neq 0$ then (5.37) has the solution

$$\dot{x}_1 = (a^{-1} - 1)x_1 + a^{-1}u \quad (5.38a)$$

$$x_2 = a^{-1}(x_1 + u) \quad (5.38b)$$

- if $a = 0$ then (5.37) has the solution

$$\dot{x}_1 = -x_1 + x_2 \quad (5.39a)$$

$$0 = x_1 + u \quad (5.39b)$$

which equivalently reads:

$$x_2 = -(u + \dot{u}) \quad (5.40a)$$

$$x_1 = -u \quad (5.40b)$$

DAE (5.39) requires that u is continuous as it appears time-differentiated. This is in contrast with (5.36) for which any integrable input profile is admissible. Additionally, we observe that (5.36) is unstable for $a = 0$ (with a pole at 0), and therefore does not satisfy the hypothesis of Theorem 8. The trajectories are illustrated in Fig. 5.1.

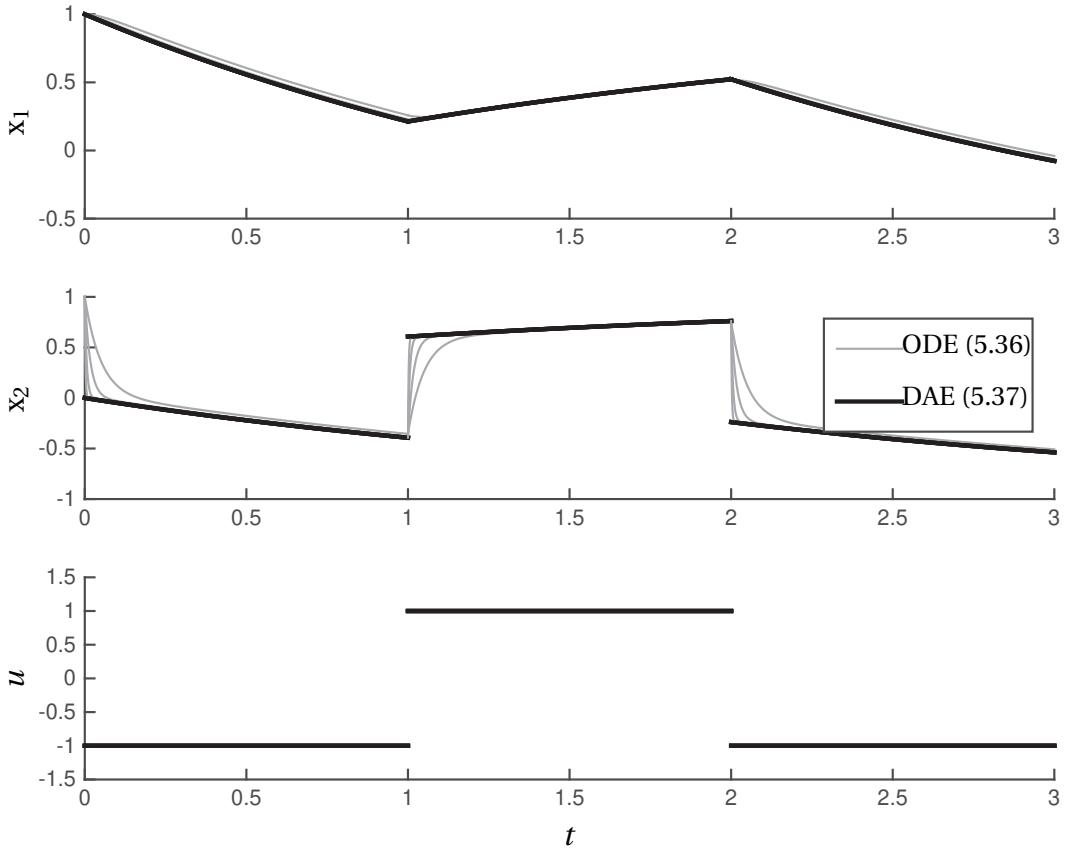


Figure 5.1: Illustration of the Tikhonov theorem for (5.36). The piecewise constant input illustrated in the third graph was selected. The grey curves represent the trajectories for ϵ ranging from 10^{-1} to $4 \cdot 10^{-3}$ with the initial conditions $x_1(0) = x_2(0) = 1$ and for $a = 2$. The black curve represents the solution to the DAE (5.37). Here we see that the Tikhonov theorem applies.

- Let us consider two masses moving horizontally (positions x_1 and x_2), subject to a velocity-dependent friction with parameter ρ , and connected by elastic links of constant K . The dynamics of this two-mass system can be modelled via the linear ODE:

$$m_1 \ddot{x}_1 = K(x_2 - x_1) - Kx_1 - \rho \dot{x}_1 + u \quad (5.41a)$$

$$m_2 \ddot{x}_2 = K(x_1 - x_2) - \rho \dot{x}_2 \quad (5.41b)$$

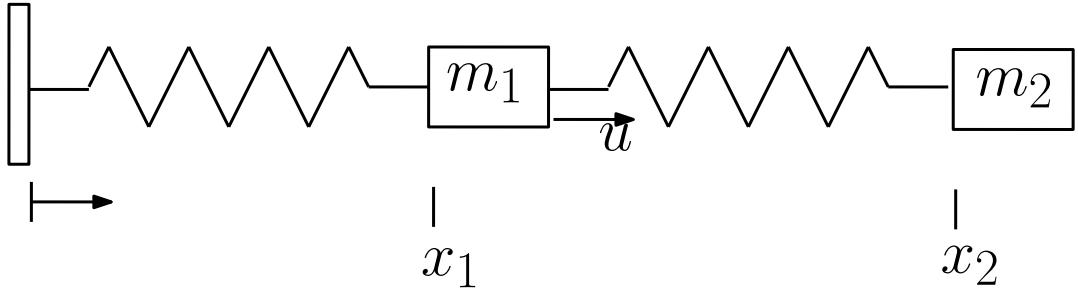


Figure 5.2: Illustration of the two mass system modelled by (5.41).

Let us assume that $m_1 \ll m_2$ and $\rho \ll K$, such that the oscillatory dynamics of mass 1 are very fast compared to the one of mass 2. Let us then consider the approximation:

$$0 = K(x_2 - x_1) - Kx_1 - \rho\dot{x}_1 + u \quad (5.42a)$$

$$m_2\ddot{x}_2 = K(x_1 - x_2) - \rho\dot{x}_2 \quad (5.42b)$$

of (5.41) where m_1 being relatively small has been set to zero. ODE (5.42) is equivalent to:

$$\frac{\rho}{K}\dot{x}_1 = (x_2 - x_1) - x_1 + u \quad (5.43a)$$

$$m_2\ddot{x}_2 = K(x_1 - x_2) - \rho\dot{x}_2 \quad (5.43b)$$

Since $\frac{\rho}{K} \ll 1$, we can form the further approximation:

$$0 = (x_2 - x_1) - x_1 + u \quad (5.44a)$$

$$m_2\ddot{x}_2 = K(x_1 - x_2) - \rho\dot{x}_2 \quad (5.44b)$$

which is a DAE since state x_2 does not appear time differentiated. We can easily observe that (5.44) boils down to eliminating \dot{x}_1 , \ddot{x}_1 from equation (5.41a). I.e. while we keep considering x_1 as time-varying, we dismiss its dynamics.

Figure 5.3 compares the trajectories of ODE (5.41) and the ones of the DAE (5.44). One can verify that for both applications of the Tikhonov theorem (5.42) and (5.44), the hypothesis apply.

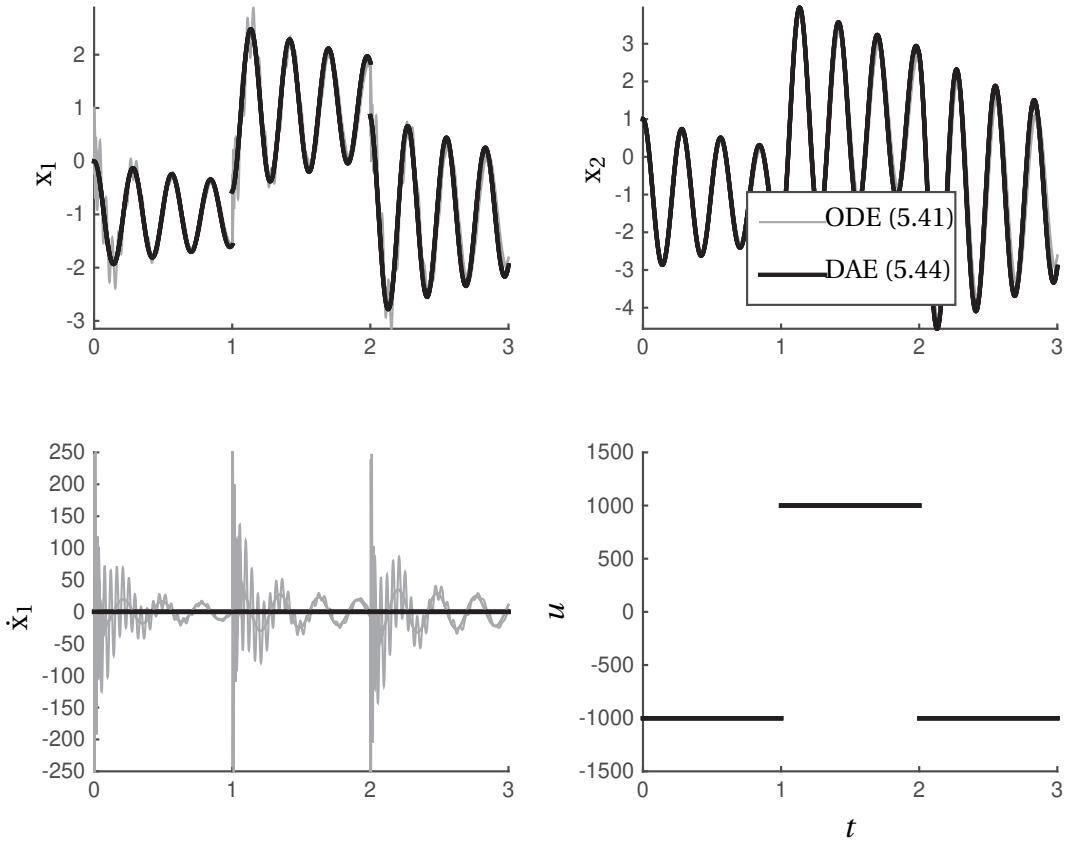


Figure 5.3: Illustration of the Tikhonov theorem for (5.41). The piecewise constant input illustrated in the lower-right graph was selected. The grey curves represent the trajectories for m_1 ranging from 10^{-1} to $4 \cdot 10^{-4}$ with the initial conditions $x_1(0) = x_2(0) = 1$ and $\dot{x}_{1,2}(0) = 0$. The parameters $\rho = 1$, $K = 10^3$ and $m_2 = 1$ were selected. The black curve represents the solution to the DAE (5.44). Here we see that the Tikhonov theorem applies.

□

To conclude these two examples, it is useful to provide an intuition behind the application of the Tikhonov theorem, which is the basis of the intuition of engineers in neglecting some of the dynamics of a system when developing a mathematical model. The intuition is easiest to build for the two-mass example. One can observe that the approximation at play here is essentially to replace the dynamics of mass 1 by its steady-state approximation. I.e. we consider that the motion of the mass decays instantaneously to its steady state provided by (5.44a). One should understand here that “steady-state” does not mean “not moving”, as e.g. (5.44a) still allows x_1 to move. However, in the approximation (5.44) the motion of x_1 is entirely dictated by the motion of x_2 moving x_1 in different steady-state positions. This principle can very often be readily applied to models holding very fast dynamics. One ought to observe that the dynamics created by the DAE (5.44) do not hold very fast dynamics (as opposed to e.g. (5.42) which still hold a very fast dynamics for x_1). This observation

will connect to Sections 6.4.

5.4 DIFFERENTIAL INDEX OF DAEs

DAEs are unlike ODEs in various ways. One very important distinction is that DAEs hold the concept of differential index, which is crucial when it comes to solving them numerically. While the notion of differential index is not straightforward, the problem of high-index DAE, which is the main focus of this section, can be construed fairly intuitively.

One ought to understand that in order for a DAE e.g. in the semi-explicit form:

$$\dot{x} = f(x, z, u) \quad (5.45a)$$

$$0 = g(x, z, u) \quad (5.45b)$$

to be "solvable", one need to be able to compute the differential state derivative \dot{x} and the algebraic states z for a given differential state x and input u . Indeed, \dot{x}, z can be readily obtained (possibly numerically) from (5.45) for any given x, u , then one can process the system trajectories and build the trajectories corresponding to (5.45). In order to build this intuition, let us consider the following example.

Example

- Let us start with a simple linear DAE:

$$F = \begin{bmatrix} \dot{x}_1 - z \\ \dot{x}_2 - x_1 \\ \dot{x}_1 + x_2 - u \end{bmatrix} = 0 \quad (5.46)$$

One can verify that, using simple algebraic manipulations, (5.46) can be rewritten as

$$z = u - x_2 \quad (5.47a)$$

$$\dot{x}_2 = x_1 \quad (5.47b)$$

$$\dot{x}_1 = u - x_2 \quad (5.47c)$$

i.e. (5.46) delivers \dot{x} and z as a function of x and u .

- Consider a linear DAE similar to but slightly different than (5.46)

$$F = \begin{bmatrix} \dot{x}_1 - z \\ \dot{x}_2 - x_1 \\ x_2 - u \end{bmatrix} = 0 \quad (5.48)$$

One can verify that no algebraic manipulation can deliver \dot{x} and z from (5.48). Indeed, (5.48) also reads as:

$$\underbrace{\begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}}_{:=M} \begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ z \end{bmatrix} = \begin{bmatrix} 0 \\ x_1 \\ x_2 - u \end{bmatrix} \quad (5.49)$$

and since matrix M is rank deficient (last row is zero), one cannot manipulate the equations defined by (5.48) in order to extract \dot{x} and z . In contrast, we observe that (5.46) can be rewritten as

$$\begin{bmatrix} 1 & 0 & -1 \\ -1 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ z \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ x_2 - u \end{bmatrix} \quad (5.50)$$

and its matrix is full rank, allowing us to build (5.47). We ought to observe that the matrix M is obtained by taking the Jacobian

$$\left[\frac{\partial F}{\partial \dot{x}} \quad \frac{\partial F}{\partial z} \right] \quad (5.51)$$

which need to be full rank in order for the DAEs to be solvable for \dot{x} and z .

□

The principles deployed in the examples above can be generalized to nonlinear DAEs, as is stated next.

Theorem 9. *a fully implicit DAE*

$$F(\dot{x}, z, x, u) = 0 \quad (5.52)$$

with function F smooth can be readily solved (i.e. solved for \dot{x}, z) if the Jacobian

$$\left[\frac{\partial F}{\partial \dot{x}} \quad \frac{\partial F}{\partial z} \right] \quad (5.53)$$

is full rank on all trajectories \dot{x}, z, x, u .

Proof. this theorem follows directly from the IFT (see Theorem 5). Indeed, for (5.53) full rank, there exist functions $\dot{x}(x, u)$ and $z(x, u)$ such that

$$F(\dot{x}(x, u), z(x, u), x, u) = 0 \quad (5.54)$$

hold locally for any x, u . It follows that (5.52) can be solved. Another way of construing this result is that if the Jacobian (5.53) is full rank, then a Newton iteration deployed on (5.52) in order to compute \dot{x}, z converges locally (because the linear system (3.8) is well-posed). □

One can readily apply Theorem 9 to semi-explicit DAEs.

Corollary 1. *a semi-explicit DAE*

$$\dot{x} = f(x, z, u) \quad (5.55a)$$

$$0 = g(x, z, u) \quad (5.55b)$$

with function g smooth can be readily solved (i.e. solved for \dot{x}, z) if the Jacobian

$$\frac{\partial g}{\partial z} \quad (5.56)$$

is full rank on all trajectories z, x, u .

Proof. recall that a semi-explicit DAE can be transformed into a full implicit one via the transformation (5.25) recalled here:

$$F(\dot{x}, x, z, u) = \begin{bmatrix} \dot{x} - f(x, z, u) \\ g(x, z, u) \end{bmatrix} = 0 \quad (5.57)$$

A direct application of Theorem 9 then yields the Jacobian:

$$\begin{bmatrix} \frac{\partial F}{\partial \dot{x}} & \frac{\partial F}{\partial z} \end{bmatrix} = \begin{bmatrix} I & \frac{\partial f}{\partial z} \\ 0 & \frac{\partial g}{\partial z} \end{bmatrix} \quad (5.58)$$

which is full rank if $\frac{\partial g}{\partial z}$ is full rank. \square

An intuitive way of construing Corollary 1 is by observing that if the Jacobian (5.56) is full rank, then the algebraic equation (5.55b) can be solved for z at any point x, u (see IFT, Theorem 5). The solution z can then be injected in (5.55a) to obtain \dot{x} .

Theorem 9 and its corollary 1 tell us that there are DAE that are “easy” to solve (those satisfying the full rankness of the Jacobians (5.53) and (5.56)) and DAEs that are “hard” to solve (those where the Jacobians (5.53) and (5.56) are rank deficient). We will see next that Theorem 9 and its corollary 1 have to do with the differential index of DAEs.

Definition 2. *the differential index of a DAE is the number of times the operator $\frac{d}{dt}$ must be applied to the equations (+ possibly an arbitrary amount of algebraic manipulations) in order to convert the DAE into an ODE.*

Definition 2 is non-trivial and can require some care when being applied. Let us make a few examples in order to gather some intuitions on how it works.

Example

- Consider the linear DAE (5.47) repeated here:

$$z = u - x_2 \quad (5.59a)$$

$$\dot{x}_2 = x_1 \quad (5.59b)$$

$$\dot{x}_1 = u - x_2 \quad (5.59c)$$

Recall that (5.47) is an “easy” DAE since it satisfies the assumption of Theorem 9. We also observe that a single application of $\frac{d}{dt}$ on (5.59a) yields:

$$\dot{z} = \dot{u} - \dot{x}_2 \quad (5.60a)$$

$$\dot{x}_2 = x_1 \quad (5.60b)$$

$$\dot{x}_1 = u - x_2 \quad (5.60c)$$

or equivalently

$$\dot{z} = \dot{u} - x_1 \quad (5.61a)$$

$$\dot{x}_2 = x_1 \quad (5.61b)$$

$$\dot{x}_1 = u - x_2 \quad (5.61c)$$

The latter differential equation is an ODE, hence (5.47) is of index 1. It is interesting to note here that by performing the transformation *DAE* → *ODE* (from (5.47) to (5.61)), we have restricted the functional space for u . Indeed, u only needs to be integrable in (5.47), while it needs to be differentiable in (5.61).

- Consider the linear DAE (5.48), repeated here

$$\dot{x}_1 - z = 0 \quad (5.62a)$$

$$\dot{x}_2 - x_1 = 0 \quad (5.62b)$$

$$x_2 - u = 0 \quad (5.62c)$$

Recall that (5.48) is a “hard” DAE since it fails the assumption of Theorem 9. A time differentiation of (5.48) (last row) yields

$$\dot{x}_1 - z = 0 \quad (5.63a)$$

$$\dot{x}_2 - x_1 = 0 \quad (5.63b)$$

$$\dot{x}_2 - \dot{u} = 0 \quad (5.63c)$$

An algebraic manipulation yields

$$\dot{x}_1 - z = 0 \quad (5.64a)$$

$$\dot{u} - x_1 = 0 \quad (5.64b)$$

$$\dot{x}_2 - \dot{u} = 0 \quad (5.64c)$$

A second time differentiation applied on (5.64b) yields

$$\dot{x}_1 - z = 0 \quad (5.65a)$$

$$\ddot{u} - \dot{x}_1 = 0 \quad (5.65b)$$

$$\dot{x}_2 - \dot{u} = 0 \quad (5.65c)$$

and an algebraic manipulation yields:

$$\ddot{u} - z = 0 \quad (5.66a)$$

$$\ddot{u} - \dot{x}_1 = 0 \quad (5.66b)$$

$$\dot{x}_2 - \dot{u} = 0 \quad (5.66c)$$

A third time differentiation applied to (5.66a) yields the ODE

$$\dot{z} = \ddot{u} \quad (5.67a)$$

$$\dot{x}_1 = \ddot{u} \quad (5.67b)$$

$$\dot{x}_2 = \dot{u} \quad (5.67c)$$

A total count of 3 time differentiations was used to transform the DAE (5.48) to the ODE (5.67). It follows that DAE (5.48) is of index 3.

□

Note that the same principles can be applied to nonlinear DAEs (see e.g. section 5.5 below).

Let us make now the connection between the concept of differential index and Theorem 9 and its corollary 9.

Theorem 10. *a fully-implicit DAE*

$$F(\dot{x}, z, x, u) = 0 \quad (5.68)$$

fulfils the assumption of Theorem 9 if and only if it is of index 1.

Proof. for an index-1 fully-implicit DAE, a single time differentiation allows a transformation to an ODE, i.e.

$$\frac{d}{dt} F = \frac{\partial F}{\partial \dot{x}} \ddot{x} + \frac{\partial F}{\partial x} \dot{x} + \frac{\partial F}{\partial z} \dot{z} + \frac{\partial F}{\partial u} \dot{u} = 0 \quad (5.69)$$

is an ODE. For the sake of clarity, let us introduce the state extension $v = \dot{x}$, such that (5.69) reads as:

$$\frac{\partial F}{\partial \dot{x}} \dot{v} + \frac{\partial F}{\partial x} v + \frac{\partial F}{\partial z} \dot{z} + \frac{\partial F}{\partial u} \dot{u} = 0 \quad (5.70)$$

or equivalently (5.69) reads as

$$\begin{bmatrix} \dot{v} \\ \dot{z} \end{bmatrix} = - \left[\begin{array}{cc} \frac{\partial F}{\partial \dot{x}} & \frac{\partial F}{\partial z} \end{array} \right]^{-1} \left(\frac{\partial F}{\partial x} v + \frac{\partial F}{\partial u} \dot{u} \right) \quad (5.71)$$

and the Jacobian matrix $\left[\begin{array}{cc} \frac{\partial F}{\partial \dot{x}} & \frac{\partial F}{\partial z} \end{array} \right]$ must be full rank. □

The same result extends to semi-explicit DAEs, i.e.

Corollary 2. *a semi-explicit DAE*

$$\dot{x} = f(x, z, u) \quad (5.72a)$$

$$0 = g(x, z, u) \quad (5.72b)$$

fulfils the assumption of Theorem 1 if and only if it is of index 1.

Proof. if (5.72) is of index 1, then a single time differentiation on the algebraic equation yields an ODE, i.e.

$$\dot{x} = f(x, z, u) \quad (5.73a)$$

$$0 = \frac{d}{dt} g(x, z, u) \quad (5.73b)$$

is an ODE, i.e. equivalently, using a chain rule

$$\dot{x} = f(x, z, u) \quad (5.74a)$$

$$0 = \frac{\partial g}{\partial x} \dot{x} + \frac{\partial g}{\partial z} \dot{z} + \frac{\partial g}{\partial u} \dot{u} \quad (5.74b)$$

is an ODE. It follows that

$$\dot{x} = f(x, z, u) \quad (5.75a)$$

$$\dot{z} = -\frac{\partial g^{-1}}{\partial z} \left(\frac{\partial g}{\partial x} \dot{x} + \frac{\partial g}{\partial u} \dot{u} \right) \quad (5.75b)$$

is an ODE, which requires the Jacobian $\frac{\partial g}{\partial z}$ to be full rank. \square

The message to take home from Theorem 10 and corollary 2 is that index-1 DAEs are “easy” to solve in the sense that the equations readily deliver \dot{x} and z . DAEs of index more than 1 are generally referred to as “high-index” DAEs, and are notoriously more difficult to handle. When approaching DAEs numerically, index-1 DAEs are clearly preferred. This does not mean than high-index DAEs cannot be treated, but they are often best treated via a so-called index-reduction procedure, which we will introduce in Section 5.5 and further discuss in Section 5.6.

5.5 CONNECTION TO LAGRANGE MECHANICS

An attentive reader will have noticed that we have already approached DAEs in Section 2.5, when discussing constrained Lagrange mechanics. Indeed equation (2.161) stemming from constrained Lagrange mechanics:

$$\frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{q}} - \frac{\partial \mathcal{L}}{\partial q} = Q \quad (5.76a)$$

$$c(q) = 0 \quad (5.76b)$$

is in fact a semi-explicit DAE, where (5.76b) is the algebraic equation and (5.76a) yields (after minor treatments) an explicit ODE.

We have seen that (5.76) does not readily deliver \ddot{q} , z , and that it ought to be modified by time-differentiating the constraint equation (5.76b) twice, delivering the model (2.166) recalled here:

$$\frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{q}} - \frac{\partial \mathcal{L}}{\partial q} = Q \quad (5.77a)$$

$$\ddot{c}(q, \dot{q}, \ddot{q}) = 0 \quad (5.77b)$$

It is interesting here to verify the differential index of (5.76) by applying definition 2. The differential state for a mechanical system reads as

$$x = \begin{bmatrix} q \\ \dot{q} \end{bmatrix} \quad (5.78)$$

The transformed Lagrange equation (5.77) can be written explicitly as:

$$\begin{bmatrix} I & 0 & 0 \\ 0 & W(q) & 0 \\ 0 & 0 & M(q) \end{bmatrix} \begin{bmatrix} \dot{x} \\ z \end{bmatrix} = \begin{bmatrix} \dot{q} \\ Q - \frac{\partial}{\partial q}(W(q)\dot{q})\dot{q} + \nabla_q \mathcal{L} \\ -\frac{\partial}{\partial q}\left(\frac{\partial c}{\partial q}\dot{q}\right)\dot{q} \end{bmatrix} \quad (5.79)$$

One can observe that (5.79) readily delivers \dot{x}, z for $W(q), M(q)$ full rank. It follows that (5.79) fulfills the assumption of Theorem 9, and is therefore of index 1. Since (5.79) has been obtained after two time differentiations of (5.76), it follows that (5.76) is an **index-3 DAE**. We can summarize this observation as:

$$\underbrace{(5.76b)}_{\text{index 3}} \xrightarrow{\frac{d^2}{dt^2}} \underbrace{(5.77b)}_{\text{index 1}} \xrightarrow{\frac{d}{dt}} \text{ODE}$$

This observation is generic, i.e. models arising from constrained Lagrange mechanics with position-dependent constraints of the form $c(q) = 0$ yield index-3 DAEs, and two time differentiations of the constraints $c(q) = 0$ yield an index-1 DAE.

The important message to take home here is that a high-index DAE such as (5.76) can be transformed into an index-1 DAE (i.e. (5.77)) via time differentiations and algebraic manipulations. The trick we have explored to transform the DAEs arising from constrained Lagrange mechanics into low-index DAEs is not limited to this special case, but it a generic approach to transform "hard" (i.e. high-index) DAEs into "easy" ones (i.e. low-index DAEs), commonly referred to as an **index-reduction** procedure. We further detail this in the next section.

5.6 INDEX REDUCTION

The index reduction of DAEs is identical to assessing their index as per definition 2, with the minor difference that the procedure ought to be stopped one step before reaching an ODE (i.e. at the index-1 step). More specifically, one ought to perform time differentiations and algebraic manipulations until an index-1 DAE emerges. It is difficult to automate this procedure in general (although some software packages offer index-reduction capabilities), as the algebraic manipulations to be performed typically require some insights into the DAE. However, we can attempt a detailed description of the procedure.

Let us consider a semi-explicit DAE of the form

$$\dot{x} = f(x, z, u) \quad (5.80a)$$

$$0 = g(x, z, u) \quad (5.80b)$$

The index-reduction procedure can be described as follows.

1. Check if the DAE system is of index 1. If yes, stop.

2. Identify a subset of algebraic equations in (5.80b) that can be solved for a subset of algebraic variables z .
3. Apply $\frac{d}{dt}$ on the remaining algebraic equations containing some differential states x_j , this leads to terms \dot{x}_j appearing in these differentiated equations.
4. Substitute the terms \dot{x}_j by the corresponding expressions $f_j(x, z, u)$, this delivers new algebraic equations to replace those differentiated in (5.80b).
5. With this new DAE system, go to step 1.

This procedure is not always straightforward to deploy. Let us consider an example to understand it better.

Example

- consider the linear semi-explicit DAE

$$\dot{x} = \begin{bmatrix} x_2 + z_2 \\ z_1 + u \end{bmatrix} \equiv f \quad (5.81a)$$

$$0 = \begin{bmatrix} x_1 - x_2 \\ z_2 - x_1 \end{bmatrix} \equiv g \quad (5.81b)$$

We observe that

$$\frac{\partial g}{\partial z} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \quad (5.82)$$

is rank deficient, such that (5.81) is of index > 1 (see Corollary 2). We observe that the second equation of (5.81b) delivers z_2 as:

$$z_2 = x_1 \quad (5.83)$$

We then apply $\frac{d}{dt}$ on the first equation in (5.81b) to obtain:

$$\dot{x}_1 - \dot{x}_2 = 0 \quad (5.84)$$

We use (5.81a) to replace $\dot{x}_{1,2}$ by their corresponding expressions, delivering a new algebraic equation:

$$x_2 + z_2 - z_1 - u = 0 \quad (5.85)$$

We obtain the new DAE where the differentiated algebraic equation has been replaced by (5.85):

$$\dot{x} = \begin{bmatrix} x_2 + z_2 \\ z_1 + u \end{bmatrix} \quad (5.86a)$$

$$0 = \begin{bmatrix} x_2 + z_2 - z_1 - u \\ z_2 - x_1 \end{bmatrix} := \tilde{g} \quad (5.86b)$$

We observe that

$$\frac{\partial \tilde{g}}{\partial z} = \begin{bmatrix} -1 & 1 \\ 0 & 1 \end{bmatrix} \quad (5.87)$$

is full rank, such that (5.86) is of index 1. This concludes the procedure. We can additionally deduce from the index reduction that (5.81) is of index 2.

- Consider the semi-explicit DAE given by:

$$\dot{x} = Ax - bz \quad (5.88a)$$

$$0 = \frac{1}{2}(x^T x - 1) \quad (5.88b)$$

for some matrices A and vector b , and where $z \in \mathbb{R}$. Here the algebraic part is:

$$g(x) = \frac{1}{2}(x^T x - 1) \quad (5.89)$$

such that $\frac{\partial g(x)}{\partial z} = 0$ is rank deficient by construction. We then take the time derivative of $g(x)$:

$$\tilde{g}(x, z) = \frac{d}{dt}g(x) = x^T \dot{x} = x^T (Ax - bz) = 0 \quad (5.90)$$

We observe that

$$\frac{\partial \tilde{g}(x, z)}{\partial z} = -x^T b \quad (5.91)$$

is full rank if $x^T b \neq 0$. Note that z is here readily delivered by $\tilde{g}(x, z) = 0$. Indeed:

$$z = \frac{x^T Ax}{x^T b} \quad (5.92)$$

The index-reduced DAE then reads as:

$$\dot{x} = Ax - bz \quad (5.93a)$$

$$0 = x^T (Ax - bz) \quad (5.93b)$$

□

Similarly to Lagrange mechanics, the reduction of the index of a DAE requires one to collect a set of consistency conditions that need to be enforced in order for the resulting index-1 DAE to match the original one. Generally speaking, when performing the index reduction, one ought to **collect all the algebraic equations on which a time differentiation is performed**, and add them to the list of consistency conditions. E.g. in the example above, the first equation in (5.81b) is time-differentiated and is therefore the (only) consistency condition.

Similarly to the Lagrange context, consistency conditions can drift numerically over long simulations, and ought to be corrected if needed. Describing the introduction of a Baumgarte stabilization in generic index-reduced in a general form can be difficult, and we will leave this question out of these notes. In simple cases (such as e.g. the Lagrange context), the principle is fairly simple.

6 EXPLICIT INTEGRATION METHODS - RUNGE-KUTTA

For the remainder of these notes we will discuss the numerical treatment of differential equations. Simulating a system in e.g. the state-space form:

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, \mathbf{u}), \quad \mathbf{x}(0) = \mathbf{x}_0 \quad (6.1)$$

over a time interval $[0, T]$ consists essentially in computing a discrete sequence of state vectors

$$\mathbf{x}_0, \dots, \mathbf{x}_N \quad (6.2)$$

that approximate the true and continuous trajectories $\mathbf{x}(t)$, $t \in [0, T]$ solution of (6.1) on a given time grid t_0, \dots, t_N sufficiently accurately that they are useful for whichever goal we have to tackle, i.e. we want e.g. the approximation:

$$\|\mathbf{x}_k - \mathbf{x}(t_k)\| \leq \epsilon, \quad k = 0, \dots, N \quad (6.3)$$

to hold. For most model equations \mathbf{f} , the sequence (6.2) can only be built numerically in a computer. For the remainder of the course, we will focus on understanding some (non-exhaustive) modern methods to build a sequence (6.2) of simulated states that are “reasonably close” to the true model trajectories.

6.1 EXPLICIT EULER

In order to get started on integration methods, it is natural to start with the most intuitive, yet *least efficient* integration method, namely the explicit Euler method. Explicit Euler ought to be avoided whenever efficiency¹⁰ matters. The explicit Euler method is, however, deceptively simple to implement and therefore sometimes a good choice when efficiency ought to be sacrificed to minimize the coding effort.

The explicit Euler method uses the rule:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta t \cdot \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \quad k = 0, \dots, N \quad (6.4)$$

where $\mathbf{u}_k = \mathbf{u}(t_k)$, starting from the given initial conditions \mathbf{x}_0 . The Euler step is illustrated in Figure 6.1, and is essentially linearly extrapolating from the model dynamics $\mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$ in order to build the next discrete state \mathbf{x}_{k+1} from the current one \mathbf{x}_k .

¹⁰i.e. the amount of computations needed to reach a given accuracy

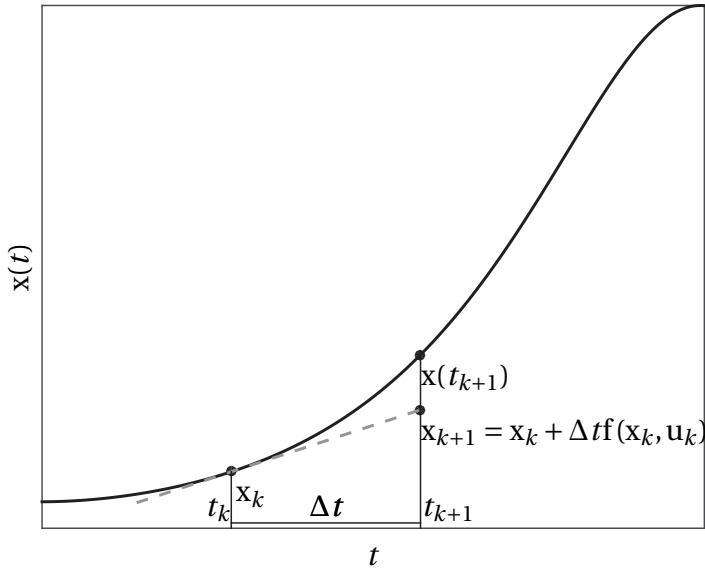


Figure 6.1: Illustration of the principle underlying the explicit Euler scheme (6.4)

It ought to be intuitively clear here that the larger the time step Δt is chosen, the larger the discrepancy between the true model trajectories $x(t)$ and the simulated ones x_0, \dots, x_N are. Indeed, Euler is essentially ignoring the “curvature” of the model trajectories between the discrete states and replacing it via a line. The longer the step, the “more curvature” is ignored.

6.1.1 ACCURACY OF THE EXPLICIT EULER METHOD

It is useful to formally analyse the discrepancy between the simulated states and the true trajectories, and relate them to the choice and parameters (e.g. Δt) of the integration method. In order to do that, it is very useful to compute the **one-step error**, i.e. assuming that $x_k = x(t_k)$ (i.e. the simulation is exact at time t_k), then what is the error at time t_{k+1} , i.e. what is $\|x_{k+1} - x(t_{k+1})\|$? This can be done via Taylor arguments. Indeed, we observe that:

$$x(t_{k+1}) = \underbrace{x(t_k) + \Delta t \cdot \dot{x}(t_k)}_{=x_{k+1}} + \frac{1}{2} \Delta t^2 \ddot{x}(\xi) \quad (6.5)$$

for some $\xi \in [t_k, t_{k+1}]$ (this is the Taylor expansion with an explicit remainder). It follows that:

$$\|x_{k+1} - x(t_{k+1})\| = \frac{\Delta t^2}{2} \|\ddot{x}(\xi)\| \leq \max_{\xi \in [t_k, t_{k+1}]} \frac{\Delta t^2}{2} \|\ddot{x}(\xi)\| \quad (6.6)$$

Two remarks can be drawn from this expression:

- The one-step error is of the order Δt^2

- The one-step error is worse when \ddot{x} is large, i.e. if the model trajectory is more “curved”.

Using this result we can then discuss the **global error** of the explicit Euler integration method, i.e. what is the final integration error $\|x_N - x(T)\|$? We can answer that question almost directly from (6.6). Indeed, we make the following observations:

- An integration up to time T with a time step Δt requires $N = \frac{T}{\Delta t}$ Euler steps
- Each step yields an error of order Δt^2
- After N step, the error will then be of the order $N\Delta t^2 = \frac{T}{\Delta t}\Delta t^2 = T\Delta t$
- After N step, the simulation error $\|x_N - x(T)\|$ will be of the order Δt . We will use the notation

$$\|x_N - x(T)\| = \mathcal{O}(\Delta t) \quad (6.7)$$

throughout the notes, which formally means:

$$\|x_N - x(T)\| \leq c\Delta t \quad (6.8)$$

holds for some constant $c > 0$ and for Δt small enough.

We illustrate next these principle on the integration of the following dynamics:

$$\dot{x} = f(x, u) = \begin{bmatrix} \sigma(x_2 - x_1) \\ x_1(\rho - x_3) - x_2 \\ x_1x_2 - \beta x_3 \end{bmatrix} \quad (6.9)$$

for $\sigma = 10$, $\beta = 3$, $\rho = 28$ and starting at the initial conditions $x_0 = [1 \ 1 \ 1]^\top$. Figures 6.2-6.4 illustrate the model trajectories obtained from the explicit Euler scheme (6.4) for different step sizes Δt , and figure 6.5 report the global error, i.e. the error observed at the end of the simulation.

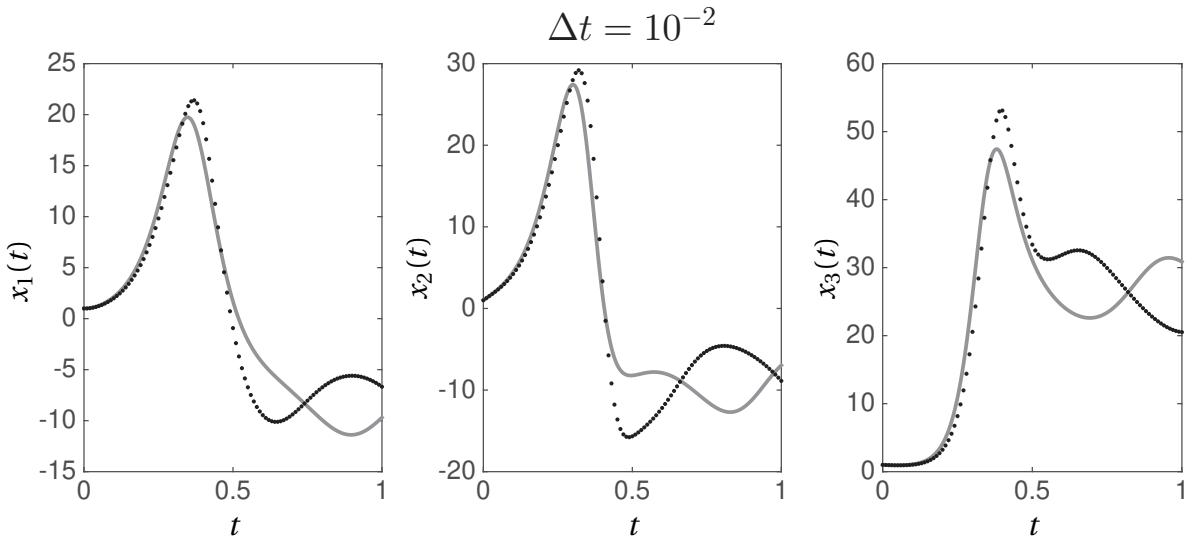


Figure 6.2: Numerical integration of the equation (6.9). In dark the trajectories from the explicit Euler scheme (6.4) using $\Delta t = 10^{-2}$, and in grey the trajectories obtained using a very high accuracy integration method.

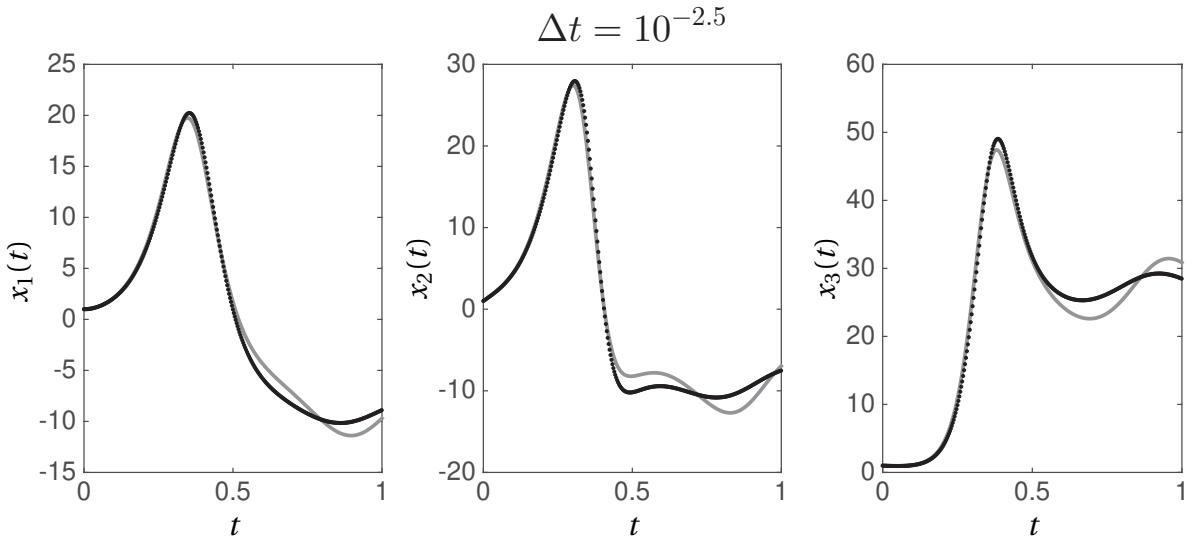


Figure 6.3: Numerical integration of the equation (6.9). In dark the trajectories from the explicit Euler scheme (6.4) using $\Delta t = 10^{-2.5}$, and in grey the trajectories obtained using a very high accuracy integration method.

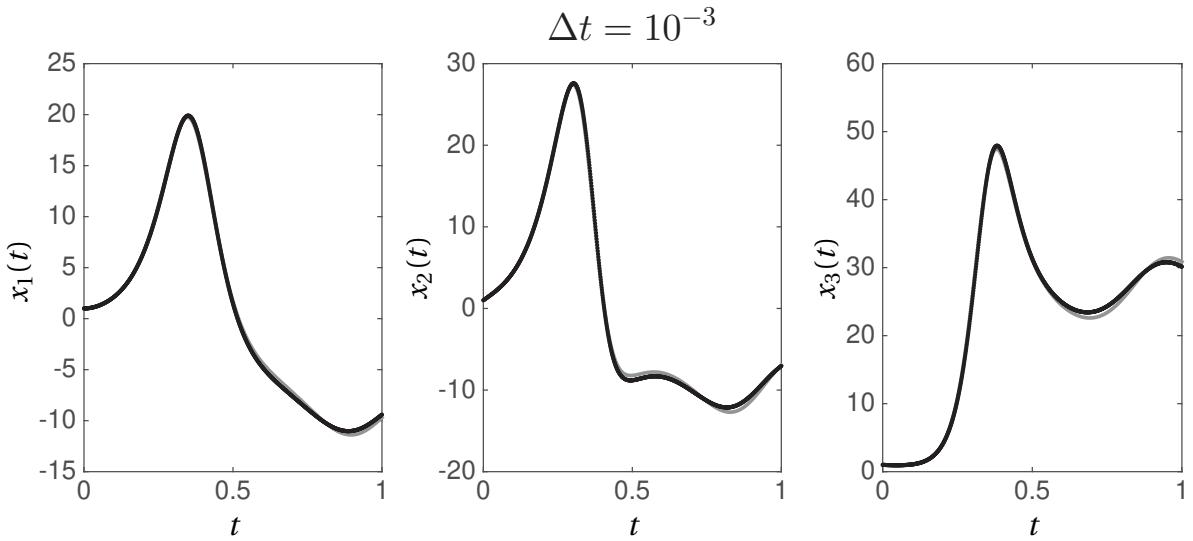


Figure 6.4: Numerical integration of the equation (6.9). In dark the trajectories from the explicit Euler scheme (6.4) using $\Delta t = 10^{-3}$, and in grey the trajectories obtained using a very high accuracy integration method.

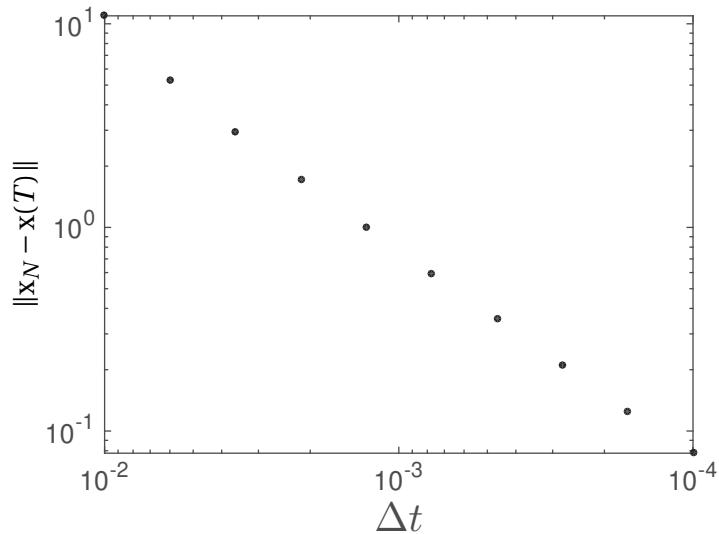


Figure 6.5: Global error vs. Δt using the explicit Euler scheme (6.4).

6.1.2 STABILITY OF THE EXPLICIT EULER METHOD

Beyond the accuracy of integration methods, another crucial aspect to discuss is their **stability**. In order to introduce this concept, let us consider the trivial stable scalar linear system:

$$\dot{x} = -\lambda x, \quad x(0) = x_0 \quad (6.10)$$

for $\lambda > 0$ and arbitrary initial conditions x_0 . The dynamics (6.10) are linear and have therefore the explicit solution:

$$x(t) = x_0 e^{-\lambda t} \quad (6.11)$$

Let us nonetheless consider deploying the explicit Euler method on (6.10). The Euler steps then read as:

$$x_{k+1} = x_k - \lambda \Delta t x_k = (1 - \lambda \Delta t) x_k \quad (6.12)$$

We observe that (6.12) is a linear discrete dynamic system, and that it becomes **unstable** for $|1 - \lambda \Delta t| > 1$. Since $\lambda, \Delta t > 0$, this happens if $1 - \lambda \Delta t < -1$, i.e. if

$$\Delta t > \frac{2}{\lambda} \quad (6.13)$$

In other words, the **explicit Euler method delivers an unstable simulation** for the dynamics (6.10) **when the time step Δt is too large** compared to the pole of the dynamics λ . We need to stress here that the dynamics (6.10) are stable, only their numerical simulation is possibly unstable. We illustrate these observations in Fig. 6.6 below.

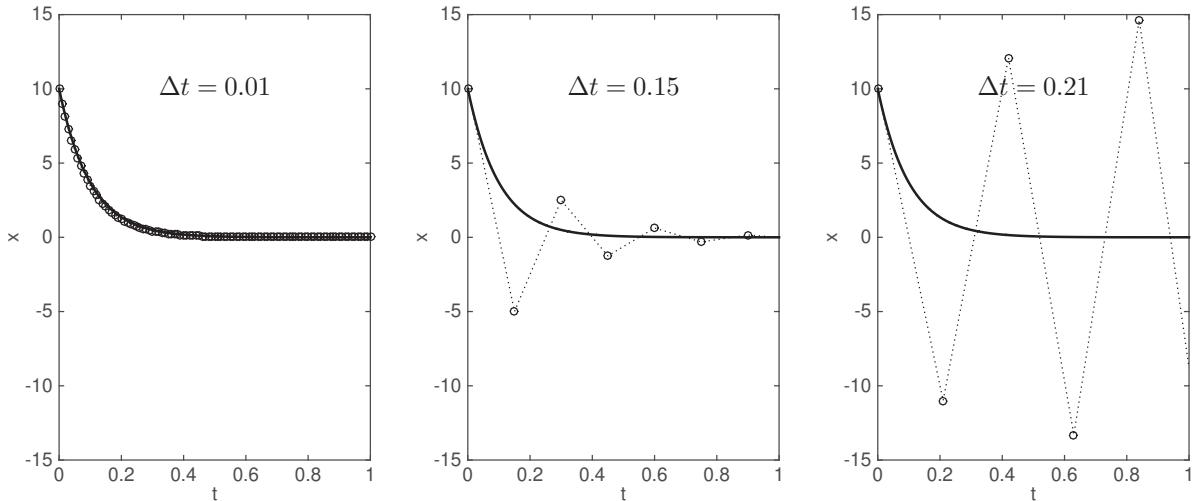


Figure 6.6: Illustration of the stability problem of the Euler method for $\lambda = 10$. Only a small enough step-size $\Delta t \leq 0.2$ allows the method to be stable (left and middle graph), while $\Delta t > 0.2$ yields an unstable simulation (right graph).

We will need to discuss the stability of each integration method we will study. This stability issue is in fact studied in a slightly more general way for various integration methods, namely the test system:

$$\dot{x} = \lambda x, \quad x(0) = x_0 \quad (6.14)$$

is used with $\lambda \in \mathbb{C}$ and the integration method stability is described in terms of **region of stability** in the complex plane, for which a given $\lambda\Delta t$ yields a stable simulation. In the specific case of the explicit Euler method the region of stability is given by a circle of radius one centered at -1 , i.e.

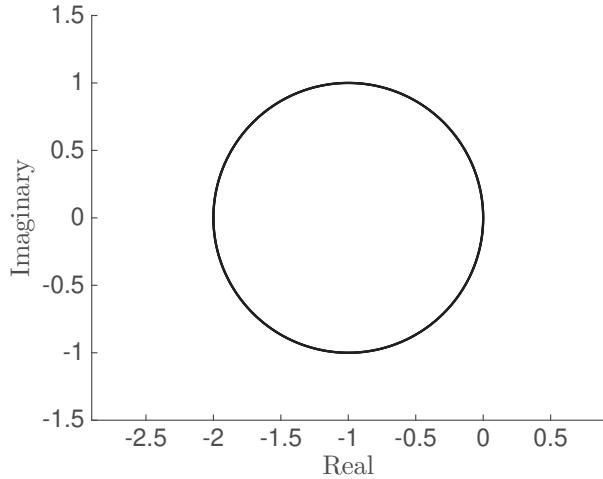


Figure 6.7: Stability region of the explicit Euler scheme (6.4) in the complex plane. All combinations of $\lambda, \Delta t$ such that $\lambda\Delta t$ lies within the circle result in a stable (though not necessarily accurate) integration.

Obviously, the larger the region of stability is the more “stable” the method is considered.

Summary

- The explicit Euler method is a method of order 1, i.e. the global error is $\|x_N - x(T)\| = \mathcal{O}(\Delta t)$ this is the lowest order that one can accept from an integration method, as one expects from a method that the error gets smaller with the step size Δt .
- The explicit Euler method can become unstable for a too large step size. Its **region of stability** for the test system (6.14) is fairly small and given in Fig. 6.7

6.2 EXPLICIT RUNGE-KUTTA 2 METHODS

We will now consider an improvement of the explicit Euler method, while remaining in the family of explicit methods (the meaning of this label will become clear later). The most widely known family of explicit integration methods (beyond Euler) is the family of Runge-Kutta methods. Let us start with studying the Runge-Kutta 2 method (RK2) for which a simple interpretation of its construction can be developed.

The exact trajectories of the model dynamics (6.1) obey the following integral relationship:

$$x(t_{k+1}) = x(t_k) + \int_{t_k}^{t_{k+1}} f(x(t), u(t)) dt \quad (6.15)$$

The key idea behind many integration methods is to try to provide good evaluations of the integral term in (6.15). In fact, one can interpret the explicit Euler method as a crude way of doing that. Indeed, explicit Euler approximates:

$$\int_{t_k}^{t_{k+1}} f(x(t), u(t)) dt \approx (t_{k+1} - t_k) f(x(t_k), u(t_k)) \approx \Delta t f(x_k, u_k) \quad (6.16)$$

essentially using a one-step rectangular quadrature on the integral. This approximation illustrated in Figure 6.8.

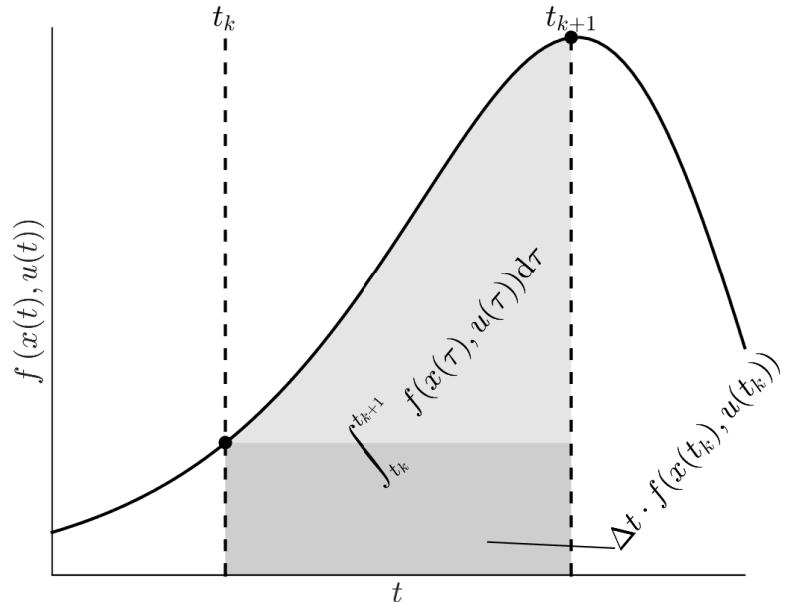


Figure 6.8: Illustration of the approximation (6.16). The light grey area is $\int_{t_k}^{t_{k+1}} f(x(t), u(t)) dt$ (including the dark grey area), which the explicit Euler scheme is approximating with $\Delta t f(x_k, u_k)$, i.e. the dark grey area.

A natural question then is how can we improve the approximation (6.16). To that end, one can e,g, use a mid-point rule, based on the approximation

$$\int_{t_k}^{t_{k+1}} f(x(t), u(t)) dt \approx \Delta t f\left(x\left(t_k + \frac{\Delta t}{2}\right), u\left(t_k + \frac{\Delta t}{2}\right)\right) \quad (6.17)$$

which tends to be better than (6.16) whenever the model trajectories have a “low curvature”. This observation is illustrated in Figure 6.9.

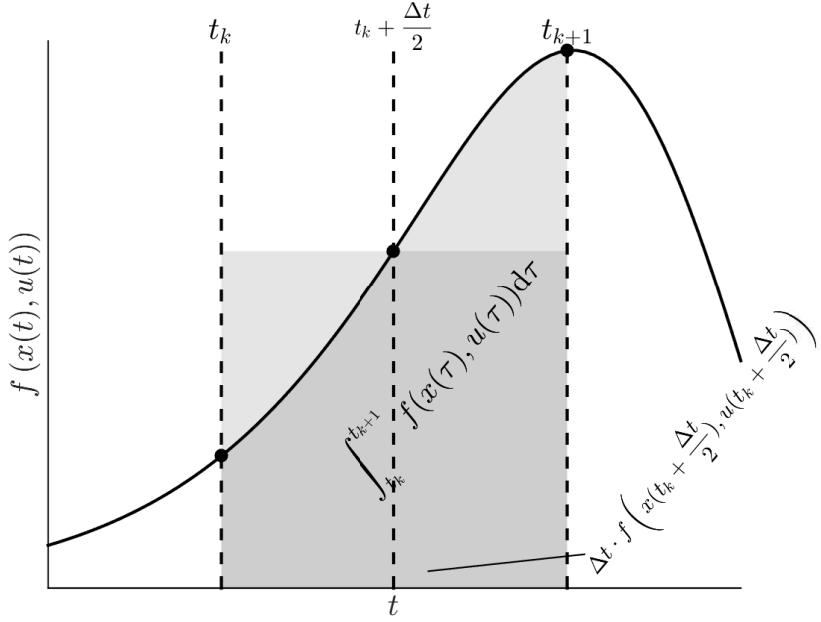


Figure 6.9: Illustration of the approximation (6.17). The light grey area is $\int_{t_k}^{t_{k+1}} f(x(t), u(t)) dt$ (including the dark grey area), which the explicit Euler scheme is approximating with $\Delta t f\left(x\left(t_k + \frac{\Delta t}{2}\right), u\left(t_k + \frac{\Delta t}{2}\right)\right)$, i.e. the rectangular grey area.

Unfortunately, $x\left(t_k + \frac{\Delta t}{2}\right)$ is unknown and needs to be itself replaced by an approximation built upon x_k . Here we rely again on an explicit Euler “half” step, i.e. we use:

$$x\left(t_k + \frac{\Delta t}{2}\right) \approx x_k + \frac{\Delta t}{2} f(x_k, u_k) \quad (6.18)$$

Combining (6.18) and (6.17), we obtain the “mid-point” RK2 integration scheme:

$$K_1 = f(x_k, u(t_k)) \quad (6.19a)$$

$$K_2 = f\left(x_k + \frac{\Delta t}{2} K_1, u\left(t_k + \frac{\Delta t}{2}\right)\right) \quad (6.19b)$$

$$x_{k+1} = x_k + \Delta t K_2 \quad (6.19c)$$

Despite the somewhat convoluted construction, we can fairly easily show that this RK2 scheme is of order 2 (hence the name). In order to compute the order, we will use a Taylor argument again. In order to reduce the complexity of the notation, let us do the calculation assuming that the model dynamics do not have an input u . Assuming again that $x_k = x(t_k)$ (i.e. the integration is exact at time t_k), we observe that:

$$x(t_{k+1}) = x(t_k) + \Delta t f(x(t_k)) + \frac{1}{2} \Delta t^2 \dot{f}(x(t_k)) + \mathcal{O}(\Delta t^3) \quad (6.20)$$

where $\dot{f} = \frac{\partial f}{\partial x} f$. We additionally observe that (6.19) in effect does:

$$x_{k+1} = x_k + \Delta t f\left(x_k + \frac{\Delta t}{2} f(x_k)\right) \quad (6.21)$$

and we finally observe that

$$f\left(x_k + \frac{\Delta t}{2} f(x(t_k))\right) = f(x(t_k)) + \frac{\Delta t}{2} \frac{\partial f}{\partial x}\Big|_{x_k} f(x(t_k)) + \mathcal{O}(\Delta t^2) \quad (6.22)$$

Hence

$$x_{k+1} = x_k + \Delta t f(x(t_k)) + \frac{\Delta t^2}{2} \frac{\partial f}{\partial x}\Big|_{x_k} f(x(t_k)) + \mathcal{O}(\Delta t^3) \quad (6.23)$$

A comparison of (6.20) and (6.23) finally yields:

$$\|x_{k+1} - x(t_{k+1})\| = \mathcal{O}(\Delta t^3) \quad (6.24)$$

We can then conclude that the **one-step error of the RK2 scheme is of order 3**, and using a similar reasoning as for the explicit Euler scheme, the **global error of the RK2 scheme is of order 2**, i.e.

$$\|x_N - x(T)\| = \mathcal{O}(\Delta t^2) \quad (6.25)$$

(see Figure 6.10 for an illustration). We will discuss the stability of RK2 later, together with the other RK methods.

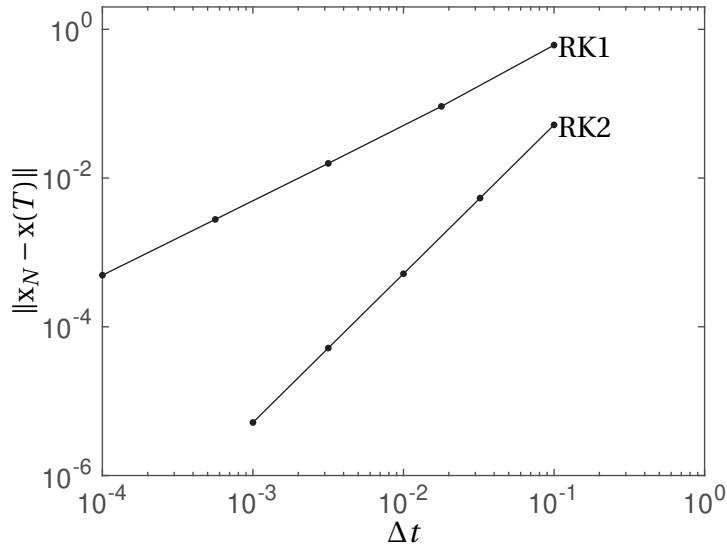


Figure 6.10: Global error of the RK2 method (6.19), and comparison with the explicit Euler scheme (6.4) (RK1), as obtained from numerical experiments. One can see the difference between having a first or second order method. The stronger slope of the RK2 method (in the log-log plot) indicates a higher power in the relationship Δt to $\|x_N - x(T)\|$.

6.3 GENERAL RK METHODS

The principles of the RK2 method can be generalized. We introduce here the general formula for RK methods, which is a generalization of (6.19).

$$K_1 = f \left(x_k + \Delta t \sum_{j=1}^s a_{1j} K_j, u(t_k + c_1 \Delta t) \right) \quad (6.26a)$$

⋮

$$K_i = f \left(x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, u(t_k + c_i \Delta t) \right) \quad (6.26b)$$

⋮

$$K_s = f \left(x_k + \Delta t \sum_{j=1}^s a_{sj} K_j, u(t_k + c_s \Delta t) \right) \quad (6.26c)$$

$$x_{k+1} = x_k + \Delta t \sum_{i=1}^s b_i K_i \quad (6.26d)$$

where the coefficients a_{ij} , b_i , c_j define the RK method, and s is the number of stages of the method. We can easily observe that both the explicit Euler method and the RK2 methods introduced above can be written in this form. Indeed, it can be verified that:

- The explicit Euler method has $s = 1$ with

$$a_{11} = 0, \quad b_1 = 1, \quad c_1 = 0 \quad (6.27)$$

- The “mid-point” RK2 method has $s = 2$ with

$$a_{11} = 0, \quad a_{12} = 0, \quad c_1 = 0 \quad (6.28a)$$

$$a_{21} = \frac{1}{2}, \quad a_{22} = 0, \quad c_2 = \frac{1}{2} \quad (6.28b)$$

$$b_1 = 0, \quad b_2 = 1 \quad (6.28c)$$

6.3.1 BUTCHER TABLEAU

The Butcher tableau is a convenient way of summarizing the coefficient of RK methods in a clean and compact way. The Butcher tableau reads as:

c_1	a_{11}	\dots	a_{1s}	
\vdots	\vdots		\vdots	
c_s	a_{s1}	\dots	a_{ss}	
	b_1	\dots	b_s	

E.g. explicit Euler has the Butcher tableau

0	0
	1

and is therefore a RK1 method (because of order 1). The “mid-point” RK2 method has the Butcher tableau

0	0 0
$\frac{1}{2}$	$\frac{1}{2} 0$
	0 1

Other RK2 methods exist, with Butcher tableaus of the same size but with different coefficients, e.g.

Ralston's RK2

0	0 0
$\frac{2}{3}$	$\frac{2}{3} 0$
	$\frac{1}{4} \frac{3}{4}$

Henn's RK2

0	0 0
1	$1 0$
	$\frac{1}{2} \frac{1}{2}$

RK methods resulting from (6.26) ought to be divided in two categories **explicit** and **implicit**. In order to make these terms clear, let us consider all RK methods having two stages ($s = 2$). In this case (6.26) reads as:

$$K_1 = f(x_k + \Delta t(a_{11}K_1 + a_{12}K_2), u(t_k + c_1\Delta t)) \quad (6.29a)$$

$$K_2 = f(x_k + \Delta t(a_{21}K_1 + a_{22}K_2), u(t_k + c_2\Delta t)) \quad (6.29b)$$

$$x_{k+1} = x_k + \Delta t(b_1K_1 + b_2K_2) \quad (6.29c)$$

We observe that for $a_{11}, a_{12}, a_{22} = 0$ in (6.29) (such as in all the RK2 we have looked at above), K_1 can be computed explicitly from $x_k, u(\cdot)$, and then K_2 can be computed explicitly from $x_k, u(\cdot), K_1$. However, for any of the a_{11}, a_{12}, a_{22} non-zero, this cannot be done as equations (6.29a)-(6.29b) become *implicit* (the unknown $K_{1,2}$ appear on both sides of the equalities and are linked through the function f). We then say that the RK scheme is implicit. One can trivially see from the Butcher tableau of an RK method whether the method is explicit or implicit. We state this next.

A Butcher tableau defines an explicit integrator if and only if only the diagonal and upper-diagonal elements are zero. I.e. if and only if $a_{ij} = 0$ for any $j \geq i$. Otherwise it defines an implicit method.

The distinction is important in practice. Indeed, while explicit RK methods require simply

computing the K_i sequentially (K_1 then K_2 , etc.), implicit RK methods require solving the equations together and numerically, typically using a Newton method. The latter is computationally more expensive because it requires solving linear systems a number of times for each integration steps. However, as we will see later, implicit RK schemes have powerful properties that make them attractive, even when computational time is important.

6.3.2 THE RK4 METHOD

It is worth discussing one of the most commonly used RK method, the explicit RK4 (order 4), which requires $s = 4$ stages. It has the Butcher tableau¹¹

0	0	0	0	0
$\frac{1}{2}$	$\frac{1}{2}$	0	0	0
$\frac{1}{2}$	0	$\frac{1}{2}$	0	0
1	0	0	1	0
	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{6}$

It is a popular approach because it has a very good trade-off between computational complexity (CPU time) and accuracy.

- The RK4 method has $s = 4$ stages
- The integration order of this scheme is 4, i.e.

$$\|\mathbf{x}_N - \mathbf{x}(T)\| \leq c \|\Delta t^4\| \quad (6.30)$$

holds for some $c > 0$ and for Δt small enough. We should probably stress here the implication of having an order 4. The consequence of the order 4 is that if one divides the time step Δt by 2, then the amount of computation required to run a simulation on a time interval $[0, T]$ is doubled, but the numerical error of the simulation is divided by $2^4 = 16$. This Δt^4 effect of the order 4 makes it "cheap" to gain accuracy (cheap in terms of computational time).

- The method requires a few lines of computer code to implement

For the sake of completeness, let us write the pseudo-code corresponding to an RK4 method.

¹¹Note that other Butcher tableaus can generate order 4 RK methods, but this one is arguably the most commonly used.

Algorithm: RK4 method for explicit ODEs

Input: Initial conditions x_0 , input profile $u(\cdot)$, step size Δt

for $k = 0, \dots, N - 1$ **do**

 Compute:

$$K_1 = f(x_k, u(t_k)) \quad (6.31a)$$

$$K_2 = f\left(x_k + \frac{\Delta t}{2} K_1, u\left(t_k + \frac{\Delta t}{2}\right)\right) \quad (6.31b)$$

$$K_3 = f\left(x_k + \frac{\Delta t}{2} K_2, u\left(t_k + \frac{\Delta t}{2}\right)\right) \quad (6.31c)$$

$$K_4 = f(x_k + \Delta t K_3, u(t_k + \Delta t)) \quad (6.31d)$$

Assemble the RK step

$$x_{k+1} = x_k + \Delta t \left(\frac{1}{6} K_1 + \frac{1}{3} K_2 + \frac{1}{3} K_3 + \frac{1}{6} K_4 \right) \quad (6.32)$$

return $x_{1,\dots,N}$

6.3.3 STAGES, ORDER & EFFICIENCY OF EXPLICIT RK METHODS

Let us discuss more the integration order of explicit RK methods. Before diving in, we should discuss more the implication of the number of stages s that the different RK methods have. It can be observed that the number of stages of an explicit RK method defines the computational complexity of evaluating one step $x_k \rightarrow x_{k+1}$ in the integrator. Indeed, each integration step (i.e. the execution of (6.26)) requires s evaluations of the model equations f and each of these evaluation translates directly into a (typically) fixed CPU time. We can deduce that the computational cost of one step of an explicit integrator is roughly proportional to s . In that sense, integrators with a low number of stages appear preferable.

On the other hand, integrators with a larger number of stages s also achieve higher orders (let's label them o), and therefore achieve a higher accuracy for a given Δt . Indeed, remember that the simulation error is bounded by:

$$\|x_N - x(T)\| \leq c \Delta t^o \quad (6.33)$$

for some $c > 0$ and Δt small. Note that we should always consider Δt small, i.e. Δt^o becomes exponentially smaller as o increases. Hence by increasing the number of stages s we pay *proportionally* more in computations but we gain *exponentially* in accuracy. The gain tends to out-weight the loss because of the "power effect" of the order.

The RK methods we have discussed so far suggest that the number of stages s of the method equals its order o . Indeed, we have seen that:

- Explicit Euler has an order $o = 1$ and is an RK method with $s = 1$ stage. If Δt is divided by 2 then the error is divided by 2.

- The RK2 methods we have seen have an order $o = 2$ and have $s = 2$ stages. If Δt is divided by 2 then the error is divided by 4.
- The RK4 method we have seen has order $o = 4$ and has $s = 4$ stages. If Δt is divided by 2 then the error is divided by 16.

From these observations, it may appear that the higher s the better. Unfortunately, this pattern $o = s$ breaks beyond $o = 4$. Let us stress this fact in the following table (these numbers are not straightforward to establish, and we will leave that question out of these notes).

Order	Stages required	Δt divided by 2 → error divided by...
RK1	1	2
RK2	2	4
RK3	3	8
RK4	4	16
RK5	6	32
RK6	7	64
RK7	9	128
RK8	11	256

We observe from this table that for a low number of stages ($s \leq 4$), we have $o = s$, hence adding a stage adds an order (see Fig. 6.11 for an illustration). However, for a higher number of stages, the increase in the order of the method “stalls”, and stops increasing as quickly as the number of stages.

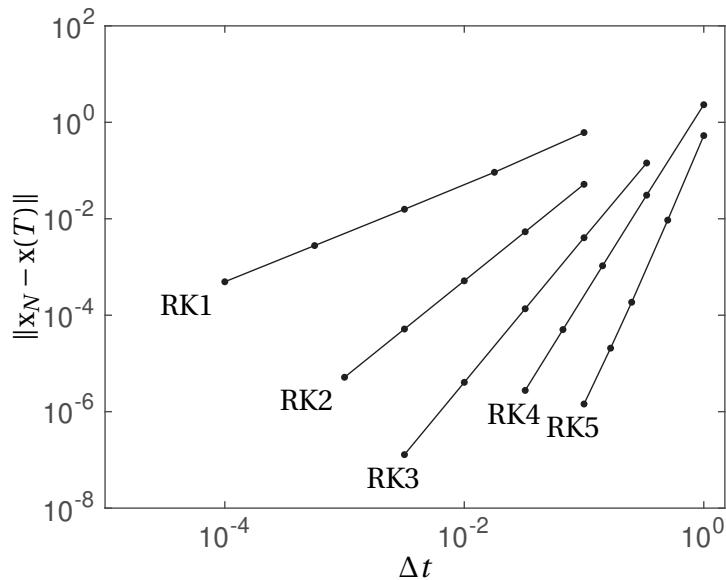


Figure 6.11: Global error of the different RK methods, as obtained from numerical experiments.

This effect has some important consequences in practice on the choice of an integration method. Indeed, the crux of any integration method is to achieve a given accuracy with a minimum computational budget. There are two ways of increasing the integrator accuracy: increasing the order o , i.e. increasing the number of stage s , or decreasing Δt . Having these two options makes the choice of integrator (choice of order and choice of time step) non-trivial.

However, using the table above, we can in fact somewhat formalize this choice. Let us assume that we want the global error to be limited to some given number Tol

- Then for an integrator of order o , we need:

$$\|\mathbf{x}_N - \mathbf{x}(T)\| \leq c\Delta t^o \leq Tol \quad (6.34)$$

such that the step size Δt is limited to:

$$\Delta t \leq \left(\frac{Tol}{c}\right)^{\frac{1}{o}} \quad (6.35)$$

- In order to carry out a simulation on the time interval $[0, T]$, the number of integrator step required is:

$$N = \frac{T}{\Delta t} \quad (6.36)$$

and for a number of stages s , the number of evaluation n of the system dynamics f is:

$$n = N \cdot s = \frac{sT}{\Delta t} \geq sT \left(\frac{Tol}{c}\right)^{-\frac{1}{o}} \quad (6.37)$$

We deduce from this simple reasoning that the computational cost per unit of simulation time is at least:

$$\frac{n}{T} \geq s \left(\frac{Tol}{c}\right)^{-\frac{1}{o}} \quad (6.38)$$

where o and s are related via the table above. It is then interesting to chart this relationship, we do that in the Figures below. Here it becomes clear that a very low or very high order is not optimum, and that the optimum is in the order 3-6.

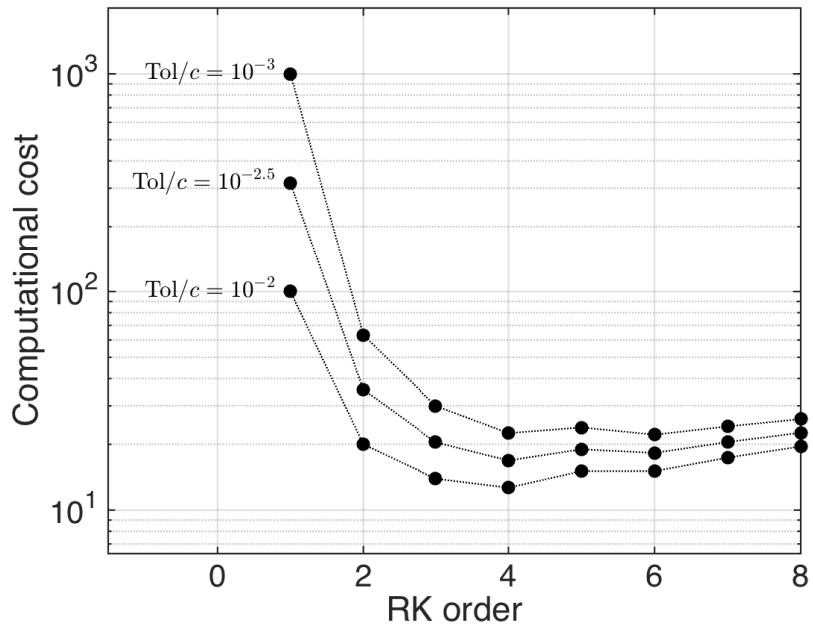


Figure 6.12: Illustration of formula (6.38) for different $\frac{\text{Tol}}{c}$. The horizontal axis represents the order o of the method, and the vertical the computational cost per simulation time T , i.e. $\frac{n}{T}$. The optimum is achieved in the middle range, and not for low nor high order methods.

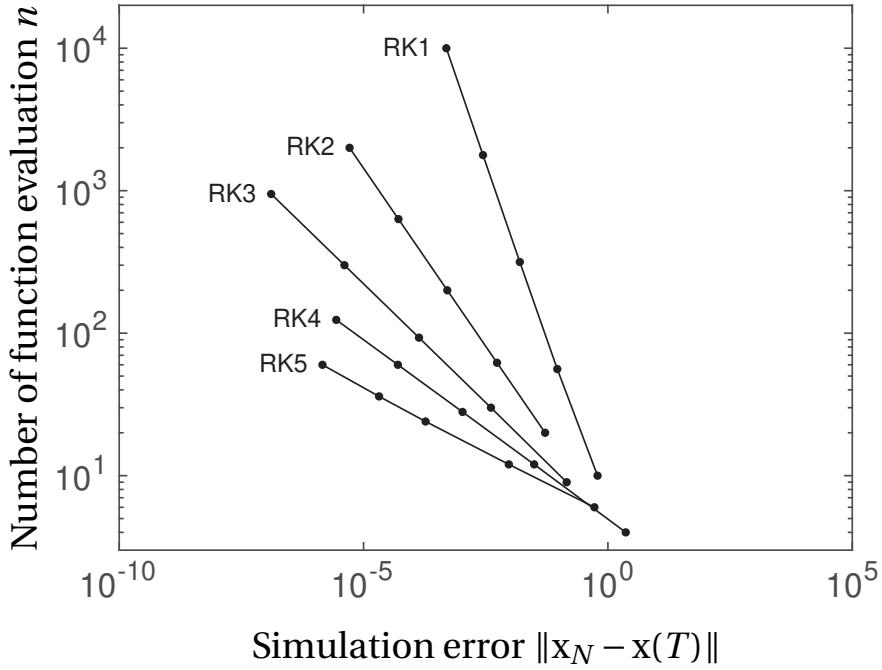


Figure 6.13: Illustration of the number of function evaluations required to reach different level of simulation accuracy for different explicit RK methods for the integration of a linear test example. Here the RK5 scheme beats the other schemes for high accuracies, and RK4 beats the other schemes for lower accuracies. This picture is model-dependent, and should just be understood as an example.

Before closing this section, a last point ought to be stressed regarding the order o achieved by different RK methods. Indeed, we have seen the relationship between the order o of different RK methods and the number of stages s required to achieve that order. We need to stress here that the order o of an RK method does not follow *simply* from using the adequate amount of stages. E.g. the RK2 schemes ($o = 2$) we have seen require $s = 2$ stages, but they also require the coefficients a, b, c to be chosen adequately. E.g. in the case of a $s = 2$ stages method, an order $o = 2$ is achieved if the coefficients satisfy:

$$b_1 + b_2 = 1, \quad b_2 c_2 = \frac{1}{2}, \quad a_{21} = c_2 \quad (6.39)$$

Similarly, for the other RK methods, the coefficients cannot take *arbitrary* value if one wants the method to achieve its highest possible order.

6.4 STABILITY OF EXPLICIT RK SCHEMES

Let us go back to the question of integration stability we investigated in section 6.1.2 in the more general context of RK methods of arbitrary order. The general formula defining the stability region of the test system

$$\dot{x} = \lambda x, \quad x(0) = x_0 \quad (6.40)$$

for an RK method of order o is:

$$S = \left\{ \lambda \Delta t \quad \text{s.t.} \quad \left| \sum_{k=0}^o \frac{(\lambda \Delta t)^k}{k!} \right| \leq 1 \right\} \quad (6.41)$$

The regions of stability for different orders are depicted in Figure 6.14. One can observe that the regions increase with the order, but their size is fairly limited.

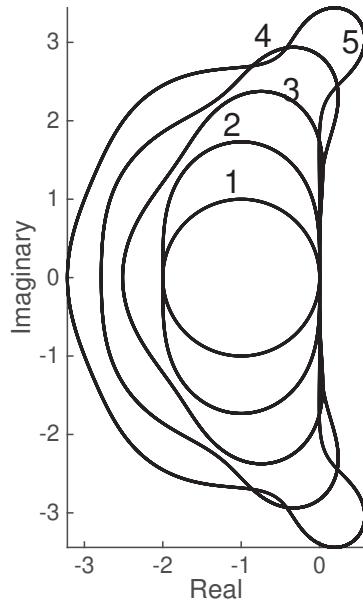


Figure 6.14: Illustration of the region os stability for the explicit RK methods 1 to 5. One can observe that the stability region grows with the order, but it nonetheless covers a limited range of admissible $\lambda \Delta t$.

6.4.1 STIFF SYSTEMS

The numerical stability considerations discussed above have important practical consequences. Indeed, while the dynamics(6.40) are clearly of no interest, they operate as a trivial test system to discuss how integration methods react to **fast dynamics**, i.e. more specifically, time constants in the dynamics that are in the same ballpark as the integration step Δt . The question of whether an integration method is capable of handling such fast dynamics is crucial in the context of stiff systems, which appear very often in models for engineering application such as mechanical and electrical systems. It is important to “detect” the stiffness of a model when trying to simulate it.

We have in fact already encountered stiff systems in the DAE section. Indeed, (5.33) essentially describes a model having very fast dynamics (for ϵ small) on the states z , and slower dynamics on the states x . DAEs allow one to approximate these fast states via their “decayed” values (i.e. their trajectories after their fast, stable dynamics have decayed). However, one may want to process these dynamics without using DAE formulations, in which

case the fast dynamics have to be dealt with. As an example of a stiff system, consider the model equations

$$\dot{\mathbf{x}} = \begin{bmatrix} 0 & 5 & 1 & 0 \\ -5 & 0 & 0 & \frac{1}{2} \\ 0 & 0 & -10^{-1}\epsilon & \epsilon \\ 0 & \epsilon & -\epsilon & -10^{-1}\epsilon \end{bmatrix} \mathbf{x} \quad (6.42)$$

for $\epsilon = 5 \cdot 10^{-4}$. The corresponding trajectories are illustrated in Fig. 6.15 and 6.16. One can observe that states $x_{1,2}$ evolve slowly (although they are influenced by the fast states and have some small oscillations in the beginning), while the states $x_{3,4}$ oscillate fast and decay to slow trajectories. This kind of system can e.g. arise in mechanical systems or electrical circuits when some part have eigenfrequencies that are much higher than other parts. They are especially expensive to treat in numerical simulations because the fast dynamics require a small Δt in order for the numerical simulation to be stable, while long simulations (T large) are required in order to see the evolution of the “slow” states.

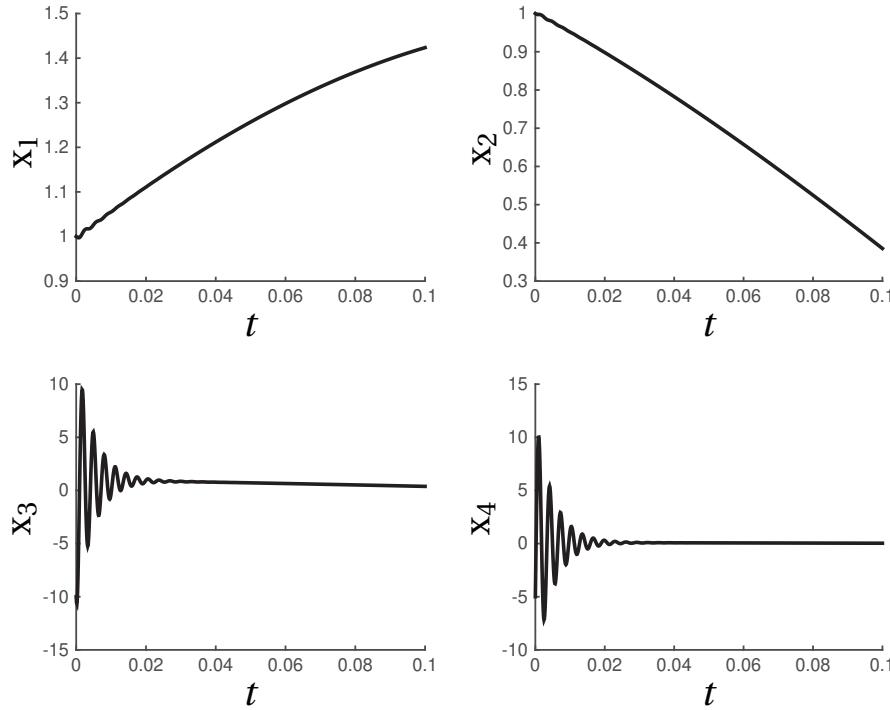


Figure 6.15: Illustration of the trajectories of the stiff system (6.42).

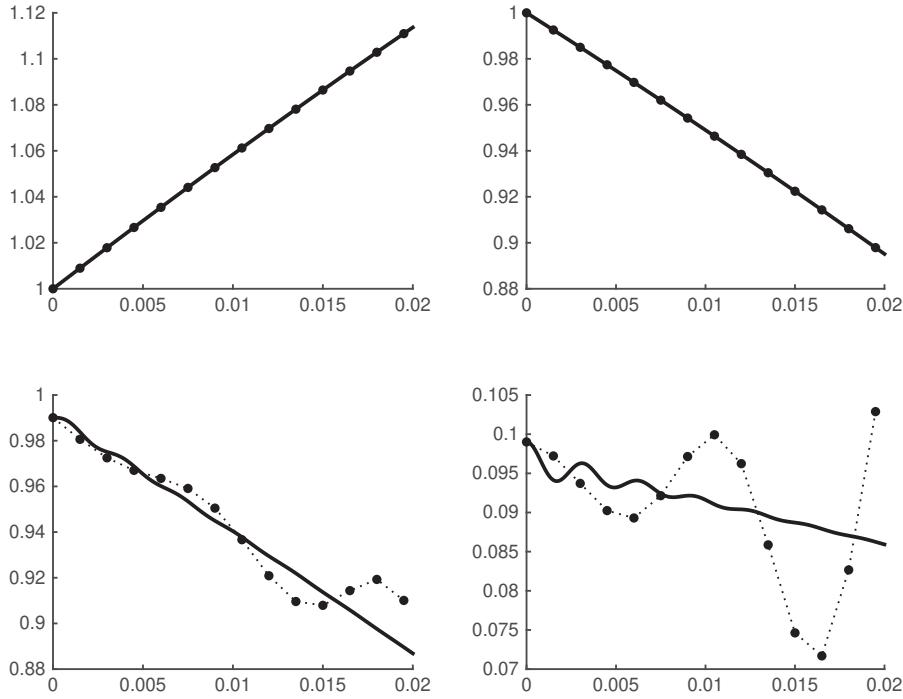


Figure 6.16: Illustration of the trajectories of the stiff system (6.42) vs. a simulation using RK4 with $\Delta t = 1.5 \cdot 10^{-3}$ (dotted curve). The initial conditions for the “fast” states are chosen such that $\dot{z} = 0$ at $t = 0$. The small step size is close to the limit for which the simulation becomes stable.

We will see later on that stiff systems are in fact one of the motivations for using implicit integration methods.

6.5 ERROR CONTROL & ADAPTIVE INTEGRATORS

We have seen above that the integrator accuracy and stability depend on the step size Δt . It can be tricky, however, to select the correct step size a priori. Moreover, as suggested in the order computation of the RK1 and 2 schemes, the accuracy depends on how “curved” the model trajectories are, and this “curvature” can vary significantly throughout the trajectories of the system. A careful approach would then be to simply use an overly small step size to ensure that the required accuracy is achieved. This conservative approach would, however, be overly expensive, as it would use very small step sizes Δt even when they are not needed. A classic example of an ODE triggering such difficulties is the Van Der Pol oscillator having the equations

$$\dot{x}_1 = (1 - x_2^2)x_1 - x_2, \quad (6.43a)$$

$$\dot{x}_2 = x_1, \quad (6.43b)$$

and e.g. the trajectories displayed in Fig. 6.17.

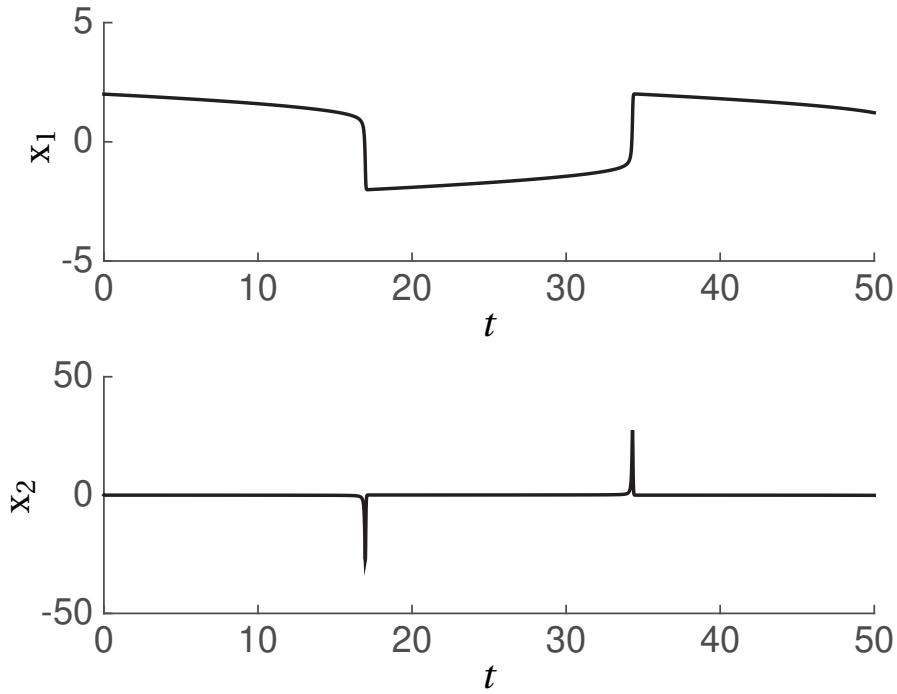


Figure 6.17: Illustration of the Van Der Pol trajectories. One can observe that the trajectories are fairly “benign” most of the time, but regularly go through very sharp changes. The Van Der Pol oscillator is a challenging ODE for numerical integration as very fine time steps Δt are needed to “survive” the sharp changes, while fairly large time steps can be used on the rest of the trajectories.

It then makes sense to use varying step sizes Δt throughout the trajectories. Adaptive integration schemes are a practical approach to do precisely that. The key idea here is to adjust the step size Δt during the simulations in order to meet the required accuracy. In practice, adaptive RK schemes perform an error control at every RK step, and when the error is larger than acceptable, the step size is reduced and a shorter step is attempted until the tolerance is met. Adaptive integrator require a baseline to assess the error at each step. A common practice is to compute the integration step using two different Butcher tableaus, and compare their outcome. If their discrepancy is above a certain level, then the step size is reduced.

Let us e.g. consider the RK45 adaptive integration scheme, which is implemented in Matlab in the function `ode45.m`. At each RK step, the procedure:

- Generates two steps from x_k , let us call them x_{k+1} and \hat{x}_{k+1} , using two different Butcher tableaus.
- Compares the two steps, i.e. computes $e = \|x_{k+1} - \hat{x}_{k+1}\|$
- Reduces Δt if the error e is above some tolerance, and computes the steps again until the tolerance is met

- Increases Δt a bit if the error e is significantly below the tolerance

We observe here that Δt has a “memory” in the sense that if a very low Δt was needed at a step k , then it will remain small for some steps. One ought to observe that this procedure is somewhat computationally expensive, as it computes at every RK steps two different steps. In order to limit the computational expenses, a common approach is to use two Butcher tableau that *differ only* in the b coefficients, i.e. in forming the linear combination (6.26d), while the coefficients a and c are the same such that equations (6.26a)-(6.26c) need to be computed only once. The two steps x_{k+1} and \hat{x}_{k+1} are therefore computed as:

$$x_{k+1} = x_k + \Delta t \sum_{i=1}^s b_i K_i \quad (6.44a)$$

$$\hat{x}_{k+1} = x_k + \Delta t \sum_{i=1}^s \hat{b}_i K_i \quad (6.44b)$$

The Butcher tableau for adaptive RK methods is then presented with two lines in the “ b ” part of the tableau, the first one for the “classic” b and the second for \hat{b} . E.g. the Butcher tableau used in the ode45.m Matlab function reads as¹²:

0								
1/5	1/5							
3/10	3/40	9/40						
4/5	44/45	-56/15	32/9					
8/9	19372/6561	-25360/2187	64448/6561	-212/729				
1	9017/3168	-355/33	46732/5247	49/176	-5103/18656			
1	35/384	0	500/1113	125/192	-2187/6784	11/84		
	35/384	0	500/1113	125/192	-2187/6784	11/84	0	
	5179/57600	0	7571/16695	393/640	-92097/339200	187/2100	1/40	

One can observe that this method has $s = 7$ stages. Both steps x_{k+1} and \hat{x}_{k+1} are of order $o = 5$. Figure 6.18 illustrates the step size Δt selected by the ode45.m when treating the Van Der Pol oscillator (6.44).

¹²the RK method coded by this tableau is referred to as the RK45 Dormand-Prince method

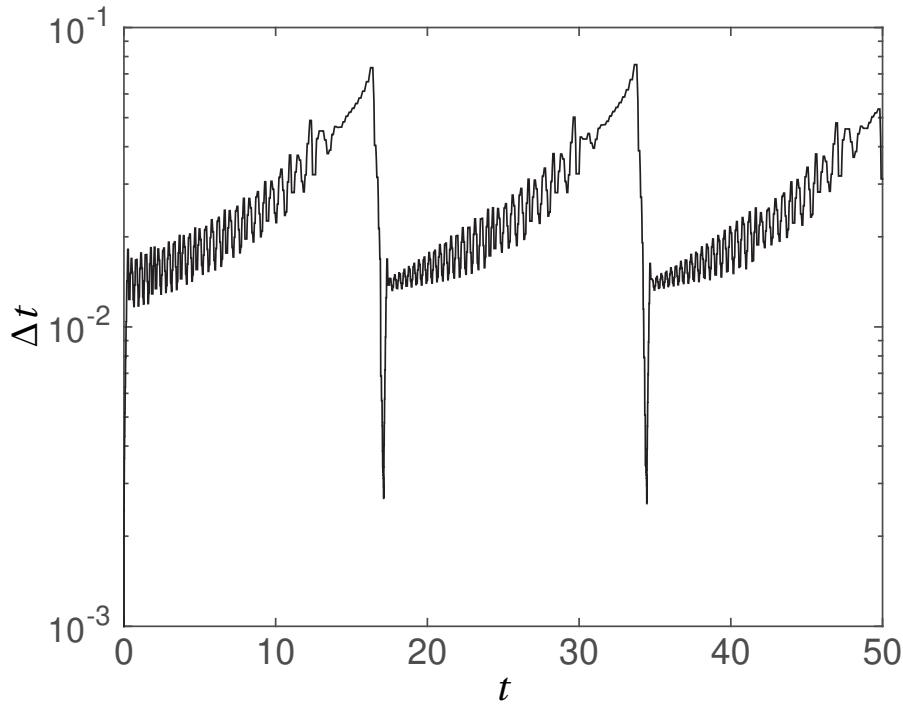


Figure 6.18: Illustration of the step sizes selected by `ode45.m` for simulating the Van Der Pol trajectories. Whenever the trajectories are going through a “sharp” turn, the step size drops drastically in order to keep the integration error under control. It then increases slowly again to reduce the computational burden.

One can finally note that the last line of the “a” part of the tableau is identical to the first line of the “b” part, and that the K_1 is given by:

$$K_1 = f(x_k, u(t_k)) \quad (6.45)$$

This implies that K_1 of step $k+1$ matches K_7 of step k such that the last K_7 can be reused for the next K_1 .

7 IMPLICIT INTEGRATION METHODS - RUNGE-KUTTA

In the explicit RK section, we have seen that a Butcher tableau can define explicit or implicit RK methods. We have discussed in details explicit RK methods. It is time now to discuss implicit RK methods. While implicit methods are more expensive computationally (one needs to solve equations at every RK step, as opposed to simply evaluate the K_i sequentially as in explicit methods), we will see that they have some striking advantages. Indeed, Implicit RK methods

- can achieve very high and systematic orders
- are stable regardless of the step size Δt
- are ideal for handling the simulation of DAEs

7.1 IMPLICIT EULER METHOD

Let us start with approaching implicit integrators from a basic point of view, by studying the implicit Euler method. Recall that the *explicit* Euler method uses the iteration:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta t \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) \quad (7.1)$$

where $\mathbf{u}_k = \mathbf{u}(t_k)$. The *implicit* Euler method uses instead:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta t \mathbf{f}(\mathbf{x}_{k+1}, \mathbf{u}_{k+1}) \quad (7.2)$$

The name “implicit” ought to be fairly clear from (7.2). Indeed, obtaining \mathbf{x}_{k+1} from (7.2) cannot be done via a simple function evaluation, but must be done via *solving* (7.2) for \mathbf{x}_{k+1} , i.e. one needs to find a solution to

$$\mathbf{r}(\mathbf{x}_{k+1}, \mathbf{x}_k, \mathbf{u}_{k+1}) := \mathbf{x}_k + \Delta t \mathbf{f}(\mathbf{x}_{k+1}, \mathbf{u}_{k+1}) - \mathbf{x}_{k+1} = 0 \quad (7.3)$$

in terms of \mathbf{x}_{k+1} . Note that $\mathbf{r} \in \mathbb{R}^n$ where n is the size of the state space, i.e. $\mathbf{x} \in \mathbb{R}^n$. Solving (7.3) is typically done via applying the Newton method covered earlier in these notes. The deployment of the implicit Euler algorithm is then detailed below.

Algorithm: Implicit Euler method

Input: Initial conditions \mathbf{x}_0 , input profile $\mathbf{u}(\cdot)$, step size Δt

for $k = 0, \dots, N - 1$ **do**

Guess \mathbf{x}_{k+1} , one can e.g. use $\mathbf{x}_{k+1} = \mathbf{x}_k$

while $\|\mathbf{r}(\mathbf{x}_{k+1}, \mathbf{x}_k, \mathbf{u}_{k+1})\| > \text{Tol}$ **do**

Compute the solution $\Delta \mathbf{x}_{k+1}$ to the linear system:

$$\frac{\partial \mathbf{r}(\mathbf{x}_{k+1}, \mathbf{x}_k, \mathbf{u}_{k+1})}{\partial \mathbf{x}_{k+1}} \Delta \mathbf{x}_{k+1} + \mathbf{r}(\mathbf{x}_{k+1}, \mathbf{x}_k, \mathbf{u}_{k+1}) = 0 \quad (7.4)$$

where \mathbf{r} is given by (7.3). Update:

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_{k+1} + \alpha \Delta \mathbf{x}_{k+1} \quad (7.5)$$

for some step size $\alpha \in]0, 1]$ (a full step $\alpha = 1$ generally works for implicit integrators)

return $\mathbf{x}_{1, \dots, N}$

Note that this procedure is significantly more complex than the simple update (7.1) of the explicit Euler method. In particular, computing the Newton step (7.4) requires computing the Jacobian $\frac{\partial \mathbf{r}(\mathbf{x}_{k+1}, \mathbf{x}_k, \mathbf{u}_{k+1})}{\partial \mathbf{x}_{k+1}}$ and forming its matrix factorization (i.e. solving the linear system). This procedure must be repeated at *each* RK step $\mathbf{x}_k \rightarrow \mathbf{x}_{k+1}$, i.e. many times in a complete simulation. Note that the Jacobian $\frac{\partial \mathbf{r}(\mathbf{x}_{k+1}, \mathbf{x}_k, \mathbf{u}_{k+1})}{\partial \mathbf{x}_{k+1}}$ reads as:

$$\frac{\partial \mathbf{r}(\mathbf{x}_{k+1}, \mathbf{x}_k, \mathbf{u}_{k+1})}{\partial \mathbf{x}_{k+1}} = \frac{\partial \mathbf{f}(\mathbf{x}_{k+1}, \mathbf{u}_{k+1})}{\partial \mathbf{x}_{k+1}} - I \quad (7.6)$$

and requires computing the Jacobian on the model dynamics $\frac{\partial f(x_{k+1}, u_{k+1})}{\partial x_{k+1}}$.

The implicit Euler method has a **global error of order 1**, i.e.

$$\|x_N - x(T)\| \leq c\Delta t \quad (7.7)$$

for some $c > 0$ and for Δt sufficiently small. Hence the order of the implicit Euler method is identical to the explicit Euler method.

7.1.1 STABILITY OF THE IMPLICIT EULER METHOD

Let us now unpack one of the main motivations for using, in some cases, implicit methods for simulation a model. Let us come back to the stability issue of the explicit Euler method (7.1), and recall that the iteration is unstable on the (stable) test system

$$\dot{x} = \lambda x \quad (7.8)$$

for

$$\Delta t \lambda > 2 \quad (7.9)$$

We can make a very similar computation for implicit Euler scheme, i.e. we observe that the Implicit Euler method (7.2) deployed on the test system (7.8) reads as

$$x_{k+1} = x_k + \Delta t \lambda x_{k+1} \quad (7.10)$$

or equivalently

$$x_{k+1} = \frac{1}{1 - \lambda \Delta t} x_k \quad (7.11)$$

Hence the Implicit Euler method (7.2) is stable if

$$|1 - \lambda \Delta t| > 1 \quad (7.12)$$

Observe that $\mathcal{R}(\lambda) < 0$ is required in order for (7.8) to be stable, and it follows that (7.12) always holds. This result is quite striking. Indeed, it entails that **the implicit Euler method is stable regardless of how “fast” the time constants of the model are**. This property is called **A-stability**, and means that the whole left-hand complex plane is stable (as opposed to the limited regions depicted in Figure 6.14). A-stability allows one to treat stiff dynamics without taking special care of the instability of the method, i.e. by taking “fairly large” steps Δt despite the fast time constant λ of the model.

Let us reuse the stiff model (6.42) and apply the implicit Euler method to simulate it. The outcome is illustrated in Fig. 7.1 for $\Delta t = 1 \cdot 10^{-3}$ and $\Delta t = 2 \cdot 10^{-2}$.

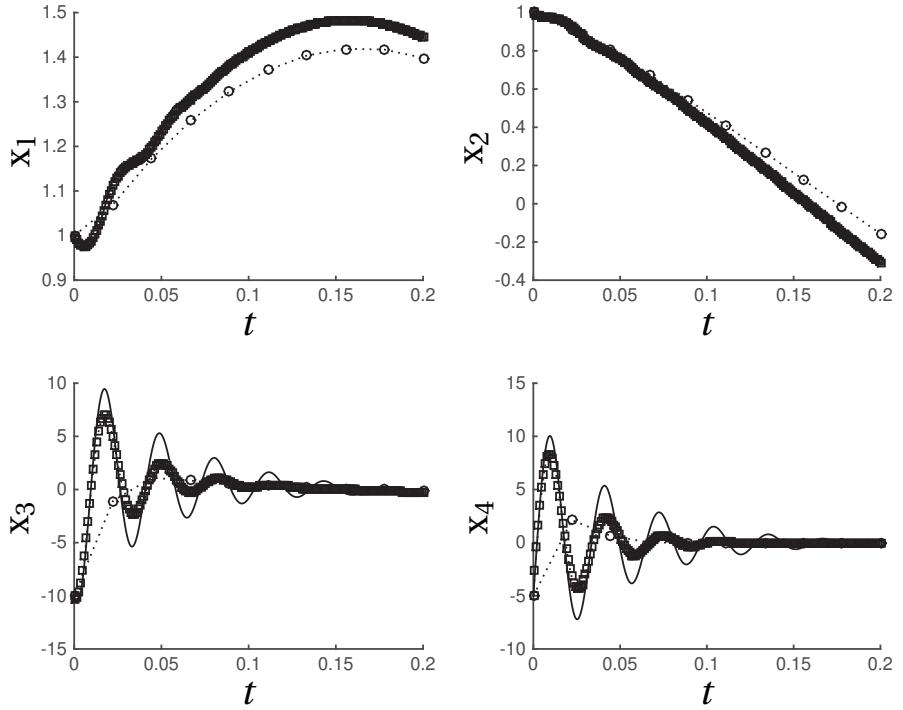


Figure 7.1: Illustration of the trajectories of the stiff system (6.42) and its simulation via the implicit Euler method using $\Delta t = 2 \cdot 10^{-2}$ (round markers) and $\Delta t = 1 \cdot 10^{-3}$ (square markers). One can observe that the implicit Euler, for “long” steps remains stable and approximates the fast dynamics by decaying to their damped values.

7.2 IMPLICIT RUNGE-KUTTA METHODS

Let us now study higher-order implicit Runge-Kutta methods. As already explained in the explicit RK section, a Butcher tableau that is not lower-diagonal describes an implicit method.

Then the RK equations (6.26) recalled here:

$$K_1 = f\left(x_k + \Delta t \sum_{j=1}^s a_{1j} K_j, u(t_k + c_1 \Delta t)\right) \quad (7.13a)$$

⋮

$$K_i = f\left(x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, u(t_k + c_i \Delta t)\right) \quad (7.13b)$$

⋮

$$K_s = f\left(x_k + \Delta t \sum_{j=1}^s a_{sj} K_j, u(t_k + c_s \Delta t)\right) \quad (7.13c)$$

$$x_{k+1} = x_k + \Delta t \sum_{i=1}^s b_i K_i \quad (7.13d)$$

are implicit in $K_{1,\dots,s}$, as $a_{ij} \neq 0$ for some $j \geq i$. In this case we need to solve (7.13a)-(7.13c) numerically, typically using the Newton method. more specifically, similarly to (7.3), we write:

$$r(K, x_k, u(.)) := \begin{bmatrix} f\left(x_k + \Delta t \sum_{j=1}^s a_{1j} K_j, u(t_k + c_1 \Delta t)\right) - K_1 \\ \vdots \\ f\left(x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, u(t_k + c_i \Delta t)\right) - K_i \\ \vdots \\ f\left(x_k + \Delta t \sum_{j=1}^s a_{sj} K_j, u(t_k + c_s \Delta t)\right) - K_s \end{bmatrix} = 0 \quad (7.14)$$

and deploy a Newton method on the “unknown”

$$K = \begin{bmatrix} K_1 \\ \vdots \\ K_s \end{bmatrix} \quad (7.15)$$

to solve (7.14). Note that $r \in \mathbb{R}^{n \cdot s}$ and $K \in \mathbb{R}^{n \cdot s}$ where n is the dimension of the state space (i.e. $x \in \mathbb{R}^n$) and s the number of stages of the IRK method. The Newton method then works as follows:

Algorithm: IRK for explicit ODEs

Input: Initial conditions x_0 , input profile $u(\cdot)$, Butcher tableau, step size Δt

for $k = 0, \dots, N - 1$ **do**

Guess for K (one can e.g. use $K_i = x_k$)

while $\|r(K, x_k, u(\cdot))\| > Tol$ **do**

Compute the solution Δw to the linear system:

$$\frac{\partial r(K, x_k, u(\cdot))}{\partial K} \Delta K + r(K, x_k, u(\cdot)) = 0 \quad (7.16)$$

with r given by (7.14). Update:

$$K \leftarrow K + \alpha \Delta K \quad (7.17)$$

for some step size $\alpha \in]0, 1]$ (a full step $\alpha = 1$ generally works for implicit integrators)

Take RK step:

$$x_{k+1} = x_k + \Delta t \sum_{i=1}^s b_i K_i \quad (7.18)$$

return x_1, \dots, x_N

The main computational complexity of this procedure is typically solving the linear system (7.16) which involves the (possibly dense) Jacobian matrix $\frac{\partial r(K, x_k, u(\cdot))}{\partial K}$ of size $\mathbb{R}^{n \cdot s \times n \cdot s}$. Forming the Jacobian matrix can also be fairly expensive in terms of computational complexity as it requires evaluating the Jacobians of the system dynamics:

$$\frac{\partial f(x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, u(t_k + c_i \Delta t))}{\partial K_l} \in \mathbb{R}^{n \times n} \quad (7.19)$$

for $i, l = 1, \dots, s$.

7.3 STABILITY & ORDER OF IRK METHODS

Similarly to the implicit Euler method, **IRK methods are A-stable**, meaning that they can “survive” (i.e. they are stable for) any fast dynamics. Some of them have even stronger form of stability (L-stability), but we will not detail this in these notes. Concerning the order of IRK methods, the picture here is fairly striking. Indeed, recall that explicit RK methods achieve an order $o = s$ (number of stages) for $s \leq 4$ and then the order “stalls” and does not increase as fast as s . This was bad news for high-order explicit RK methods in terms of efficiency (complexity vs. accuracy).

In contrast, **implicit RK methods can achieve $o = 2s$ for any number of stages s** .

We can then revisit the “order table” that we saw in the explicit RK section, comparing the explicit RK methods (ERK) to the implicit RK methods (IRK)

Order	Stages required	Order	Stages required
ERK2	2	IRK2	1
ERK4	4	IRK4	2
ERK6	7	IRK6	3
ERK8	11	IRK8	4

and readily see that IRK methods require dramatically less stages to achieve the same order as ERK methods. Let us illustrate this statement in Figure 7.2.

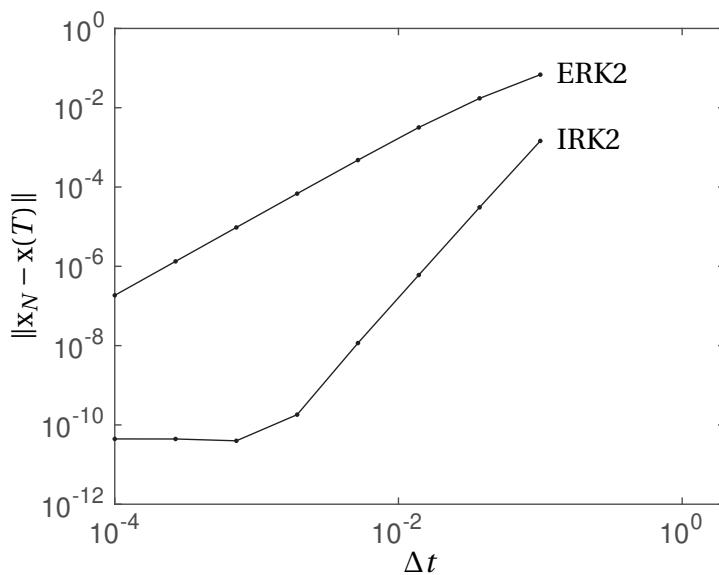


Figure 7.2: Comparison of the global error of the explicit RK2 method (6.19) and an implicit RK2 method, as obtained from numerical experiments. One can see the high order $o = 2 \cdot s$ of the IRK method to the order $o = s$ of the ERK method. One can also observe how the accuracy of the IRK method stops decreasing at about $10^{-10} - 10^{-11}$ where it reaches the accuracy of the linear algebra.

Unfortunately, this striking difference does not mean that IRK methods are necessarily better than the ERKs. In order to unpack that statement, let us investigate the efficiency of IRK methods (complexity for a given accuracy), using similar lines as what we did for the ERK methods. Let us assume that we want the global error to be limited to some given number Tol

- Then for an integrator of order o (even), we need:

$$\|x_N - x(T)\| \leq c\Delta t^o \leq \text{Tol} \quad (7.20)$$

such that the step size Δt is limited to:

$$\Delta t \leq \left(\frac{\text{Tol}}{c} \right)^{\frac{1}{o}} \quad (7.21)$$

- In order to carry out a simulation on the time interval $[0, T]$, the number of integrator step required per simulation time T is:

$$N = \frac{T}{\Delta t} \geq T \left(\frac{\text{Tol}}{c} \right)^{-\frac{1}{o}} \quad (7.22)$$

- The similarity with ERK methods stops here. Indeed, the complexity of IRK methods is typically dominated by solving the linear system (7.16). The complexity of solving this system is typically in the order of the cube of the size of the Jacobian matrix if the matrix is dense¹³, which is $n \cdot s$, where n is the state size. This system must typically be solved several times, say m in order to reach a good accuracy in solving the IRK equations. The complexity per time step is then in the order of:

$$\mathcal{C} = \mathcal{O}(mn^3s^3) \quad (7.23)$$

We deduce from this simple reasoning that the computational cost per unit of simulation time is of the order

$$\frac{\mathcal{C}}{T} = \mathcal{O}(Nm n^3 s^3) = \mathcal{O}\left(\left(\frac{\text{Tol}}{c}\right)^{-\frac{1}{o}} mn^3 s^3\right) = \mathcal{O}\left(\left(\frac{\text{Tol}}{c}\right)^{-\frac{1}{o}} \frac{mn^3 o^3}{8}\right) \quad (7.24)$$

One can put (7.24) in contrast with the complexity of explicit RK methods(6.38), recalled here:

$$\frac{n}{T} \geq s \left(\frac{\text{Tol}}{c} \right)^{-\frac{1}{o}} \quad (7.25)$$

but the comparison is arguably a bit difficult, as we are trying to compare evaluations of the model dynamics f in explicit methods to forming and solving linear systems in implicit methods. We can, however, compare the explicit and implicit approaches via numerical experiments. E.g. Fig. 7.3 depicts the computational time vs. integration accuracy for explicit and implicit RK methods of various orders (for Δt fixed) for the Van Der Pol oscillator (6.43). The observations we make in this specific example are fairly consistent for different models. Implicit methods can be a bit more computationally expensive than explicit ones, though the difference is typically mild. This picture is dramatically changed for stiff systems, where implicit integrators are typically required to have a numerically stable integration scheme.

¹³the cube can be lowered if the Jacobian is sparse or structured, which can be the case for IRK methods. Linear Algebra packages are very efficient at finding such structures, and lower the complexity of solving the linear system. Hence in numerical experiments, one often observe a lower complexity.

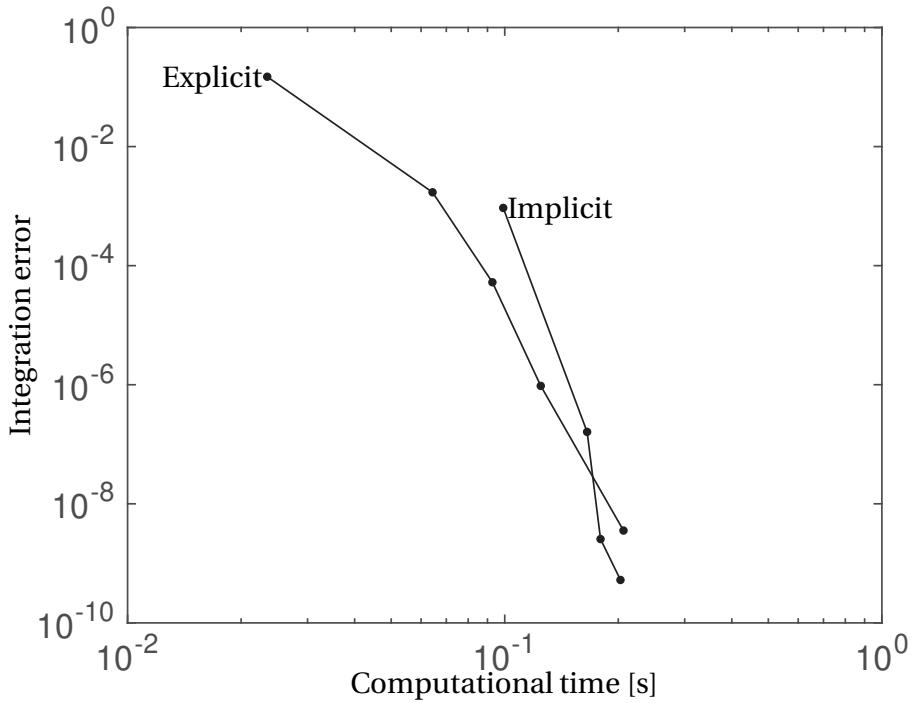


Figure 7.3: Illustration of the computational time vs. accuracy for IRK and ERK methods for the VanDerPol system (6.44) with $\Delta t = 10^{-2}$. One can observe the different behaviors of the implicit and explicit methods, and that even though implicit methods reach very high order for few stages, the computational cost of solving linear systems at every RK step makes the IRK method typically more expensive than the ERK ones. Note that these results are only illustrative, and may change depending on the coding language, computer architecture, system, etc.

7.4 COLLOCATION METHODS

There are different families of implicit RK methods, but in this course we will investigate the family of collocation methods, which have very strong properties, and which have a fairly intuitive interpretation. The key idea behind collocation methods is to approximate the model trajectories via polynomials. This is intrinsically an interpolation problem.

7.4.1 POLYNOMIAL INTERPOLATION

Interpolation based on polynomials typically uses a linear combination:

$$p(\tau, K) = \sum_{i=1}^s K_i \ell_i(\tau) \quad (7.26)$$

of a collection of polynomials $\ell_i \in \mathbb{R}$ built on the interval $\tau \in [0, 1]$ and weighted by the parameters $K_1, \dots, s \in \mathbb{R}^n$. Note that $p \in \mathbb{R}^n$ as well. A common choice of polynomial collection

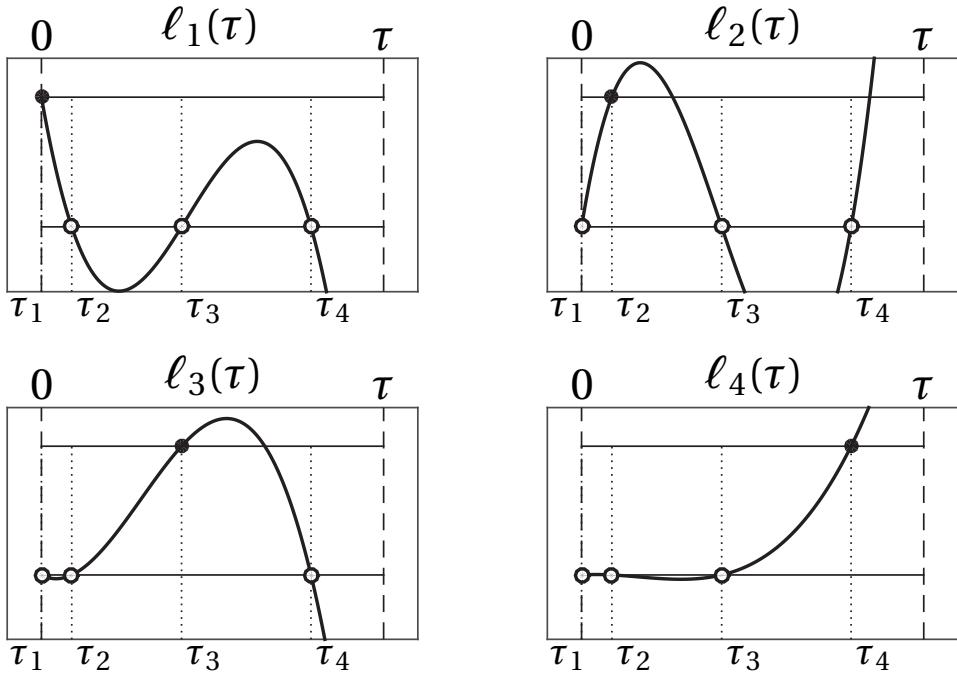


Figure 7.4: Illustration of the Lagrange polynomials for $s = 4$. The grid $\tau_{1,\dots,4}$ is displayed via the vertical dotted lines (note that $\tau_1 = 0$). Property (7.29) is readily visible.

$\ell_{1,\dots,s}$ is the so called Lagrange polynomials, which are constructed as:

$$\ell_i(\tau) = \prod_{j \neq i} \frac{\tau - \tau_j}{\tau_i - \tau_j} \quad (7.27)$$

based on a grid $\tau_{1,\dots,s} \in [0, 1]$ (see Figure 7.4). These polynomials have two interesting features

- They are orthogonal, i.e.

$$\int_0^1 \ell_i(\tau) \ell_j(\tau) d\tau = 0 \quad \text{if } i \neq j \quad (7.28)$$

- They satisfy:

$$\ell_i(\tau_j) = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases} \quad (7.29)$$

resulting in the following property on p:

$$p(\tau_i, K) = K_i, \quad i = 1, \dots, s \quad (7.30)$$

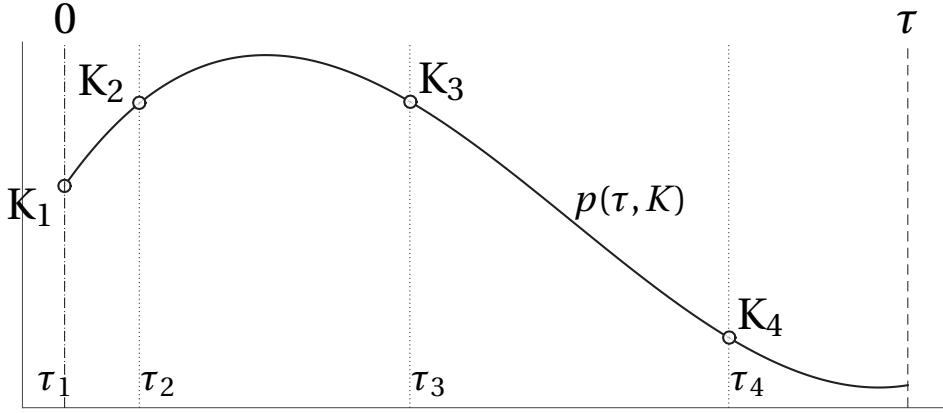


Figure 7.5: Illustration of the polynomial $p(\tau, K) = \sum_{i=1}^s K_i \ell_i(\tau)$ for $s = 4$ and for an arbitrary coefficient vector $K \in \mathbb{R}^4$, and $K_i \in \mathbb{R}$. One can observe the property (7.30), i.e. $p(\tau_i, v_k) = K_i$.

7.4.2 INTERPOLATION OF THE TRAJECTORIES

The key idea behind the collocation method is to approximate¹⁴ the true model trajectories by using¹⁵:

$$\dot{x}(t_k + \tau \cdot \Delta t) = p(\tau, K) = \sum_{i=1}^s K_i \ell_i(\tau) \quad (7.31)$$

on the interval $[t_k, t_{k+1}]$ (corresponding to $\tau \in [0, 1]$) by selecting the coefficients $K_{1,\dots,s} \in \mathbb{R}^n$. Note that each interval $[t_k, t_{k+1}]$ (for $k = 0, \dots, N - 1$) will have its own, possibly different coefficients

$$K = \begin{bmatrix} K_1 \\ \vdots \\ K_s \end{bmatrix} \quad (7.32)$$

where $K_i \in \mathbb{R}^n$.

We should now formulate equations that let us determine the parameters K . The key idea here is to enforce the dynamics on the grid points τ_1, \dots, τ_s , i.e. we want to determine the $K_{1,\dots,s}$ such that:

$$\dot{x}(t_k + \Delta t \cdot \tau_j) = p(\tau_j, K) = f(x(t_k + \Delta t \cdot \tau_j), u(t_k + \Delta t \cdot \tau_j)) \quad (7.33)$$

¹⁴sometimes the method is presented by considering that the trajectories $x(t)$ themselves are approximated by the polynomials instead of their time derivatives $\dot{x}(t)$. This does not change much the mathematics of the collocation schemes.

¹⁵Note that in the following developments we use the label $x(\cdot)$ to denote the trajectories provided by the integrator.

holds for $j = 1, \dots, s$ (i.e. on each grid point τ_j). We first observe that

$$\dot{x}(t_k + \tau_j \cdot \Delta t) = p(\tau_j, K) = K_j \quad (7.34)$$

holds by construction from (7.30). In order to further specify the equations above, we need to relate $x(t_k + \tau_j \cdot \Delta t)$ (i.e. the integral of (7.31)) to the coefficients K . We do this next. We first observe that:

$$x(t_k + \tau \cdot \Delta t) = x_k + \int_0^{\tau \cdot \Delta t} \dot{x}(t_k + \nu) d\nu \quad (7.35)$$

where x_k is the initial state of interval $[t_k, t_{k+1}]$. We can then make a change of variable $\nu = \xi \cdot \Delta t$, yielding $d\nu = \Delta t \cdot d\xi$, we can then rewrite (7.35) as:

$$x(t_k + \tau \cdot \Delta t) = x_k + \Delta t \int_0^\tau \dot{x}(t_k + \Delta t \cdot \xi) d\xi \quad (7.36)$$

We then use (7.31) in (7.36) to get:

$$x(t_k + \tau \cdot \Delta t) \approx x_k + \Delta t \int_0^\tau \sum_{i=1}^s K_i \ell_i(\xi) d\xi = \quad (7.37a)$$

$$= x_k + \Delta t \sum_{i=1}^s K_i \int_0^\tau \ell_i(\xi) d\xi \quad (7.37b)$$

$$= x_k + \Delta t \sum_{i=1}^s K_i L_i(\tau) \quad (7.37c)$$

where we define:

$$L_i(\tau) = \int_0^\tau \ell_i(\xi) d\xi \quad (7.38)$$

Note that $L_i(\tau)$ can be easily computed explicitly as it is simply the integration of the polynomial ℓ_i . We can then use (7.37c) in (7.33) to get:

$$\underbrace{\dot{x}(t_k + \Delta t \cdot \tau_j)}_{=K_j} = f \left(x_k + \Delta t \sum_{i=1}^s K_i L_i(\tau_j), u(t_k + \Delta t \cdot \tau_i) \right) \quad (7.39)$$

which should hold for all grid points $j = 1, \dots, s$. It follows that on a given interval $[t_k, t_{k+1}]$

the **collocation equations** read as:

$$K_1 = f \left(x_k + \Delta t \sum_{i=1}^s K_i L_i(\tau_1), u(t_k + \Delta t \cdot \tau_1) \right) \quad (7.40a)$$

$$\vdots \quad (7.40b)$$

$$K_j = f \left(x_k + \Delta t \sum_{i=1}^s K_i L_i(\tau_j), u(t_k + \Delta t \cdot \tau_j) \right) \quad (7.40c)$$

$$\vdots \quad (7.40d)$$

$$K_s = f \left(x_k + \Delta t \sum_{i=1}^s K_i L_i(\tau_s), u(t_k + \Delta t \cdot \tau_s) \right) \quad (7.40e)$$

$$x_{k+1} = x_k + \Delta t \cdot \sum_{i=1}^s K_i L_i(1) \quad (7.40f)$$

At this stage, it is very useful to compare (7.40) to the RK equations (6.26) and realise that they are identical if one defines:

$$a_{ji} = L_i(\tau_j), \quad b_i = L_i(1), \quad c_j = \tau_j \quad (7.41)$$

We also observe that if one picks the grid points τ_1, \dots, τ_s , then the polynomials $\ell_{1,\dots,s}(\tau)$ are defined and so is their integrals $L_i(\tau)$. **It follows that the coefficients (7.41) of the Butcher tableau are entirely defined via the grid points τ_1, \dots, τ_s .** We also construe here the role of the variables $K_{1,\dots,s}$ in RK methods. Indeed, (7.39) shows us that these variables are holding the state derivatives \dot{x} at the grid points τ_1, \dots, τ_s .

Note that not all IRK methods are collocation methods, as one could choose a dense Butcher tableau with coefficients that do not satisfy (7.41) for any choice of grid points τ_1, \dots, τ_s . Indeed, one can easily see that the Butcher tableau has $s^2 + 2s$ degrees of freedom, while the grid point selection offers only s degrees of freedom, such that not all choice of coefficients a, b, c can arise from (7.41).

We now need to briefly discuss how to choose the grid points τ_1, \dots, τ_s effectively. There are a few possible choices for a given number of stages s . The most commonly used for treating ODEs is the Gauss-Legendre method, which chooses the grid points τ_1, \dots, τ_s as the roots of the polynomial:

$$P_s(\tau) = \frac{1}{s!} \frac{d^s}{d\tau^s} \left[(\tau^2 - 1)^s \right] \quad (7.42)$$

i.e. we select the grid points τ_1, \dots, τ_s such that:

$$P_s(\tau_j) = 0, \quad j = 1, \dots, s \quad (7.43)$$

This selection rule may sound mysterious. Its motivation, though, has solid and deep roots in the Gauss quadrature theory, but we will not discuss this further here. Equipped with this rule, we have a procedure to build IRK methods (i.e. their Butcher tableau):

1. Select the number of stages s
2. Find the roots of (7.42) to get τ_1, \dots, τ_s
3. Build the polynomials $\ell_1(\tau), \dots, \ell_s(\tau)$ according to (7.27)
4. Build their integrals $L_i(\tau)$ according to (7.38)
5. Evaluate $L_i(\tau_j)$, $b_i = L_i(1)$ for $i, j = 1, \dots, s$
6. Build a computer code to solve the collocation equations (7.40)

This procedure **builds an IRK method with order $o = 2s$, and A-stable.**

An additional benefit of using polynomial interpolation for simulation models is that the polynomials provide us with an approximation of the state trajectories at any point in time, i.e. unlike other E/IRK methods, which deliver only the states on the time grid $t_{0, \dots, N}$, collocation methods can be prompted to deliver the states in-between, using (7.36), i.e. the true state trajectories (say x^{true}) is approximated at any time point in $[0, T]$ by:

$$x^{\text{true}}(t_k + \Delta t \tau) \approx x_k + \Delta t \int_0^\tau \dot{x}(t_k + \Delta t \cdot \xi) d\xi = x_k + \Delta t \sum_{i=1}^s K_i \int_0^\tau \ell_i(\xi) d\xi \quad (7.44)$$

Note that the order of approximation $o = 2s$ unfortunately only holds on the grid $t_{0, \dots, N}$, such that this approximation is typically a bit worse than order $2s$.

7.5 RK METHODS FOR IMPLICIT ODES

We have so far discussed numerical integration methods for models in the explicit ODE form:

$$\dot{x} = f(x, u) \quad (7.45)$$

Some models, however, are easier to treat in an implicit ODE form:

$$F(\dot{x}, x, u) = 0 \quad (7.46)$$

This is e.g. the case for models of complex mechanical system arising in Lagrange mechanics, taking the form (see (2.40)), where we use $v \equiv \dot{q}$:

$$\begin{bmatrix} I & 0 \\ 0 & W(q) \end{bmatrix} \underbrace{\begin{bmatrix} \dot{q} \\ \dot{v} \end{bmatrix}}_{\equiv \dot{x}} = \begin{bmatrix} v \\ Q + \nabla_q \mathcal{L} - \frac{\partial}{\partial q} (W(q)v) v \end{bmatrix} \quad (7.47)$$

and where the symbolic inverse of matrix $W(q)$ is very complex to write. In that case, it is best to *avoid* trying to form the explicit version of (7.47), i.e.

$$\begin{bmatrix} \dot{q} \\ \dot{v} \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & W(q) \end{bmatrix}^{-1} \begin{bmatrix} v \\ Q + \nabla_q \mathcal{L} - \frac{\partial}{\partial q} (W(q)v) v \end{bmatrix} \quad (7.48)$$

and work with the implicit form (7.47) directly. One trivial, but not necessarily effective approach would be to use an explicit RK method while forming the matrix inverse numerically in (7.48) at every integrator step.

A usually more effective approach is to use an implicit RK method directly on the implicit equation (7.47). Note that we can easily write (7.47) as (7.46), using:

$$F(\dot{x}, x) = \begin{bmatrix} I & 0 \\ 0 & W(q) \end{bmatrix} \dot{x} - \begin{bmatrix} v \\ \nabla_q \mathcal{L} - \frac{\partial}{\partial q}(W(q)v)v \end{bmatrix} = 0 \quad (7.49)$$

We should then focus on treating (7.46) numerically in implicit integrators. Interestingly enough, the IRK equations (6.26) apply almost directly to implicit models, using only a minor modification. Recall that (6.26) for an explicit model of the form $\dot{x} = f(x, u)$ reads as:

$$K_1 = f\left(x_k + \Delta t \sum_{j=1}^s a_{1j} K_j, u(t_k + c_1 \Delta t)\right) \quad (7.50a)$$

⋮

$$K_i = f\left(x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, u(t_k + c_i \Delta t)\right) \quad (7.50b)$$

⋮

$$K_s = f\left(x_k + \Delta t \sum_{j=1}^s a_{sj} K_j, u(t_k + c_s \Delta t)\right) \quad (7.50c)$$

$$x_{k+1} = x_k + \Delta t \sum_{j=1}^s b_j K_j \quad (7.50d)$$

Recall also that the variables $K_{1,\dots,s}$ in RK method are holding the state derivatives \dot{x} at the grid points τ_1, \dots, τ_s . The modification of the IRK equations for treating (7.46) therefore read as:

$$F\left(K_1, x_k + \Delta t \sum_{j=1}^s a_{1j} K_j, u(t_k + c_1 \Delta t)\right) = 0 \quad (7.51a)$$

⋮

$$F\left(K_i, x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, u(t_k + c_i \Delta t)\right) = 0 \quad (7.51b)$$

⋮

$$F\left(K_s, x_k + \Delta t \sum_{j=1}^s a_{sj} K_j, u(t_k + c_s \Delta t)\right) = 0 \quad (7.51c)$$

$$x_{k+1} = x_k + \Delta t \sum_{j=1}^s b_j K_j \quad (7.51d)$$

The equations to solve are then

$$r(K, x_k, u(.)) := \begin{bmatrix} F\left(K_1, x_k + \Delta t \sum_{j=1}^s a_{1j} K_j, u(t_k + c_1 \Delta t)\right) \\ \vdots \\ F\left(K_i, x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, u(t_k + c_i \Delta t)\right) \\ \vdots \\ F\left(K_s, x_k + \Delta t \sum_{j=1}^s a_{sj} K_j, u(t_k + c_s \Delta t)\right) \end{bmatrix} = 0 \quad (7.52)$$

The rest works as IRK schemes for explicit ODEs (see algorithm “IRK for explicit ODEs”).

Here it is useful to reflect on the IRK method in comparison to using an ERK method for solving the implicit ODE (7.46). Indeed, since the ODE does not readily deliver the state derivative \dot{x} , even when using an explicit method, one would have to solve the implicit equation for \dot{x} at every stage of the RK step, typically via a Newton method. We can then observe that

- Deploying an explicit RK method would then requiring solving s implicit equations of size n (the size of the state space), in order to get an order $o = s$ (for $s \leq 4$).
- Deploying an implicit IRK method requires solving one implicit equation of size $n \cdot s$, in order to get an order $o = 2s$ (for collocation methods).

Comparing formally the efficiency of the two approaches can be tricky, but the bottom line here is that if the ODE model is implicit, using an implicit RK method to perform the simulation can be a very good choice, independently of questions of stiffness.

7.6 RK METHODS FOR IMPLICIT DAEs

In order to close this section, it remains to be discussed how to treat DAEs numerically. We have in fact already all the tools required to tackle this problem, we only need to clarify a few points carefully. Let us consider DAEs in the fully implicit form (5.22) recalled here:

$$F(\dot{x}, z, x, u) = 0 \quad (7.53)$$

DAEs can be treated very similarly to implicit ODEs, but we need to understand how the algebraic variables z are treated here.

At every time step, the algebraic variables z ought to be considered as “free” variables that need to be determined independently of the other time steps, and they need to be adjusted so that (7.53) holds. More specifically, while an implicit ODE ought to be treated via imposing

$$F\left(K_i, x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, u(t_k + c_i \Delta t)\right) = 0 \quad (7.54)$$

for $i = 1, \dots, s$ of each RK step $\mathbf{x}_k \rightarrow \mathbf{x}_{k+1}$, equivalently for a DAE of the form (7.53), we ought to impose

$$F\left(K_i, z_i, \mathbf{x}_k + \Delta t \sum_{j=1}^s a_{ij} K_j, \mathbf{u}(t_k + c_i \Delta t)\right) = 0 \quad (7.55)$$

for $i = 1, \dots, s$ of each RK step $\mathbf{x}_k \rightarrow \mathbf{x}_{k+1}$. Note that here the variables $z_i \in \mathbb{R}^m$ operate as "unknown" that must be determined alongside the variables $K_i \in \mathbb{R}^n$. The complete IRK equations for DAEs then read as:

$$\mathbf{r}(\mathbf{w}, \mathbf{x}_k, \mathbf{u}(.)) := \begin{bmatrix} F\left(K_1, z_1, \mathbf{x}_k + \Delta t \sum_{j=1}^s a_{1j} K_j, \mathbf{u}(t_k + c_1 \Delta t)\right) \\ \vdots \\ F\left(K_s, z_s, \mathbf{x}_k + \Delta t \sum_{j=1}^s a_{sj} K_j, \mathbf{u}(t_k + c_s \Delta t)\right) \\ \vdots \\ F\left(K_s, z_s, \mathbf{x}_k + \Delta t \sum_{j=1}^s a_{sj} K_j, \mathbf{u}(t_k + c_s \Delta t)\right) \end{bmatrix} = 0 \quad (7.56)$$

where we gathered the K_i, z_i variables in:

$$\mathbf{w} = \begin{bmatrix} K_1 \\ \vdots \\ K_s \\ z_1 \\ \vdots \\ z_s \end{bmatrix} \quad (7.57)$$

Note that if $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{z} \in \mathbb{R}^m$, then function F returns vectors of dimension \mathbb{R}^{n+m} . It follows that $\mathbf{w} \in \mathbb{R}^{s(n+m)}$. Moreover, the "residual" function r in (7.56) returns a vector of dimension $\mathbb{R}^{s(n+m)}$, and that solving (7.56) must provide the variables w . The set of equations (7.56) must be solved at each RK step $\mathbf{x}_k \rightarrow \mathbf{x}_{k+1}$ in order to build the complete simulation on the time interval $[0, T]$, and is done as previously using the Newton method.

For the sake of clarity, let us write down explicitly the pseudo-code required to perform the numerical simulation of a fully implicit DAE model.

Algorithm: IRK for Fully Implicit DAEs

Input: Initial conditions x_0 , input profile $u(\cdot)$, Butcher tableau, step size Δt

for $k = 0, \dots, N - 1$ **do**

 Guess for w (one can e.g. use $K_i = x_k, z_i = 0$)

while $\|r(w, x_k, u(\cdot))\| > Tol$ **do**

 Compute the solution Δw to the linear system:

$$\frac{\partial r(w, x_k, u(\cdot))}{\partial w} \Delta w + r(w, x_k, u(\cdot)) = 0 \quad (7.58)$$

 with r given by (7.56). Update:

$$w \leftarrow w + \alpha \Delta w \quad (7.59)$$

 for some step size $\alpha \in]0, 1]$ (a full step $\alpha = 1$ generally works for implicit integrators)

 Take RK step:

$$x_{k+1} = x_k + \Delta t \sum_{j=1}^s b_j K_j \quad (7.60)$$

return $x_{1,\dots,N}$

Note that the remarks at the end of section 7.5 also hold here, i.e. in order to deploy an explicit RK method on the DAE (7.53), one would have to solve it for \dot{x} and z at every stage of the RK steps, typically requiring the deployment of a Newton method.

7.6.1 RK METHOD FOR SEMI-EXPLICIT DAE MODELS

Adapting the approach presented above to the semi-explicit DAE case is fairly straightforward. Let us do it here nonetheless. Recall that semi-explicit DAEs are in the form:

$$\dot{x} = f(z, x, u) \quad (7.61a)$$

$$0 = g(z, x, u) \quad (7.61b)$$

Here we ought to impose at each $i = 1, \dots, s$:

$$K_i = f\left(z_i, x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, u(t_k + c_i \Delta t)\right) \quad (7.62a)$$

$$0 = g\left(z_i, x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, u(t_k + c_i \Delta t)\right) \quad (7.62b)$$

hence we form the RK equations:

$$r(w, x_k, u(.)) := \begin{bmatrix} f(z_1, x_k + \Delta t \sum_{j=1}^s a_{1j} K_j, u(t_k + c_1 \Delta t) - K_1) \\ g(z_1, x_k + \Delta t \sum_{j=1}^s a_{1j} K_j, u(t_k + c_1 \Delta t)) \\ \vdots \\ f(z_i, x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, u(t_k + c_i \Delta t) - K_i) \\ g(z_i, x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, u(t_k + c_i \Delta t)) \\ \vdots \\ f(z_s, x_k + \Delta t \sum_{j=1}^s a_{sj} K_j, u(t_k + c_s \Delta t) - K_s) \\ g(z_s, x_k + \Delta t \sum_{j=1}^s a_{sj} K_j, u(t_k + c_s \Delta t)) \end{bmatrix} = 0 \quad (7.63)$$

And apply the same procedure as described several times above to perform the simulation.

Let us know make a connection between the problem of DAE differential index and numerical simulations. Recall that a semi-explicit DAE (7.61) is of index 1 (i.e. an "easy" DAE) if the Jacobian $\frac{\partial g}{\partial z}$ is full rank (i.e. invertible) on the system trajectories. In order to make a simple point, let us consider an IRK method having a single stage $s = 1$ to treat the DAE. In that simple case, we observe that the RK equations (7.63) boil down to:

$$r(w, x_k, u(.)) = \begin{bmatrix} f(z_1, x_k + \Delta t a_{11} K_1, u(t_k + c_1 \Delta t) - K_1) \\ g(z_1, x_k + \Delta t a_{11} K_1, u(t_k + c_1 \Delta t)) \end{bmatrix} = 0 \quad (7.64)$$

and that the Jacobian of these equations read as:

$$\frac{\partial r(w, x_k, u(.))}{\partial w} = \left[\begin{array}{cc} \frac{\partial r(w, x_k, u(.))}{\partial K_1} & \frac{\partial r(w, x_k, u(.))}{\partial z_1} \end{array} \right] = \left[\begin{array}{cc} \Delta t a_{11} \frac{\partial f(x, u)}{\partial x} - I & \frac{\partial f(x, u)}{\partial z} \\ \Delta t a_{11} \frac{\partial g(x, u)}{\partial x} & \frac{\partial g(x, u)}{\partial z} \end{array} \right] \Bigg| \begin{array}{l} x = x_k + \Delta t a_{11} K_1 \\ z = z_1 \\ u = u(t_k + c_1 \Delta t) \end{array} \quad (7.65)$$

Suppose then that we let $\Delta t \rightarrow 0$, i.e. we investigate the behavior of the method for an "infinite" accuracy (i.e. very high). Then the Jacobian matrix tends to:

$$\lim_{\Delta t \rightarrow 0} \frac{\partial r(w, x_k, u(.))}{\partial w} = \begin{bmatrix} -I & \frac{\partial f}{\partial z} \\ 0 & \frac{\partial g}{\partial z} \end{bmatrix} \quad (7.66)$$

which is full rank (invertible) if and only if $\frac{\partial g(x, u)}{\partial z}$ is full rank. Recall that we need to deploy a Newton method in order to solve $r(w, x_k, u(.)) = 0$, where we need the Jacobian $\frac{\partial r(w, x_k, u(.))}{\partial w}$ to be invertible, and therefore requires $\frac{\partial g(x, u)}{\partial z}$ to be full rank. We can therefore conclude from this simple example that

The IRK methods we have investigated should not be deployed on DAEs of index above 1.

Note that if $\Delta t > 0$, the rank of the Jacobian matrix may be nonetheless full, and one may conclude that the method is nonetheless fine. However, the small analysis above shows us that “something is intrinsically wrong” when deploying the IRK schemes we have studied on high-index DAEs (index above 1). We will leave out of this course further analysis of this question. Before closing this section, it ought to be underlined here that one can deploy some of the RK schemes we have investigated on high-index DAEs, and get a trajectory computed from the method. I.e. there may not be a clear sign when deploying the method that “something is not right”. The trajectories obtained, however, are typically non-sensical. Hence one ought to check the index of the DAE before trusting the output of a classic RK method.

8 SENSITIVITY OF SIMULATIONS

Our last section will deal with a topic that is crucial in many fields of engineering that make use of numerical simulations. In many problems it is important not only to compute accurate and reliable simulations for an ODE or a DAE, but also to get some information on how this simulation would be affected by a change in the model parameters. To be a bit more formal, let us consider an ODE depending on a fixed parameter p :

$$\dot{x} = f(x, p), \quad x(0) = x_0 \quad (8.1)$$

Note that we will omit the inputs in our developments here, and discuss in the end what role they play. Suppose that we have computed a simulation of (8.1), i.e. we have a sequence:

$$x_k \approx x(t_k) \quad (8.2)$$

on a given time grid t_0, \dots, t_N . Clearly, if one was to modify “a bit” either the initial conditions x_0 and/or the parameters p , then the sequence $x_{1,\dots,N}$ would be affected. For nonlinear dynamics f , one cannot assess the impact of changing the parameters and/or initial conditions x_0 on the simulation without redoing the entire simulations. However, one can use first-order arguments and compute the approximate effect of changing the parameters and initial conditions i.e. one can try to compute:

$$\frac{\partial x_k}{\partial x_0}, \quad \frac{\partial x_k}{\partial p} \quad (8.3)$$

for $k = 1, \dots, N$, evaluated at x_0, p given. The sensitivity of simulations then deals with computing these derivatives.

8.1 VARIATIONAL APPROACH

Before giving the answer to the question raised above, let us investigate the derivative of the trajectories $x(t)$ generated by the ODE (8.1) with respect to p, x_0 , i.e. let us forget for

now the numerical aspect of the problem. The question of computing

$$\frac{\partial \mathbf{x}(t)}{\partial p} \quad \frac{\partial \mathbf{x}(t)}{\partial \mathbf{x}_0} \quad (8.4)$$

can appear tricky, but let us consider the following construction. Suppose that we know the trajectory $\mathbf{x}(t)$ associated to (8.1). We observe that:

$$\frac{d}{dt} \frac{\partial \mathbf{x}(t)}{\partial \mathbf{x}_0} = \frac{\partial \dot{\mathbf{x}}(t)}{\partial \mathbf{x}_0} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}(t)}{\partial \mathbf{x}_0} \quad (8.5a)$$

$$\frac{d}{dt} \frac{\partial \mathbf{x}(t)}{\partial p} = \frac{\partial \dot{\mathbf{x}}(t)}{\partial p} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}(t)}{\partial p} + \frac{\partial f}{\partial p} \quad (8.5b)$$

where all Jacobians are evaluated at $\mathbf{x}(t)$ and p . We moreover observe that

$$\frac{\partial \mathbf{x}(0)}{\partial \mathbf{x}_0} = I, \quad \frac{\partial \mathbf{x}(0)}{\partial p} = 0 \quad (8.6)$$

Let us label:

$$A(t) = \frac{\partial \mathbf{x}(t)}{\partial \mathbf{x}_0}, \quad B(t) = \frac{\partial \mathbf{x}(t)}{\partial p} \quad (8.7)$$

We can then write the dynamic system for A and B using (8.5)-(8.6):

$$\dot{A}(t) = \frac{\partial f}{\partial \mathbf{x}} A(t) \quad A(0) = I \quad (8.8a)$$

$$\dot{B}(t) = \frac{\partial f}{\partial \mathbf{x}} B(t) + \frac{\partial f}{\partial p} \quad B(0) = 0 \quad (8.8b)$$

We then observe that (8.8) is a linear time-varying system with a matrix-based state space (A and B). It is often referred to as the *variational equations* associated to (8.1). In principle, one could answer the sensitivity question raised above by performing the integration of the dynamic (8.1) together with (8.8) using any numerical integration method.

Such an approach, however, would not deliver exact derivatives. Indeed, the variational equations are expressing the sensitivity of the trajectories of the system, whereas the sequence $\mathbf{x}_{1,\dots,N}$ arises from numerical simulation and it approximates the trajectories. Moreover, a numerical integration of (8.8) is also an approximation of the trajectories $A(t)$, $B(t)$. As a result the sensitivities generated can be quite inaccurate.

For these reasons, the variational approach is not very often used to treat the sensitivity question, and rather replaced by so called *Algorithmic Differentiation* methods, which seek to compute *exact derivatives* of the inexact simulation of the dynamics. Let us briefly study these techniques next.

8.2 ALGORITHMIC DIFFERENTIATION OF THE EXPLICIT EULER SCHEME

As usual, it will be easiest to first consider the problem of sensitivity computation for the explicit Euler scheme:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta t f(\mathbf{x}_k, p) \quad (8.9)$$

Computing (8.3) for the simulation (8.9) boils down to a careful application of the chain rule on (8.9). We observe that:

$$\frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{x}_0} = \frac{\partial \mathbf{x}_k}{\partial \mathbf{x}_0} + \Delta t \frac{\partial f}{\partial \mathbf{x}_k} \frac{\partial \mathbf{x}_k}{\partial \mathbf{x}_0}, \quad \frac{\partial \mathbf{x}_0}{\partial \mathbf{x}_0} = I \quad (8.10a)$$

$$\frac{\partial \mathbf{x}_{k+1}}{\partial p} = \frac{\partial \mathbf{x}_k}{\partial p} + \Delta t \left(\frac{\partial f}{\partial \mathbf{x}_k} \frac{\partial \mathbf{x}_k}{\partial p} + \frac{\partial f}{\partial p} \right), \quad \frac{\partial \mathbf{x}_0}{\partial p} = 0 \quad (8.10b)$$

where all Jacobians are evaluated on the baseline simulation $\mathbf{x}_{1,\dots,N}$ and for the parameter value p . We observe that (8.10) is in fact a discrete linear dynamic system, assigning dynamics to the matrices (8.3). Indeed, let us label

$$A_k = \frac{\partial \mathbf{x}_k}{\partial \mathbf{x}_0}, \quad B_k = \frac{\partial \mathbf{x}_k}{\partial p} \quad (8.11)$$

And rewrite (8.10) as:

$$A_{k+1} = \left(I + \Delta t \frac{\partial f}{\partial \mathbf{x}_k} \right) A_k, \quad A_0 = I \quad (8.12a)$$

$$B_{k+1} = \left(I + \Delta t \frac{\partial f}{\partial \mathbf{x}_k} \right) B_k + \Delta t \frac{\partial f}{\partial p}, \quad B_0 = 0 \quad (8.12b)$$

This procedure is labelled Algorithmic Differentiation. An explicit Euler integrator with sensitivity generation can then be implemented using the following code:

Algorithm: Explicit Euler with sensitivity generation

Input: Initial conditions \mathbf{x}_0 , parameter p , Δt

Set $A_0 = I$, $B_0 = 0$

for $k = 0, \dots, N - 1$ **do**

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta t f(\mathbf{x}_k, p) \quad (8.13a)$$

$$A_{k+1} = \left(I + \Delta t \frac{\partial f}{\partial \mathbf{x}_k} \right) \Big|_{\mathbf{x}_k, p} A_k \quad (8.13b)$$

$$B_{k+1} = \left(I + \Delta t \frac{\partial f}{\partial \mathbf{x}_k} \right) \Big|_{\mathbf{x}_k, p} B_k + \Delta t \frac{\partial f}{\partial p} \Big|_{\mathbf{x}_k, p} \quad (8.13c)$$

return $\mathbf{x}_{1,\dots,N}$, $A_{1,\dots,N}$ and $B_{1,\dots,N}$

Note that if only the final state \mathbf{x}_N is needed, together with the associated sensitivities A_N and B_N , then a lot of memory can be saved by reusing the memory cells assigned to storing

\mathbf{x} , A and B .

In contrast to the variational approach, the sensitivities delivered by processing the discrete dynamics (8.8) provides the exact (up to machine precision) derivative of the sequence $\mathbf{x}_{1,\dots,N}$ delivered by the explicit Euler scheme, even though that sequence is not an exact (up to machine precision) approximation of the true trajectories.

An attentive reader, however, may have observed that if one treats both (8.1)-(8.8) using an explicit Euler scheme, then one recovers the algorithm above. Hence in this specific case, the variational equations and the Algorithmic Differentiation coincide. This observation is not true in general.

8.3 ALGORITHMIC DIFFERENTIATION OF EXPLICIT RUNGE-KUTTA METHODS

The explicit RK equations for (8.1) arising from a lower-diagonal Butcher tableau can be written as:

$$K_1 = f(\mathbf{x}_k, \mathbf{p}) \quad (8.14a)$$

⋮

$$K_i = f\left(\mathbf{x}_k + \Delta t \sum_{j=1}^{i-1} a_{ij} K_j, \mathbf{p}\right) \quad (8.14b)$$

⋮

$$K_s = f\left(\mathbf{x}_k + \Delta t \sum_{j=1}^{s-1} a_{sj} K_j, \mathbf{p}\right) \quad (8.14c)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta t \sum_{j=1}^s b_j K_j \quad (8.14d)$$

(observe that the summations in the formation of the $K_{1,\dots,s}$ entail by construction an explicit method). Although it is a bit tedious, one can apply chain rules to (8.14) in order to extract the sensitivities. Indeed, one can observe that:

- The differentiation of (8.14d) reads as:

$$\frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{x}_0} = \left(I + \Delta t \sum_{j=1}^s b_j \frac{\partial K_j}{\partial \mathbf{x}_k} \right) \frac{\partial \mathbf{x}_k}{\partial \mathbf{x}_0} \quad (8.15a)$$

$$\frac{\partial \mathbf{x}_{k+1}}{\partial \mathbf{p}} = \left(I + \Delta t \sum_{j=1}^s b_j \frac{\partial K_j}{\partial \mathbf{x}_k} \right) \frac{\partial \mathbf{x}_k}{\partial \mathbf{p}} + \Delta t \sum_{j=1}^s b_j \frac{\partial K_j}{\partial \mathbf{p}} \quad (8.15b)$$

- The differentiation of (8.14a)-(8.14c) reads as:

$$\frac{\partial K_i}{\partial x_k} = \frac{\partial f}{\partial x} \Big|_{x_k + \Delta t \sum_{j=1}^{i-1} a_{sj} K_j, p} \cdot \left(I + \Delta t \sum_{j=1}^{i-1} b_j \frac{\partial K_j}{\partial x_k} \right) \quad (8.16a)$$

$$\frac{\partial K_i}{\partial p} = \frac{\partial f}{\partial x} \Big|_{x_k + \Delta t \sum_{j=1}^{i-1} a_{sj} K_j, p} \cdot \left(I + \Delta t \sum_{j=1}^{i-1} b_j \frac{\partial K_j}{\partial p} \right) + \frac{\partial f}{\partial p} \Big|_{x_k + \Delta t \sum_{j=1}^{i-1} a_{sj} K_j, p} \quad (8.16b)$$

A pseudo-code deploying the above principle reads as follows.

Algorithm: Explicit RK scheme with sensitivity generation

Input: Initial conditions x_0 , parameter p , Δt

Set $A_0 = I$, $B_0 = 0$

for $k = 0, \dots, N - 1$ **do**

for $i = 1, \dots, s$ **do**

$$K_i = f \left(x_k + \Delta t \sum_{j=1}^{i-1} a_{ij} K_j, p \right) \quad (8.17a)$$

$$\frac{\partial K_i}{\partial x_k} = \frac{\partial f}{\partial x} \Big|_{x_k + \Delta t \sum_{j=1}^{i-1} a_{sj} K_j, p} \cdot \left(I + \Delta t \sum_{j=1}^{i-1} b_j \frac{\partial K_j}{\partial x_k} \right) \quad (8.17b)$$

$$\frac{\partial K_i}{\partial p} = \frac{\partial f}{\partial x} \Big|_{x_k + \Delta t \sum_{j=1}^{i-1} a_{sj} K_j, p} \cdot \left(I + \Delta t \sum_{j=1}^{i-1} b_j \frac{\partial K_j}{\partial p} \right) + \frac{\partial f}{\partial p} \Big|_{x_k + \Delta t \sum_{j=1}^{i-1} a_{sj} K_j, p} \quad (8.17c)$$

$$x_{k+1} = x_k + \Delta t \sum_{j=1}^s b_j K_j \quad (8.18a)$$

$$A_{k+1} = \left(I + \Delta t \sum_{j=1}^s b_j \frac{\partial K_j}{\partial x_k} \right) A_k \quad (8.18b)$$

$$B_{k+1} = \left(I + \Delta t \sum_{j=1}^s b_j \frac{\partial K_j}{\partial x_k} \right) B_k + \Delta t \sum_{j=1}^s b_j \frac{\partial K_j}{\partial p} \quad (8.18c)$$

return $x_{1,\dots,N}$, $A_{1,\dots,N}$ and $B_{1,\dots,N}$

Note that these computations are tedious to derive and to program. Fortunately, very efficient Algorithmic Differentiation tools such as e.g. CasADi can perform these operations automatically.

8.4 SENSITIVITY OF IMPLICIT RUNGE-KUTTA STEPS

Note that the algorithm above can only work for explicit RK methods, i.e. RK methods having a lower-diagonal Butcher tableau. Recall that the implicit RK equations generally read as:

$$K_1 = f\left(x_k + \Delta t \sum_{j=1}^s a_{1j} K_j, p\right) \quad (8.19a)$$

⋮

$$K_i = f\left(x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, p\right) \quad (8.19b)$$

⋮

$$K_s = f\left(x_k + \Delta t \sum_{j=1}^s a_{sj} K_j, p\right) \quad (8.19c)$$

$$x_{k+1} = x_k + \Delta t \sum_{j=1}^s b_j K_j \quad (8.19d)$$

Interestingly enough, the sensitivities of an implicit RK scheme are fairly straightforward to formulate and compute, i.e. easier than an explicit RK scheme. Recall that the implicit equations (8.19a)-(8.19c) are solved via a Newton method, i.e. we write (7.14), recalled here:

$$r(K, x_k, p) := \begin{bmatrix} f\left(x_k + \Delta t \sum_{j=1}^s a_{1j} K_j, u(t_k + c_1 \Delta t)\right) - K_1 \\ \vdots \\ f\left(x_k + \Delta t \sum_{j=1}^s a_{ij} K_j, u(t_k + c_i \Delta t)\right) - K_i \\ \vdots \\ f\left(x_k + \Delta t \sum_{j=1}^s a_{sj} K_j, u(t_k + c_s \Delta t)\right) - K_s \end{bmatrix} = 0 \quad (8.20)$$

and deploy a Newton method on the “unknown”

$$K = \begin{bmatrix} K_1 \\ \vdots \\ K_s \end{bmatrix} \quad (8.21)$$

i.e. we iterate (see algorithm “IRK for explicit ODEs”):

$$K \leftarrow K - \alpha \left(\frac{\partial r(K, x_k, p)}{\partial K} \right)^{-1} r(K, x_k, p) \quad (8.22)$$

Recall the Implicit Function theorem 5, and in particular equation (3.40). Here one ought to view

$$r(K, x_k, p) = 0 \quad (8.23)$$

as making K implicitly a function of x_k, p , and use equation (3.40) to deduce that:

$$\frac{\partial K}{\partial p} = - \frac{\partial r(K, x_k, p)}{\partial K}^{-1} \frac{\partial r(K, x_k, p)}{\partial p} \quad (8.24a)$$

$$\frac{\partial K}{\partial x_k} = - \frac{\partial r(K, x_k, p)}{\partial K}^{-1} \frac{\partial r(K, x_k, p)}{\partial x_k} \quad (8.24b)$$

One ought to observe that the matrix inverse (or rather the matrix factorization for the linear-algebra savy among the readers) needed in (8.22) can be reused in (8.24), such that the latter is inexpensive to compute. We then observe that (8.15) can be readily used to compute the sensitivities of x_{k+1} as in the explicit RK method (see (8.15)). We can summarize these observations in the following pseudo-code.

Algorithm: IRK for explicit ODEs with sensitivities

Input: Initial conditions x_0 , input profile $u(\cdot)$, Butcher tableau, step size Δt

Set $A_0 = I$, $B_0 = 0$

for $k = 0, \dots, N - 1$ **do**

 Guess for K (one can e.g. use $K_i = x_k$)

while $\|r(K, x_k, u(\cdot))\| > Tol$ **do**

 Compute the solution Δw to the linear system:

$$\frac{\partial r(K, x_k, u(\cdot))}{\partial K} \Delta K + r(K, x_k, u(\cdot)) = 0 \quad (8.25)$$

 with r given by (7.14). Update:

$$K \leftarrow K + \alpha \Delta K \quad (8.26)$$

 for some step size $\alpha \in]0, 1]$ (a full step $\alpha = 1$ generally works for implicit integrators)

 Reuse the matrix factorization required to solve (8.25) to compute:

$$\frac{\partial K}{\partial p} = -\frac{\partial r(K, x_k, p)}{\partial K}^{-1} \frac{\partial r(K, x_k, p)}{\partial p} \quad (8.27a)$$

$$\frac{\partial K}{\partial x_k} = -\frac{\partial r(K, x_k, p)}{\partial K}^{-1} \frac{\partial r(K, x_k, p)}{\partial x_k} \quad (8.27b)$$

Take RK step with sensitivities:

$$x_{k+1} = x_k + \Delta t \sum_{j=1}^s b_j K_j \quad (8.28a)$$

$$A_{k+1} = \left(I + \Delta t \sum_{j=1}^s b_j \frac{\partial K_j}{\partial x_k} \right) \Bigg| A_k \quad (8.28b)$$

$$B_{k+1} = \left(I + \Delta t \sum_{j=1}^s b_j \frac{\partial K_j}{\partial p} \right) \Bigg| B_k + \Delta t \sum_{j=1}^s b_j \frac{\partial K_j}{\partial p} \quad (8.28c)$$

return

8.5 SENSITIVITY WITH RESPECT TO INPUTS

Let us finish these lecture notes with discussing the computation of the sensitivity of a simulation with respect to the input u applied to the dynamics $\dot{x} = f(p, u)$. Since $u(\cdot)$ is, in general, a profile (i.e. a function of time in the interval $[0, T]$) as opposed to a finite set of parameters p , discussing sensitivities in the form explored here is a priori inadequate. However, the classic approach in numerical simulations is to consider that the input profile $u(\cdot)$ is in fact parameterized by a finite set of parameters, i.e. we consider that the input

is given by:

$$u(t, p) \quad (8.29)$$

Any derivative in the algorithms above is then simply replaced by its chain-rule expansion:

$$\frac{\partial \cdot}{\partial p} = \frac{\partial \cdot}{\partial u} \frac{\partial u}{\partial p} \quad (8.30)$$

e.g.

$$\frac{\partial f}{\partial p} = \frac{\partial f}{\partial u} \frac{\partial u}{\partial p} \quad (8.31)$$

The input parametrization can take many forms. It can e.g. be convenient to use an approach similar to the one used in collocation-based IRK schemes. In that context, the input profile is e.g. given by:

$$u(t_k + \tau \Delta t, p) = \sum_{i=1}^s p_{k,i} \ell_i(\tau) \quad (8.32)$$

where we have:

$$p = \begin{bmatrix} p_{0,1} \\ \vdots \\ p_{0,s} \\ \vdots \\ p_{N-1,1} \\ \vdots \\ p_{N-1,s} \end{bmatrix} \quad (8.33)$$

Hence the input profile is piece-wise smooth, and interpolates the points $p_{k,i}$ on the time intervals $[t_k, t_{k+1}]$ for $k = 0, N - 1$. This principle is illustrated in Figure 8.1. A wide-spread alternative input parametrization is to use the piecewise constant approach, which uses:

$$u(t_k + \tau \Delta t) = p_k, \quad k = 0, \dots, N - 1 \quad (8.34)$$

This principle is illustrated in Figure 8.2. All the sensitivity principles we studied above can be readily applied.

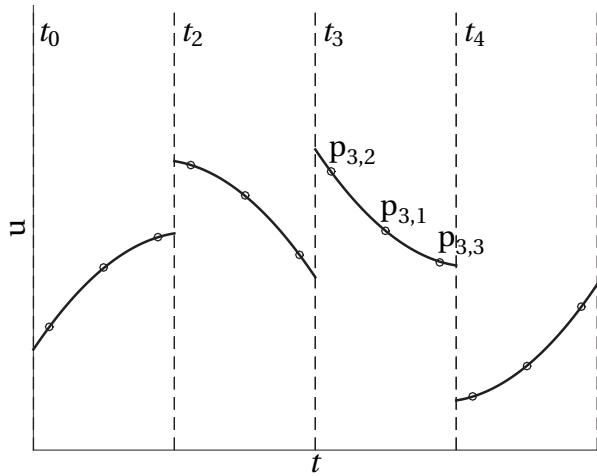


Figure 8.1: Illustration of the “collocation-based” parametrization of the inputs, for $N = 5$, and $s = 3$. Note that the input profile is piecewise-smooth.

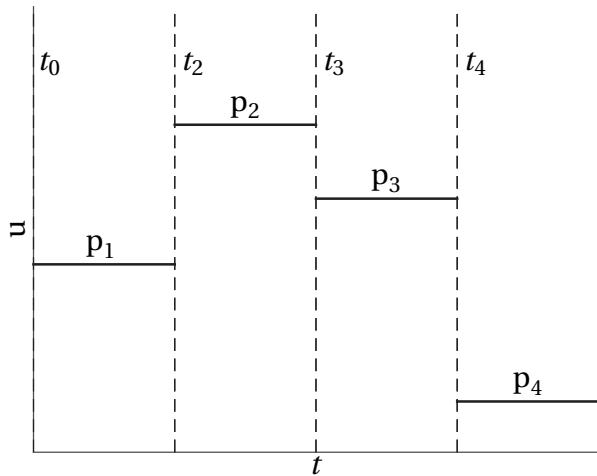


Figure 8.2: Illustration of the piecewise-constant approach to the input parametrization for $N = 4$.