

Chapter 1

Ordinary Differential Equations

In mathematics, an *ordinary differential equation* is a *differential equation* containing one or more functions of one independent variable and its derivatives. The term *ordinary* is used in contrast with the term *partial* so to specify that the differential equations has only *one* independent variable. This mathematical tool is widely used to describe dynamical systems (physical, social, economical, biological and so on and so forth). The first mathematicians to study and apply ordinary differential equations include important names like Newton, Euler, Leibniz, the Bernoulli family, Riccati, Clairaut and d'Alembert. An example of this kind of equations is probably one of the most known formulas in the world, that is trivially the Newton's second law of motion:

$$\mathcal{F}(z(t)) = m \frac{d^2}{dt^2} z(t) \quad (1.1)$$

where \mathcal{F} is a function of the unknown position function $z(t)$ at time t .

1.1 An introduction to Initial Value Problems

Before getting into any numerical solution scheme let us briefly recall the essential formal definitions and properties of ordinary differential equations and initial value problems.

Definition 1 (Ordinary differential equation). An *ordinary differential equation* (ODE) is an equation for a function $z(t)$, defined on an interval $I \subset \mathbb{R}$ and with values in the real or complex numbers or in the space \mathbb{R}^d (or \mathbb{C}^d), of the form:

$$F\left(t, z(t), z'(t), z''(t) \dots, z^{(n)}(t)\right) = 0 \quad (1.2)$$

Here F represents an arbitrary function of its arguments. The *order* n of a differential equation is the highest derivative which occurs. If the dimension d of the value range of $z(t)$ is higher than one, we talk about *systems of differential equations*.

Definition 2. An *explicit* differential equation of first order is a equation of the form:

$$z'(t) = f(t, z(t)) \quad \text{or shortly:} \quad z' = f(t, z) \quad (1.3)$$

A differential equation of order n is called explicit, if it is of the form:

$$z^{(n)}(t) = F\left(t, z(t), z'(t), z''(t) \dots, z^{(n-1)}(t)\right) \quad (1.4)$$

Definition 3. A differential equation of the form (1.3) is called *autonomous*, if the right hand side f is not explicitly dependent on t , namely:

$$z'(t) = f(z(t)) \quad (1.5)$$

Lemma 0.1. *Every differential equation of higher order can be written as a system of first-order differential equations. If the equation is explicit, then the system is explicit.*

Proof. By the introduction of additional variables $z_0(t) = z(t)$, $z_1(t) = z'(t)$ to $z_{n-1}(t) = z^{(n-1)}(t)$, each differential equation of order n can be transformed into a system of n differential equations of first order. This system has the form:

$$\begin{cases} z'_0(t) - z_1(t) \\ z'_1(t) - z_2(t) \\ \vdots \\ z'_{n-2}(t) - z_{n-1}(t) \\ F(t, z_0(t), z_1(t), \dots, z_{n-1}(t), z'_{n-1}(t)) \end{cases} = \begin{cases} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{cases} \quad (1.6)$$

In the case of an explicit equation, the system has the form:

$$\begin{cases} z'_0(t) \\ z'_1(t) \\ \vdots \\ z'_{n-2}(t) \\ z'_{n-1}(t) \end{cases} = \begin{cases} z_1(t) \\ z_2(t) \\ \vdots \\ z_{n-1}(t) \\ F(t, z_0(t), z_1(t), \dots, z_{n-1}(t)) \end{cases} \quad (1.7)$$

□

We will see that most of times the ODEs do not come alone. Almost always they are coupled with an *initial condition*, thus forming the so-called *initial value problem* (also known as *Cauchy problem*). The initial condition give information about the value of function z at a given point t_0 in the domain.

Definition 4 (Initial value problem). Given a point $(t_0, u_0) \in \mathbb{R} \times \mathbb{R}^d$. Furthermore, let the function $f(t, z(t))$ with values in \mathbb{R}^d be defined in a neighborhood $I \times U \subset \mathbb{R} \times \mathbb{R}^d$ of the initial value. Then an *initial value problem* (IVP) is defined as follows: find a function $z(t)$, such that:

$$\begin{cases} z'(t) = f(t, z(t)) \\ z(t_0) = z_0 \end{cases} \quad (1.8)$$

Definition 5 (Local solution). We call a continuously differentiable function $z(t)$ with $z(t_0) = z_0$ a *local solution* of the IVP, if there exists a neighborhood J of the point in time t_0 in which $z(t)$ and $f(t, z(t))$ are defined and if the equation $z'(t) = f(t, z(t))$ holds for all $t \in J$.

Definition 6 (Linear ODE). A differential equation is said to be *linear* if F can be written as a linear combination of the derivatives of $z(t)$:

$$z^{(n)}(t) = \sum_{i=1}^{n-1} a_i(t) z^{(i)}(t) + r(t) \quad (1.9)$$

with $a_i(t)$ and $r(t)$ continuous functions in t . If $r(t) = 0$ then we call the linear differential equation *homogeneous* otherwise we call it *inhomogeneous*.

1.1.1 Well-posedness of the Initial Value Problems

Definition 7. A mathematical problem is called well-posed if the following *Hadamard conditions* are satisfied:

- A solution exists.
- The solution is unique.
- The solution is continuously dependent on the data.

In the specific case of IVPs, the third condition is often dropped and substituted with the so-called Lipschitz continuity, which is more a quantitative condition.

Definition 8. The function $f(t, z)$ satisfies on its domain $D = I \times \Omega \subset \mathbb{R} \times \mathbb{R}^d$ an uniformly continuous Lipschitz condition if it is Lipschitz continuous with regard to z . In other words it exists a positive constant L , such that:

$$\forall t \in I; x, z \in \Omega : |f(t, x) - f(t, z)| \leq L|x - z| \quad (1.10)$$

It satisfies a local Lipschitz condition if the same holds true for all compact subsets of D .

As we will see the *Peano's existence theorem* proofs that if f is a continuous map than there exists at least a solution to the ODE. It must be pointed out that the Peano theorem does not proof that the solution is unique.

Theorem 1 (Peano's existence theorem). *Let the function $f(t, z)$ be continuous on the closed set*

$$\overline{D} = \{(t, z) \in \mathbb{R} \times \mathbb{R}^d \mid |t - t_0| \leq \alpha, |z - z_0| \leq \beta\} \quad (1.11)$$

where $\alpha, \beta > 0$. Then there exists a solution $z \in C^1(I)$ on the interval $I = [t_0 - T, t_0 + T]$ with:

$$T = \min\left(\alpha, \frac{\beta}{M}\right), M = \max_{(t, u) \in \overline{D}} |f(t, z)| \quad (1.12)$$

Theorem 2 (Peano's continuation theorem). *Let the assumptions of Peano's existence theorem hold. Then, the solution can be extended to an interval $I_m = [t_-, t_+]$ such that the points $(t_-, z(t_-))$ and $(t_+, z(t_+))$ are on the boundary of D . Neither the values of t , nor of $z(t)$ need to be bounded as long as f remains bounded.*

To proof the uniqueness of solution we employ the aforementioned *Lipschitz continuity* together with the *Grönwall's lemma*.

Lemma 2.2 (Grönwall's theorem). *Let be $w(t)$, $a(t)$ and $b(t)$ be nonnegative, integrable functions, such that $a(t)w(t)$ is integrable. Furthermore, let $b(t)$ be monotonically nondecreasing and let $w(t)$ satisfy the integral inequality:*

$$w(t) \leq b(t) + \int_{t_0}^t a(s)w(s)ds, \quad t \geq t_0 \quad (1.13)$$

Then, for almost all $t \geq t_0$ there holds:

$$w(t) \leq b(t) \exp\left(\int_{t_0}^t a(s)ds\right) \quad (1.14)$$

The existence and uniqueness of solution can be now proofed thanks to the *Picard-Lindelöf theorem* (also called *Picard's existence theorem*, *Cauchy-Lipschitz theorem*, or *existence and uniqueness theorem*):

Theorem 3 (Picard-Lindelöf theorem). *Let $f(t, z)$ be continuous on a cylinder $D = \{(t, z) \in \mathbb{R} \times \mathbb{R}^d \mid |t - t_0| \leq a, |z - z_0| \leq b\}$. Let f be bounded such that there is a constant $M = \max_D |f|$ and satisfy the Lipschitz condition with constant L on D . Then the IVP:*

$$\begin{cases} z'(t) = f(t, z(t)) \\ z(t_0) = z_0 \end{cases} \quad (1.15)$$

is uniquely solvable on the interval $I = [t_0 - T, t_0 + T]$ where $T = \min\{a, b/M\}$.

1.2 Numerical Solutions of Initial Value Problems

In most cases an analytical solution to IVPs cannot be found or it is simply impractical (complex integrals appear in the solution). Thus, a set of numerical schemes for solving IVPs are now presented. It must be pointed out that all the schemes can be naturally extended to systems of differential equations. Moreover, since higher order differential equations can be rewritten as a system of first order differential equations, we will only concentrate on numerical methods for these last ones.

Definition 9 (Time step). On a time interval $I = [t_0, t_0 + T]$, we define a partitioning in n subintervals, also known as *time steps*. The time steps $I_k = [t_{k-1}, t_k]$ have the step size $h_k = t_k - t_{k-1}$. A partitioning in n time steps implies $t_n = T$. The term k -th time step is used for both the interval I_k and for the point in time t_k , but it should always be clear through context which one is meant. Very often, we will consider evenly spaced time steps, in which case we denote the step size by h and $h_k = h$ for all k .

Numerical methods can be subdivided into two main categories: *one-step methods* and *multi-step methods*. In the first category the step z_{k+1} is determined only by z_k , where as in the second category the the step z_{k+1} is determined also by z_{k-1}, \dots, z_{k-n}

Definition 10 (Explicit method). An *explicit method* (forward integration) is a method which, given z_0 at t_0 computes a sequence of approximations z_1, \dots, z_n to the solution of an IVP in the time steps t_1, \dots, t_n using an update formula of the form:

$$z_k = z_{k-1} + h_k F(t_{k-1}, z_{k-1}, \dots, z_{k-p}) \quad (1.16)$$

where F is called *increment function* and p is the order is the numerical multi-step method. If $p = 1$ we call the method as one-step method.

Definition 11 (Implicit method). An *implicit method* (backward integration) is a method which, given z_0 at t_0 computes a sequence of approximations z_1, \dots, z_n to the solution of an IVP in the time steps t_1, \dots, t_n using an update formula of the form:

$$z_k = z_{k-1} + h_k F(t_k, z_k, \dots, z_{k-p}) \quad (1.17)$$

Notice that the increment function F depends on z_k and the previous equation must be solved for z_k .

1.2.1 Euler Method

The first and simplest on-step numerical method is the Euler method. There are many ways of deriving this method. To derive this method let us consider the following IVP:

$$\begin{cases} z'(t) = f(t, z) \\ z(t_0) = z_0 \end{cases} \quad (1.18)$$

The truncated Taylor series of the solution $z(t)$ centered in t_0 , that is:

$$z(t_0 + h) = z_0 + h z_1 + o(h^2) \quad (1.19)$$

where $z_1(t = t_0) = z'(t = t_0)$ and since $z(t)$ is considered to be the solution is the IVPs $z_1(t = t_0) = f(t_0, z_0)$. We get an approximation of $z(t_0 + h)$, namely:

$$(1.20)$$

Notice that the difference between $\tilde{z}(t_0 + h)$ and $z(t_0 + h)$ is proportional to h^2 .

If we consider a generic interval $I = [t_0, t_0 + T]$, we can generalize the equation (1.20) to the k -th step and repeat it on each k -th subinterval $I_k = [t_{k-1}, t_k]$. Thus, in general the k -th step of the *explicit Euler method* can be written as:

$$z(t_k) = z_{k-1} + h f(t_{k-1}, z_{k-1}), \quad k = 1, \dots, n \quad (1.21)$$

whereas the the k -th step of the *implicit Euler method* can be written as:

$$z(t_k) = z_{k-1} + hf(t_k, z_k), \quad k = 1, \dots, T/h \quad (1.22)$$

It can be prooved that both explicit Euler method and implicit Euler method converges to the exact solution as $h \rightarrow 0$ and the error at any time $t \in I = [t_0, t_0 + T]$ can be bounded by Ch , where C is a positive constant. Since C is proportional to e^{LT} , where the Lipschitz constant L of f may be very large. Such problem, which are commonly referred to be *stiff*, are characterized by a large changes in at very different time scales and high sensitivity to changes in the initial condition. If we apply the Euler method to a stiff problems it turns out that the explicit method is not able to properly find the solution due to its *stability*. On the other hand implicit Euler method works large number of applications with a rate of converge of $o(h)$.

1.2.2 Runge-Kutta Methods

Even if the Euler method works fine for most of applications, if the dimension d is large, the end time T is large, the error tolerance ε is small or more importantly the ODE is stiff we might want to improve the solution method. To pursue these improvements we can use the so-called *Runge-Kutta methods*.

Explicit Runge-Kutta Methods

The generic *explicit Runge-Kutta method* is a one-step method with the representation:

$$\begin{aligned} g_i &= z_k + h \sum_{j=1}^{i-1} a_{ij} k_j \\ k_i &= f(hc_i, g_i) \\ g_{k+1} &= z_k + h \sum_{i=1}^s b_i k_i \end{aligned} \quad (1.23)$$

where $i = 1, \dots, s$. In this method the values hc_i are the quadrature points on the interval $[0, h]$. The values k_i are approximations to function values of the integrand in these points and the values g_i constitute approximations to the solution $z(hc_i)$ in the quadrature points. This method uses s intermediate values and is thus called an s -stage method.

Definition 12 (Butcher tableau). It is customary to write Runge-Kutta methods in the form of a *Butcher tableau*, containing only the coefficients of equation 1.23 in the following matrix form:

$$\begin{array}{c|cccccc} 0 & 0 & 0 & 0 & 0 & 0 \\ c_2 & a_{21} & 0 & 0 & 0 & 0 \\ c_3 & a_{31} & a_{32} & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & 0 & 0 \\ c_s & a_{s1} & a_{s2} & \dots & a_{s,s-1} & 0 \\ \hline & b_1 & b_2 & \dots & b_{s-1} & b_s \end{array} \quad (1.24)$$

Explicit Euler A more generalized formulation of the explicit Euler method.

$$\begin{array}{c|c} 0 & 0 \\ \hline & 1 \end{array} \quad (1.25)$$

Midpoint method The midpoint method is a variation of the explicit Euler method, also known as Collatz method. It is a second-order method with two stages.

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 1/2 & 1/2 & 0 \\ \hline & 0 & 1 \end{array} \quad (1.26)$$

Heun's method The Heun's method is also known as the explicit trapezoid rule since it improves the standard explicit Euler method. It is a second-order method with two stages.

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 1 & 1 & 0 \\ \hline & 1/2 & 1/2 \end{array} \quad (1.27)$$

Ralston's method It is a second-order method with two stages and a minimum local error bound.

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 2/3 & 2/3 & 0 \\ \hline & 1/4 & 3/4 \end{array} \quad (1.28)$$

Generic second-order method

$$\begin{array}{c|cc} 0 & 0 & 0 \\ \alpha & \alpha & 0 \\ \hline & 1 - \frac{1}{2\alpha} & \frac{1}{2\alpha} \end{array} \quad (1.29)$$

Runge-Kutta's third-order method (RK3)

$$\begin{array}{c|ccc} 0 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 \\ 1 & -1 & 2 & 0 \\ \hline & 1/6 & 2/3 & 1/6 \end{array} \quad (1.30)$$

Heun's third-order method

$$\begin{array}{c|ccc} 0 & 0 & 0 & 0 \\ 1/3 & 1/3 & 0 & 0 \\ 1/3 & 0 & 2/3 & 0 \\ \hline & 1/4 & 0 & 3/4 \end{array} \quad (1.31)$$

Ralston's third-order method

$$\begin{array}{c|ccc}
0 & 0 & 0 & 0 \\
1/2 & 1/2 & 0 & 0 \\
3/4 & 0 & 3/4 & 0 \\
\hline
& 2/9 & 1/3 & 4/9
\end{array} \tag{1.32}$$

Strong Stability Preserving Runge-Kutta (SSPRK3)

$$\begin{array}{c|ccc}
0 & 0 & 0 & 0 \\
1 & 1 & 0 & 0 \\
1/2 & 1/4 & 1/4 & 0 \\
\hline
& 1/6 & 1/6 & 2/3
\end{array} \tag{1.33}$$

Generic third-order method With $\alpha \neq 0, 2/3, 1$:

$$\begin{array}{c|ccc}
0 & 0 & 0 & 0 \\
\alpha & \alpha & 0 & 0 \\
1 & 1 + \frac{1-\alpha}{\alpha(3\alpha-2)} & -\frac{1-\alpha}{\alpha(3\alpha-2)} & 0 \\
\hline
& \frac{1}{2} - \frac{1}{6\alpha} & \frac{1}{6\alpha(1-\alpha)} & \frac{2-3\alpha}{6\alpha(1-\alpha)}
\end{array} \tag{1.34}$$

Classical Runge-Kutta's method (RK4) The classical fourth-order Runge-Kutta:

$$\begin{array}{c|cccc}
0 & 0 & 0 & 0 & 0 \\
1/2 & 1/2 & 0 & 0 & 0 \\
1/2 & 0 & 1/2 & 0 & 0 \\
1 & 0 & 0 & 1 & 0 \\
\hline
& 1/6 & 2/6 & 2/6 & 1/6
\end{array} \tag{1.35}$$

Ralston's fourth-order method It is a second-order method with two stages and a minimum truncation error.

$$\begin{array}{c|ccccc}
0 & 0 & 0 & 0 & 0 \\
.4 & .4 & 0 & 0 & 0 \\
.45573725 & .29697761 & .15875964 & 0 & 0 \\
1 & .21810040 & -3.05096516 & 3.83286476 & 0 \\
\hline
& .17476028 & -.55148066 & 1.20553560 & .17118478
\end{array} \tag{1.36}$$

3/8-rule fourth-order method It is a second-order method with two stages and a minimum truncation error.

$$\begin{array}{c|cccc}
0 & 0 & 0 & 0 & 0 \\
1/3 & 1/3 & 0 & 0 & 0 \\
2/3 & -1/3 & 1 & 0 & 0 \\
1 & 1 & -1 & 1 & 0 \\
\hline
& 1/8 & 3/8 & 3/8 & 1/8
\end{array} \tag{1.37}$$

Implicit Runge-Kutta Methods

All the presented method are *explicit Runge-Kutta methods*. If we deal with *stiff* equations the *implicit Runge-Kutta method* should be preferred.

The generic *implicit Runge-Kutta method* is a one-step method with the representation:

$$\begin{aligned}
g_i &= z_k + h \sum_{j=1}^s a_{ij} k_j \\
k_i &= f(hc_i, g_i) \\
g_{k+1} &= z_k + h \sum_{i=1}^s b_i k_i
\end{aligned} \tag{1.38}$$

The butcher tableau for (1.38) will be then of the form:

$$\begin{array}{c|cccc}
c_1 & a_{11} & a_{12} & \dots & a_{1s} \\
c_2 & a_{21} & a_{22} & \dots & a_{2s} \\
\vdots & \vdots & \vdots & \ddots & \vdots \\
c_s & a_{s1} & a_{s2} & \dots & a_{ss} \\
\hline
& b_1 & b_2 & \dots & b_s
\end{array} \tag{1.39}$$

Implicit Euler method A more generalized formulation of the implicit Euler method.

$$\begin{array}{c|c}
1 & 1 \\
\hline
& 1
\end{array} \tag{1.40}$$

Implicit midpoint method The implicit midpoint method is a variation of the explicit Euler method, also known as Collatz method. It is a second-order method with two stages.

$$\begin{array}{c|c}
1/2 & 1/2 \\
\hline
& 1
\end{array} \tag{1.41}$$

Crank-Nicolson method The Crank-Nicolson method corresponds to the implicit trapezoidal rule.

$$\begin{array}{c|cc}
0 & 0 & 0 \\
1 & 1/2 & 1/2 \\
\hline
& 1/2 & 1/2
\end{array} \tag{1.42}$$

Gauss-Legendre methods These methods are based on the points of Gauss-Legendre quadrature.

Gauss-Legendre second-order method

$$\begin{array}{c|cc}
 \frac{1}{2} - \frac{\sqrt{3}}{6} & \frac{1}{4} & \frac{1}{4} - \frac{\sqrt{3}}{6} \\
 \frac{1}{2} + \frac{\sqrt{3}}{6} & \frac{1}{4} + \frac{\sqrt{3}}{6} & \frac{1}{4} \\
 \hline
 & \frac{1}{2} & \frac{1}{2} \\
 & \frac{1}{2} + \frac{\sqrt{3}}{2} & \frac{1}{2} - \frac{\sqrt{3}}{2}
 \end{array} \quad (1.43)$$

Gauss-Legendre sixth-order method

$$\begin{array}{c|ccc}
 \frac{1}{2} - \frac{\sqrt{15}}{10} & \frac{5}{36} & \frac{2}{9} - \frac{\sqrt{15}}{15} & \frac{5}{36} - \frac{\sqrt{15}}{30} \\
 \frac{1}{2} & \frac{5}{36} + \frac{\sqrt{15}}{24} & \frac{2}{9} & \frac{5}{36} - \frac{\sqrt{15}}{24} \\
 \frac{1}{2} + \frac{\sqrt{15}}{10} & \frac{5}{36} + \frac{\sqrt{15}}{30} & \frac{2}{9} + \frac{\sqrt{15}}{15} & \frac{5}{36} \\
 \hline
 & \frac{5}{18} & \frac{4}{9} & \frac{5}{18} \\
 & -\frac{5}{6} & \frac{8}{3} & -\frac{5}{6}
 \end{array} \quad (1.44)$$

Lobatto methods

Lobatto IIIA fourth-order method

$$\begin{array}{c|ccc}
 0 & 0 & 0 & 0 \\
 1/2 & 5/24 & 1/3 & -1/24 \\
 1 & 1/6 & 2/3 & 1/6 \\
 \hline
 & 1/6 & 2/3 & 1/6
 \end{array} \quad (1.45)$$

Lobatto IIIB fourth-order method

$$\begin{array}{c|ccc}
 0 & 1/6 & -1/6 & 0 \\
 1/2 & 1/6 & 1/3 & 0 \\
 1 & 1/6 & 5/6 & 0 \\
 \hline
 & 1/6 & 2/3 & 1/6
 \end{array} \quad (1.46)$$

Lobatto IIIC fourth-order method

$$\begin{array}{c|ccc}
 0 & 1/6 & -1/3 & 1/6 \\
 1/2 & 1/6 & 5/12 & -1/12 \\
 1 & 1/6 & 2/3 & 1/6 \\
 \hline
 & 1/6 & 2/3 & 1/6
 \end{array} \quad (1.47)$$

Lobatto IIIC* fourth-order method

0	0	0	0
1/2	1/4	1/4	0
1	0	1	0
	1/6	2/3	1/6

(1.48)

Radau methods**Radau IA fifth-order method**

0	$\frac{1}{9}$	$\frac{-1 - \sqrt{6}}{18}$	$\frac{-1 + \sqrt{6}}{18}$
$\frac{3}{5} - \frac{\sqrt{6}}{10}$	$\frac{1}{9}$	$\frac{11}{45} + \frac{7\sqrt{6}}{360}$	$\frac{11}{45} - \frac{43\sqrt{6}}{360}$
$\frac{3}{5} + \frac{\sqrt{6}}{10}$	$\frac{1}{9}$	$\frac{11}{45} + \frac{43\sqrt{6}}{360}$	$\frac{11}{45} - \frac{7\sqrt{6}}{360}$
	$\frac{1}{9}$	$\frac{4}{9} + \frac{\sqrt{6}}{36}$	$\frac{4}{9} - \frac{\sqrt{6}}{36}$

(1.49)

Radau II fifth-order method

$\frac{2}{5} - \frac{\sqrt{6}}{10}$	$\frac{11}{45} - \frac{7\sqrt{6}}{360}$	$\frac{37}{225} - \frac{169\sqrt{6}}{1800}$	$-\frac{2}{225} + \frac{\sqrt{6}}{75}$
$\frac{2}{5} + \frac{\sqrt{6}}{10}$	$\frac{37}{225} + \frac{169\sqrt{6}}{1800}$	$\frac{11}{45} + \frac{7\sqrt{6}}{360}$	$-\frac{2}{225} - \frac{\sqrt{6}}{75}$
1	$\frac{4}{9} - \frac{\sqrt{6}}{36}$	$\frac{4}{9} + \frac{\sqrt{6}}{36}$	$\frac{1}{9}$
	$\frac{4}{9} - \frac{\sqrt{6}}{36}$	$\frac{4}{9} + \frac{\sqrt{6}}{36}$	$\frac{1}{9}$

(1.50)

Linear Differential Algebraic Equations

2.1 Linear Differential Algebraic Equations

In this chapter linear implicit differential equations or DAE s of the form

$$\mathbf{E}\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{f}(t) \quad (2.1)$$

are studied. In the simplest case, $\mathbf{E}, \mathbf{A} \in \mathbb{R}^{n \times n}$ are square constant matrices and $\mathbf{f}(t) : \mathbb{R} \rightarrow \mathbb{R}^n$ is some external perturbation independent of $\mathbf{x}(t)$. Such equations occur for example by linearization of autonomous non-linear problems with respect to constant (or critical) solutions, here $\mathbf{f}(t)$ is a perturbation that pushes away the solution from the nominal trajectory. Other cases are modelling of linear electrical circuits, simple mechanical systems or, in general, linear systems with additional linear algebraic constraints. The first question that we need to answer, regards existence and uniqueness of a solution $\mathbf{x}(t) : \mathbb{R} \rightarrow \mathbb{R}^n$. Of course the only interesting case is when \mathbf{E} is singular (not invertible). Otherwise it is sufficient to invert \mathbf{E} and the DAE s can be written as a ODE s as follows,

$$\dot{\mathbf{x}}(t) = \mathbf{E}^{-1}\mathbf{A}\mathbf{x}(t) + \mathbf{E}^{-1}\mathbf{f}(t). \quad (2.2)$$

For eq. (3.2) it is sufficient to assume $\mathbf{f}(t)$ piecewise continuous¹ and existence and uniqueness are ensured, see for example [4]. Conversely when matrix \mathbf{E} is singular, it is not a priori clear if the solution exists or not. A singular matrix \mathbf{E} means that ODE s are mixed with linear algebraic constraints. At this point is important consider some regularity conditions for the pair (\mathbf{E}, \mathbf{A}) , as shown in the following example.

Example 1. Consider the DAE s with $\mathbf{x}(t) = (x_1(t); x_2(t)) \in \mathbb{R}^2$, $\mathbf{f}(t) = 0, \forall t \in [0, \infty)$ and matrices \mathbf{E}, \mathbf{A} given by

$$\mathbf{E} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{A} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}. \quad (2.3)$$

This DAE s corresponds to the scalar differential equation $\dot{x}_1(t) = x_2(t)$. Fixing an initial condition $\mathbf{x}_0 = (0; 0) \in \mathbb{R}^2$ a solution is provided by

$$\mathbf{x}(t) = \int_0^t \begin{bmatrix} x_2(\tau) \\ \alpha(\tau) \end{bmatrix} d\tau, \quad (2.4)$$

where $\alpha(\tau)$ represents an arbitrary function that satisfies the initial condition. For example let us consider $\alpha(t) = \sin(t)$ which is compatible with the initial condition ($\sin(0) = 0$),

¹Weaker conditions for $\mathbf{f}(t)$ can be used, for example it is sufficient that $\mathbf{f}(t)$ is Lebesgue measurable.

then the solution $\mathbf{x}(t)$ is

$$\mathbf{x}(t) = \begin{bmatrix} -\sin(t) \\ -\cos(t) \end{bmatrix}. \quad (2.5)$$

Nevertheless other choices are possible, i.e. $\alpha(t) = t$, then the solution is

$$\mathbf{x}(t) = \begin{bmatrix} t^3/6 \\ t^2/2 \end{bmatrix}. \quad (2.6)$$

It is therefore clear that the solution is not unique and that we must assume some regularity assumptions for the pair (\mathbf{E}, \mathbf{A}) . \circ

In order to clarify the existence and uniqueness of a solution for eq. (3.1) we need the following definition.

Definition 13. Matrix pair (\mathbf{E}, \mathbf{A}) is said to be a *regular pencil* if there exist a scalar $\lambda_0 \in \mathbb{R}$, such that

$$\det(\mathbf{A} - \lambda_0 \mathbf{E}) \neq 0. \quad (2.7)$$

Example 2. The pair of matrices $\mathbf{E}, \mathbf{A} \in \mathbb{R}^{3 \times 3}$ defined below is a non regular (singular) pencil.

$$\mathbf{E} = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 0 & 0 \\ 0 & -2 & 0 \end{bmatrix}, \quad \mathbf{A} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad (2.8)$$

Since computing the determinant of the pencil is identically zero for every value of $\lambda_0 \in \mathbb{R}$.

$$\det(\mathbf{A} - \lambda_0 \mathbf{E}) = \det \left(\begin{bmatrix} 1 - \lambda_0 & 0 & 0 \\ 2\lambda_0 & 1 & 0 \\ 0 & 2\lambda_0 & 0 \end{bmatrix} \right) = 0 \quad (2.9)$$

\circ

The key idea is that if pair (\mathbf{E}, \mathbf{A}) is a regular pencil, then under an appropriate change of coordinates it is possible to simplify the problem. The idea is made more precise by the following theorem.

Theorem 4 (Weistrass canonical form). *Matrix pair $\mathbf{E}, \mathbf{A} \in \mathbb{R}^{n \times n}$ is a regular pencil if, and only if, there exist invertible matrices $\mathbf{P} \in \mathbb{R}^{n \times n}$ and $\mathbf{Q} \in \mathbb{R}^{n \times n}$ such that*

$$(\mathbf{PEQ}, \mathbf{PAQ}) = \left(\begin{bmatrix} \mathbf{N} & 0 \\ 0 & \mathbf{I} \end{bmatrix}, \begin{bmatrix} \mathbf{I} & 0 \\ 0 & \mathbf{J} \end{bmatrix} \right), \quad (2.10)$$

where \mathbf{N} is an upper triangular nilpotent matrix; \mathbf{J} is a matrix in Jordan canonical form and \mathbf{I} is the identity matrix of suitable dimensions.

Just for the convenience of the reader we recall the definition of nilpotent matrix.

Definition 14. A square matrix \mathbf{N} is said to be nilpotent if there exists an integer $\nu \in \mathbb{N}$ such that $\mathbf{N}^\nu = 0$.

Remark 1. The smallest ν that satisfies $\mathbf{N}^\nu = 0$ is usually called *degree* of \mathbf{N} . Notice also that any upper triangular matrix with zeros along the main diagonal is nilpotent. \lrcorner

Example 3. Matrix $\mathbf{N} \in \mathbb{R}^{3 \times 3}$ as follows is nilpotent with degree $\nu = 3$.

$$\mathbf{N} = \begin{bmatrix} 0 & 1 & 3 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \end{bmatrix}, \quad \mathbf{N}^2 = \begin{bmatrix} 0 & 0 & 2 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad \mathbf{N}^3 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad (2.11)$$

◦

The decoupling form in eq. (3.10) is called *Weierstrass canonical form*, for further details see [5]. Moreover the two matrices have the following structure;

$$\mathbf{N} = \begin{bmatrix} \mathbf{N}_1 & 0 & \dots & 0 \\ 0 & \mathbf{N}_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{N}_{n_1} \end{bmatrix}, \quad \mathbf{J} = \begin{bmatrix} \mathbf{J}_1 & 0 & \dots & 0 \\ 0 & \mathbf{J}_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{J}_{n_2} \end{bmatrix}. \quad (2.12)$$

Here \mathbf{N}_i with $i = 1, \dots, n_1$ and \mathbf{J}_i with $i = 1, \dots, n_2$ are Jordan blocks. The blocks \mathbf{N}_i are associated to zero eigenvalues of matrix \mathbf{E} , therefore they appear as follows

$$\mathbf{N}_i = \begin{bmatrix} 0 & 1 & & & \\ & 0 & 1 & & \\ & & \ddots & \ddots & \\ & & & \ddots & 1 \\ & & & & 0 \end{bmatrix}, \quad i = 1, \dots, n_1, \quad (2.13)$$

and the dimension of each block \mathbf{N}_i is equal to the difference between the algebraic and the geometric multiplicity of the corresponding eigenvalue. Conversely blocks \mathbf{J}_i with $i = 1, \dots, n_2$ are classical Jordan blocks with the following structure,

$$\mathbf{J}_k = \begin{bmatrix} \lambda_k & 1 & & & \\ & \lambda_k & 1 & & \\ & & \ddots & \ddots & \\ & & & \ddots & 1 \\ & & & & \lambda_k \end{bmatrix}, \quad k = 1, \dots, n_2. \quad (2.14)$$

We now sketch a proof for the Weierstrass canonical form. The proof will be constructive and based on a sequential multiplication of non-singular matrices. Here we prove only the sufficiency part, for the necessity see [5].

Proof. Let us consider the regular pencil pair (\mathbf{E}, \mathbf{A}) , surely we can write $(\mathbf{E}, \mathbf{A} + \lambda_0 \mathbf{E} - \lambda_0 \mathbf{E})$. Now let us define the matrix $\mathbf{T}_0 = (\mathbf{A} - \lambda_0 \mathbf{E})$, which is non singular by regular pencil assumption. Therefore the quantity $\mathbf{T}_0^{-1}(\mathbf{E}, \mathbf{A} + \lambda_0 \mathbf{E} - \lambda_0 \mathbf{E}) = (\mathbf{T}_0^{-1} \mathbf{E}, \mathbf{I} + \lambda_0 \mathbf{T}_0^{-1} \mathbf{E})$. Now let us call \mathbf{T}_1 the non singular transformation matrix such that $\mathbf{T}_1^{-1} \mathbf{T}_0^{-1} \mathbf{E} \mathbf{T}_1$ is in Jordan canonical form. Then the pair $\mathbf{T}_1^{-1} (\mathbf{T}_0^{-1} \mathbf{E}, \mathbf{I} + \lambda_0 \mathbf{T}_0^{-1} \mathbf{E}) \mathbf{T}_1$ has the following structure,

$$\left(\begin{bmatrix} \tilde{\mathbf{N}} & 0 \\ 0 & \tilde{\mathbf{J}} \end{bmatrix}, \begin{bmatrix} \mathbf{I} + \lambda_0 \tilde{\mathbf{N}} & 0 \\ 0 & \mathbf{I} + \lambda_0 \tilde{\mathbf{J}} \end{bmatrix} \right). \quad (2.15)$$

Here the block $\tilde{\mathbf{N}}$ is associated to the zero eigenvalues of the singular matrix $\mathbf{T}_0^{-1} \mathbf{E}$. Therefore $\tilde{\mathbf{N}}$ is upper triangular with zero elements on the main diagonal. Thanks to this

property the matrix $(\mathbf{I} + \lambda_0 \tilde{\mathbf{N}})$ is instead invertible. Therefore we can consider the non singular transformation matrix \mathbf{T}_2 defined as follows,

$$\mathbf{T}_2 = \begin{bmatrix} (\mathbf{I} + \lambda_0 \tilde{\mathbf{N}}) & 0 \\ 0 & \tilde{\mathbf{J}} \end{bmatrix}. \quad (2.16)$$

Now left multiplying eq. (3.15) by \mathbf{T}_2^{-1} we obtain the following structure,

$$\left(\begin{bmatrix} (\mathbf{I} + \lambda_0 \tilde{\mathbf{N}})^{-1} \tilde{\mathbf{N}} & 0 \\ 0 & \mathbf{I} \end{bmatrix}, \begin{bmatrix} \mathbf{I} & 0 \\ 0 & \tilde{\mathbf{J}}^{-1}(\mathbf{I} + \lambda_0 \tilde{\mathbf{J}}) \end{bmatrix} \right). \quad (2.17)$$

Finally it is sufficient to consider matrices \mathbf{T}_{31} and \mathbf{T}_{32} such that $\mathbf{T}_{31}^{-1}(\mathbf{I} + \lambda_0 \tilde{\mathbf{N}})^{-1} \tilde{\mathbf{N}} \mathbf{T}_{31}$ and $\mathbf{T}_{32}^{-1} \tilde{\mathbf{J}}^{-1}(\mathbf{I} + \lambda_0 \tilde{\mathbf{J}}) \mathbf{T}_{32}$ are in Jordan canonical form and build matrix \mathbf{T}_3 as follows,

$$\mathbf{T}_3 = \begin{bmatrix} \mathbf{T}_{31} & 0 \\ 0 & \mathbf{T}_{32} \end{bmatrix}. \quad (2.18)$$

Thus the final change of coordinates yields,

$$\mathbf{T}_3^{-1} \left(\begin{bmatrix} (\mathbf{I} + \lambda_0 \tilde{\mathbf{N}})^{-1} \tilde{\mathbf{N}} & 0 \\ 0 & \mathbf{I} \end{bmatrix}, \begin{bmatrix} \mathbf{I} & 0 \\ 0 & \tilde{\mathbf{J}}^{-1}(\mathbf{I} + \lambda_0 \tilde{\mathbf{J}}) \end{bmatrix} \right) \mathbf{T}_3 = \left(\begin{bmatrix} \mathbf{N} & 0 \\ 0 & \mathbf{I} \end{bmatrix}, \begin{bmatrix} \mathbf{I} & 0 \\ 0 & \mathbf{J} \end{bmatrix} \right). \quad (2.19)$$

It remains to prove that matrix \mathbf{N} is nilpotent and to do that is sufficient to show that matrix $(\mathbf{I} + \lambda_0 \tilde{\mathbf{N}})^{-1} \tilde{\mathbf{N}}$ is nilpotent as well. First notice that matrices $(\mathbf{I} + \lambda_0 \tilde{\mathbf{N}})$ and $\tilde{\mathbf{N}}$ commute since,

$$(\mathbf{I} + \lambda_0 \tilde{\mathbf{N}}) \tilde{\mathbf{N}} = \tilde{\mathbf{N}} + \lambda_0 \tilde{\mathbf{N}}^2 = \tilde{\mathbf{N}}(\mathbf{I} + \lambda_0 \tilde{\mathbf{N}}). \quad (2.20)$$

Then also $(\mathbf{I} + \lambda_0 \tilde{\mathbf{N}})^{-1}$ and $\tilde{\mathbf{N}}$ commute since,

$$(\mathbf{I} + \lambda_0 \tilde{\mathbf{N}})^{-1} \tilde{\mathbf{N}} = (\mathbf{I} + \lambda_0 \tilde{\mathbf{N}})^{-1} \tilde{\mathbf{N}} (\mathbf{I} + \lambda_0 \tilde{\mathbf{N}}) (\mathbf{I} + \lambda_0 \tilde{\mathbf{N}})^{-1} = \tilde{\mathbf{N}} (\mathbf{I} + \lambda_0 \tilde{\mathbf{N}})^{-1} \quad (2.21)$$

where we used the fact that $\tilde{\mathbf{N}}(\mathbf{I} + \lambda_0 \tilde{\mathbf{N}}) = (\mathbf{I} + \lambda_0 \tilde{\mathbf{N}}) \tilde{\mathbf{N}}$. Now since we know that commutative property for these matrices holds, called ν the index such that $\tilde{\mathbf{N}}^\nu = 0$ the following is true

$$\left[(\mathbf{I} + \lambda_0 \tilde{\mathbf{N}}) \tilde{\mathbf{N}} \right]^\nu = (\mathbf{I} + \lambda_0 \tilde{\mathbf{N}})^\nu \tilde{\mathbf{N}}^\nu = 0, \quad (2.22)$$

and this complete the proof. Just to summarize, we notice that the overall transformation matrices (\mathbf{P}, \mathbf{Q}) are given by $\mathbf{P} = (\mathbf{T}_0 \mathbf{T}_1 \mathbf{T}_2 \mathbf{T}_3)^{-1}$ and $\mathbf{Q} = \mathbf{T}_1 \mathbf{T}_3$. \square

The Weistrass canonical form is quite useful since we can re-write system eq. (3.1) using the transformation matrices (\mathbf{P}, \mathbf{Q}) just introduced, thus

$$\mathbf{P} \mathbf{E} \mathbf{Q} \mathbf{Q}^{-1} \dot{\mathbf{x}}(t) = \mathbf{P} \mathbf{A} \mathbf{Q} \mathbf{Q}^{-1} \mathbf{x}(t) + \mathbf{P} \mathbf{f}(t). \quad (2.23)$$

Now considering the following change of coordinates,

$$\mathbf{z}(t) = \mathbf{Q}^{-1} \mathbf{x}(t), \quad (2.24a)$$

$$\mathbf{g}(t) = \mathbf{P} \mathbf{f}(t), \quad (2.24b)$$

we obtain the following equivalent DAE s

$$\begin{bmatrix} \mathbf{N} & 0 \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \dot{\mathbf{z}}_1(t) \\ \dot{\mathbf{z}}_2(t) \end{bmatrix} = \begin{bmatrix} \mathbf{I} & 0 \\ 0 & \mathbf{J} \end{bmatrix} \begin{bmatrix} \mathbf{z}_1(t) \\ \mathbf{z}_2(t) \end{bmatrix} + \begin{bmatrix} \mathbf{g}_1(t) \\ \mathbf{g}_2(t) \end{bmatrix}. \quad (2.25)$$

Here we partitioned the $\mathbf{z}(t)$ and $\mathbf{g}(t)$ vectors according respectively to the dimensions of \mathbf{J} and \mathbf{N} . Obviously $\mathbf{z}(t) = (\mathbf{z}_1(t); \mathbf{z}_2(t))$ and $\mathbf{g}(t) = (\mathbf{g}_1(t); \mathbf{g}_2(t))$. It is also interesting, notice that the two sets of equations corresponding to $\mathbf{z}_1(t)$ and $\mathbf{z}_2(t)$ are decoupled

$$\mathbf{N}\dot{\mathbf{z}}_1 = \mathbf{z}_1 + \mathbf{g}_1(t), \quad (2.26a)$$

$$\dot{\mathbf{z}}_2 = \mathbf{J}\mathbf{z}_2 + \mathbf{g}_2(t). \quad (2.26b)$$

It is straightforward to see that the equation associated to $\mathbf{z}_2(t)$ is just an ODE then the unique solution from the initial condition $\mathbf{z}_2(0)$ is given by

$$\mathbf{z}_2(t) = e^{\mathbf{J}t} \mathbf{z}_2(0) + \int_0^t e^{\mathbf{J}(t-\tau)} \mathbf{g}_2(\tau) d\tau, \quad \forall t \geq 0, \quad (2.27)$$

(for further details see for example [6]). eq. (3.26a) is instead again a DAE, but since matrix \mathbf{N} is nilpotent we can manipulate it. It is convenient to re-write eq. (3.26a) as

$$\mathbf{z}_1(t) = \mathbf{N}\dot{\mathbf{z}}_1(t) - \mathbf{g}_1(t). \quad (2.28)$$

Now substituting eq. (3.28) inside itself and differentiating yields,

$$\begin{aligned} \mathbf{z}_1(t) &= \mathbf{N} \frac{d}{dt} [\mathbf{N}\dot{\mathbf{z}}_1(t) - \mathbf{g}_1(t)] - \mathbf{g}_1(t) = \mathbf{N}^2 \frac{d^2 \mathbf{z}_1(t)}{dt^2} - \mathbf{N} \frac{d\mathbf{g}_1(t)}{dt} - \mathbf{g}_1(t) \\ &= \mathbf{N}^2 \frac{d^2}{dt^2} [\mathbf{N}\dot{\mathbf{z}}_1(t) - \mathbf{g}_1(t)] - \mathbf{N} \frac{d\mathbf{g}_1(t)}{dt} - \mathbf{g}_1(t) = \mathbf{N}^3 \frac{d^3 \mathbf{z}_1(t)}{dt^3} - \mathbf{N}^2 \frac{d^2 \mathbf{g}_1(t)}{dt^2} - \dots \\ &\vdots \\ &= \mathbf{N}^\nu \frac{d^\nu \mathbf{z}_1(t)}{dt^\nu} - \mathbf{N}^{\nu-1} \frac{d^{\nu-1} \mathbf{g}_1(t)}{dt^{\nu-1}} - \dots - \mathbf{N} \frac{d\mathbf{g}_1(t)}{dt} - \mathbf{g}_1(t), \end{aligned}$$

establishing differential properties of the solutions. Notice that the procedure stops when we reach the degree ν of matrix \mathbf{N} , indeed $\mathbf{N}^\nu = 0$ and we may find a unique solution expressed as a finite summation of increasing derivatives of $\mathbf{g}(t)$ as follows:

$$\mathbf{z}_1(t) = - \sum_{k=0}^{\nu-1} \mathbf{N}^k \frac{d^k \mathbf{g}_1(t)}{dt^k}. \quad (2.30)$$

Remark 2. Sometimes index ν is also called *Kronecker index*. Notice also that ν is the number of times that we differentiated function $\mathbf{g}_1(t)$ in order to obtain the solution $\mathbf{z}_1(t)$. \lrcorner

Finally we are ready to express the unique solution to the original problem in eq. (3.1). It is sufficient to invert the change of coordinates in eq. (3.24a) and we obtain

$$\mathbf{x}(t) = \mathbf{Q}\mathbf{z}(t) = \mathbf{Q} \begin{bmatrix} - \sum_{k=0}^{\nu-1} \mathbf{N}^k \frac{d^k \mathbf{g}_1(t)}{dt^k} \\ e^{\mathbf{J}t} \mathbf{z}_2(0) + \int_0^t e^{\mathbf{J}(t-\tau)} \mathbf{g}_2(\tau) d\tau \end{bmatrix} \quad (2.31)$$

Remark 3. Notice that the initial condition $\mathbf{z}_1(0)$ cannot be chosen freely, but has to satisfy the following,

$$\mathbf{z}_1(0) = - \sum_{k=0}^{\nu-1} \mathbf{N}^k \frac{d^k \mathbf{g}_1(0)}{dt^k}. \quad (2.32)$$

Therefore also the initial condition $\mathbf{x}_0 = \mathbf{x}(t=0)$ for eq. (3.1) cannot be chosen completely freely. This intuitively expresses the fact that also the initial condition has to be consistent with constraints expressed by \mathbf{E} . \lrcorner

Summarizing what we discussed we have the following,

Theorem 5. Problem in eq. (3.1) with initial condition $\mathbf{x}_0 = \mathbf{x}(t=0)$ has a solution if the following conditions holds:

- function $\mathbf{f}(t)$ is piece-wise continuous,
- the pair (\mathbf{E}, \mathbf{A}) is a regular pencil,
- the initial condition \mathbf{x}_0 satisfies $\mathbf{x}_0 = \mathbf{Q} \begin{bmatrix} -\sum_{k=0}^{\nu-1} \mathbf{N}^k \frac{d^k \mathbf{g}_1(0)}{dt^k} \\ \mathbf{z}_2(0) \end{bmatrix}$.

Moreover the corresponding solutions is given by eq. (3.31).

2.2 Review of Gaussian Elimination

The Weirstrass canonical form is very useful from the theoretical point of view, however is computationally demanding, because it requires to go through several Jordan canonical forms. A numerical more efficient algorithm to reduce a DAE to an ODE is presented in the next section. The idea is to use the celebrated Gauss's elimination method combined with differentiation in order to perform the desired reduction. Before doing so, we want recall some elementary concepts and definitions useful for the purpose.

Definition 15. Given a matrix $\mathbf{M} \in \mathbb{R}^{n \times m}$ we call elementary operations the followings,

- permutation of two rows (resp. columns),
- linear combinations of rows (resp. columns).

These operations can be also regarded as matrix multiplications.

Definition 16. A non-singular matrix $\mathbf{P}(i, j) \in \{0, 1\}^{n \times n}$ (resp. $\mathbf{P}(i, j) \in \{0, 1\}^{m \times m}$) is said to be a rows (resp. columns) permutation matrix if it is obtained by the identity matrix $\mathbf{I} \in \mathbb{R}^{n \times n}$ (resp. $\mathbf{I} \in \mathbb{R}^{m \times m}$) exchanging the column i with column j .

Permutation matrices have the following remarkable properties,

- $\mathbf{P}(i, j)^{-1} = \mathbf{P}(i, j)^\top$, therefore $\mathbf{P}(i, j)$ is an orthogonal matrix.
- $\mathbf{P}(i, j)^2 = \mathbf{I}$, therefore $\mathbf{P}(i, j)$ is an involution matrix.
- $\mathbf{P}(i, j)\mathbf{M}$ exchanges row i with row j .
- $\mathbf{M}\mathbf{P}(i, j)$ exchanges column i with column j .

Example 4. Consider the following permutation matrix $\mathbf{P}(1, 2) \in \{1, 0\}^{3 \times 3}$ and a matrix $\mathbf{M} \in \mathbb{R}^{3 \times 3}$ then $\mathbf{P}(1, 2)\mathbf{M}$ exchanges rows (1, 2) of matrix \mathbf{M} .

$$\mathbf{P}(1, 2)\mathbf{M} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} = \begin{bmatrix} 4 & 5 & 6 \\ 1 & 2 & 3 \\ 7 & 8 & 9 \end{bmatrix} \quad (2.33)$$

In the same way $\mathbf{M}\mathbf{P}(1, 2)$ exchanges columns (1, 2) of matrix \mathbf{M} ,

$$\mathbf{M}\mathbf{P}(1, 2) = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 1 & 3 \\ 5 & 4 & 6 \\ 8 & 7 & 9 \end{bmatrix} \quad (2.34)$$

◦

In analogous way it is possible to define a matrix that performs linear combinations of rows (resp. columns). These matrices are a special class of Frobenius matrices and usually are called elementary matrices.

Definition 17. A matrix $\mathbf{F}(i, j, \alpha) \in \mathbb{R}^{n \times n}$ (resp. $\mathbf{F}(i, j, \alpha) \in \mathbb{R}^{m \times m}$) is said to be a rows (resp. columns) elementary matrix if it has the following structure,

$$\mathbf{F}(i, j, \alpha) = \begin{bmatrix} 1 & \dots & 0 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \dots & 1 & & 0 & & 0 \\ \vdots & & \vdots & \ddots & \vdots & & \vdots \\ 0 & \dots & \alpha & \dots & 1 & \dots & 0 \\ \vdots & & \vdots & & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 & \dots & 1 \end{bmatrix} \quad (2.35)$$

with the (i, j) element equal to α .

Remark 4. An elementary matrix $\mathbf{F}(i, j, \alpha) \in \mathbb{R}^{n \times n}$ (resp. $\mathbf{F}(i, j, \alpha) \in \mathbb{R}^{m \times m}$) is equal to the identity matrix of dimension n (resp. m) for all the entries except the (i, j) position which is equal to α . \square

Combination matrices have the following remarkable properties:

- $\mathbf{F}(i, j, \alpha)\mathbf{M}$ add the row j multiplied by α to row i ,
- $\mathbf{M}\mathbf{F}(i, j, \alpha)$ add the column i multiplied by α to column j .

Theorem 6. Given a matrix $\mathbf{M} \in \mathbb{R}^{n \times m}$ with $n \leq m$ and $\text{rank } \mathbf{M} = r \leq n$, there exist a pair of non singular matrices, $\mathbf{L} \in \mathbb{R}^{n \times n}$ and $\mathbf{U} \in \mathbb{R}^{m \times m}$ such that,

$$\mathbf{LMU} = \begin{bmatrix} \mathbf{\Lambda} & 0 \\ 0 & 0 \end{bmatrix}, \quad (2.36)$$

with $\mathbf{\Lambda} \in \mathbb{R}^{r \times r}$ a non singular diagonal matrix. Blocks of zeros have dimensions compatible with the partitioning.

Proof. The proof is quite algorithmic and based on sequentially Gauss's eliminations. We define $\mathbf{M}^{(k)}$ the matrix \mathbf{M} at the k iteration. Parenthesis are used to stress out that $\mathbf{M}^{(k)}$ is not the matrix \mathbf{M} to the power k . So let us initialize the algorithm with $\mathbf{M}^{(1)} = \mathbf{M}$ and $k = 1$. Now the idea is to look for the first non zero element inside the minor $\mathbf{M}(k \dots n, k \dots m)$; obviously at the first iteration we are exploring all the entries. Suppose to move first along rows and next along columns; then unless the matrix is completely full of zeros we will find a pair of indices such that $m_{pq}^{(k)} \neq 0$. Then the idea is to move this non zero element in the (k, k) position. This can be done with a rows and columns permutations as follows,

$$\overline{\mathbf{M}}^{(k)} = \mathbf{P}(k, p)\mathbf{M}^{(k)}\mathbf{P}(q, k). \quad (2.37)$$

Notice that the $\overline{m}_{k,k}^{(k)} = m_{p,q}^{(k)}$, roughly speaking the non zero element in position (p, q) of matrix $\mathbf{M}^{(k)}$ is moved to the (k, k) position of matrix $\overline{\mathbf{M}}^{(k)}$. Now we can use this non zero element to set to zero all the element along the k column and the k row. This can be done using a composition of linear combination matrices. Specifically the matrix

$\bar{\mathbf{L}}^{(k)} \in \mathbb{R}^{n \times n}$ defined as follows is designed to set to zero all the elements in the indices range $(k \dots n, k \dots k)$,

$$\begin{aligned} \bar{\mathbf{L}}^{(k)} &= \mathbf{F} \left(k+1, k, -\frac{\bar{m}_{k+1,k}}{\bar{m}_{k,k}} \right) \dots \mathbf{F} \left(n-1, k, -\frac{\bar{m}_{n-1,k}}{\bar{m}_{k,k}} \right) \mathbf{F} \left(n, k, -\frac{\bar{m}_{n,k}}{\bar{m}_{k,k}} \right) \\ &= \prod_{h=k+1}^n \mathbf{F} \left(h, k, -\frac{\bar{m}_{h,k}}{\bar{m}_{k,k}} \right). \end{aligned} \quad (2.38)$$

The same operation can be performed along columns, in this case matrix $\bar{\mathbf{U}}^{(k)}$ performs linear combinations to set to zero all the element in the indices range $(k \dots k, k \dots m)$,

$$\begin{aligned} \bar{\mathbf{U}}^{(k)} &= \mathbf{F} \left(k, k+1, -\frac{\bar{m}_{k,k+1}}{\bar{m}_{k,k}} \right) \dots \mathbf{F} \left(k, n-1, -\frac{\bar{m}_{k,n-1}}{\bar{m}_{k,k}} \right) \mathbf{F} \left(k, n, -\frac{\bar{m}_{k,n}}{\bar{m}_{k,k}} \right) \\ &= \prod_{h=k+1}^m \mathbf{F} \left(k, h, -\frac{\bar{m}_{k,h}}{\bar{m}_{k,k}} \right). \end{aligned} \quad (2.39)$$

Now just for notational convenience let us define the following quantities,

$$\mathbf{L}^{(k)} = \bar{\mathbf{L}}^{(k)} \mathbf{P}(k, p), \quad (2.40a)$$

$$\mathbf{U}^{(k)} = \mathbf{P}(q, k) \bar{\mathbf{U}}^{(k)}. \quad (2.40b)$$

Then using simultaneously the two matrices $\mathbf{L}^{(k)} \in \mathbb{R}^{n \times n}$ and $\mathbf{U}^{(k)} \in \mathbb{R}^{m \times m}$ we can set to zero all the elements along the k row and the k column, obviously except the (k, k) position. Then the matrix \mathbf{M} at next iteration is given by,

$$\mathbf{M}^{(k+1)} = \mathbf{L}^{(k)} \mathbf{M}^{(k)} \mathbf{U}^{(k)}. \quad (2.41)$$

Is now easy to see that iterating the process; after exactly r iterations, matrix $\mathbf{M}^{(r+1)} \in \mathbb{R}^{n \times m}$ has exactly the following structure,

$$\mathbf{M}^{(r+1)} = \begin{bmatrix} \mathbf{\Lambda} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}. \quad (2.42)$$

Moreover the matrices $\mathbf{L} \in \mathbb{R}^{n \times n}$ and $\mathbf{U} \in \mathbb{R}^{m \times m}$ that realize that structure are given by,

$$\mathbf{L} = \mathbf{L}^{(r)} \mathbf{L}^{(r-1)} \dots \mathbf{L}^{(1)}, \quad (2.43a)$$

$$\mathbf{U} = \mathbf{U}^{(1)} \mathbf{U}^{(2)} \dots \mathbf{U}^{(r)}. \quad (2.43b)$$

Notice also that \mathbf{L} and \mathbf{U} are non singular since are obtained through product of non singular matrices. \square

To make things clear we propose the following example of the algorithm.

Example 5. Let us consider the matrix $\mathbf{M} \in \mathbb{R}^{3 \times 4}$ with rank $\mathbf{M} = r = 2$ defined as follows,

$$\mathbf{M} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 4 & 2 & 0 \end{bmatrix}. \quad (2.44)$$

First of all we set $\mathbf{M}^{(1)} = \mathbf{M}$ and $k = 1$ and we look for the first non zero element starting from position $(k, k) = (1, 1)$ and moving along the first column and then along the second and so forth. At the first iteration the first non zero element is in the position $(p, q) = (2, 2)$,

so the following permutation matrices $\mathbf{P}(1, 2) \in \{0, 1\}^{3 \times 3}$ and $\mathbf{P}(2, 1) \in \{0, 1\}^{4 \times 4}$ provide a way to move it in the $(1, 1)$ position.

$$\overline{\mathbf{M}}^{(1)} = \mathbf{P}(1, 2)\mathbf{M}^{(1)}\mathbf{P}(2, 1) = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \mathbf{M}^{(1)} \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 4 & 0 & 2 & 0 \end{bmatrix} \quad (2.45)$$

Now we need to set to zero all the non zero elements in column and row 1, and the non only non zero element is $\overline{m}_{3,1}^{(1)} = 4$. Therefore the row combination matrix $\mathbf{C}(3, 1, -4/2) \in \mathbb{R}^{3 \times 3}$ reaches the goal,

$$\mathbf{M}^{(2)} = \mathbf{C}(3, 1, -4/2)\overline{\mathbf{M}}^{(1)} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -2 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 4 & 0 & 2 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 2 & 0 \end{bmatrix} \quad (2.46)$$

Now $k = 2$ and we perform again a search for the first non zero element in $\mathbf{M}^{(2)}$. The answer is $(p, q) = (2, 3)$ so we can use again permutation matrices, $\mathbf{P}(2, 2) = \mathbf{I} \in \{0, 1\}^{3 \times 3}$ and $\mathbf{P}(3, 1) \in \{0, 1\}^{4 \times 4}$.

$$\overline{\mathbf{M}}^{(2)} = \mathbf{M}^{(2)}\mathbf{P}(3, 1) = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 2 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 2 & 0 & 0 \end{bmatrix} \quad (2.47)$$

Finally we can use the rows combination matrix $\mathbf{C}(3, 2, -2/1) \in \mathbb{R}^{3 \times 3}$ to complete the procedure,

$$\mathbf{M}^{(3)} = \mathbf{C}(3, 2, -2/1)\overline{\mathbf{M}}^{(2)} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -2 & 1 \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 2 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad (2.48)$$

Here obviously the diagonal block $\mathbf{\Lambda} \in \mathbb{R}^{2 \times 2} = \text{diag}(2, 1)$. Notice also that is now clear that the rank of matrix \mathbf{M} is equal $r = 2$. \circ

2.3 Gauss like method to compute DAEs' index

We are finally ready to discuss why theorem 9 is useful in computing the index of a DAE, and how it can be used to transform a generic linear DAE into an ODE. So let us consider a generic linear DAE in the form shown in eq. (3.1), and we notice that can be also written as,

$$\begin{bmatrix} \mathbf{E} & -\mathbf{A} \end{bmatrix} \begin{bmatrix} \dot{\mathbf{x}} \\ \mathbf{x} \end{bmatrix} = \mathbf{f}. \quad (2.49)$$

Just for notational brevity we dropped the time dependency. Here matrix $\begin{bmatrix} \mathbf{E} & -\mathbf{A} \end{bmatrix} \in \mathbb{R}^{n \times 2n}$. Let us call

$$\text{rank } \mathbf{E} =: r \leq n. \quad (2.50)$$

Then let us consider a pair of matrices $\mathbf{L}^{(0)} \in \mathbb{R}^{n \times n}$ and $\mathbf{U}^{(0)} \in \mathbb{R}^{n \times n}$ such that,

$$\mathbf{U}^{(0)} \mathbf{E} \mathbf{L}^{(0)} = \mathbf{E}^{(0)} = \begin{bmatrix} \mathbf{\Lambda}^{(0)} & 0 \\ 0 & 0 \end{bmatrix}.$$

As done in the previous section we used the superscript $^{(0)}$ to indicate the iteration number; it will be clear soon that this number is also the Kronecker index of the DAE. Notice that at this first iteration $\mathbf{\Lambda}^{(0)} \in \mathbb{R}^{r \times r}$. If $r = n$ then we done since the DAE has index $\nu = 0$ and matrix \mathbf{E} was full rank and invertible. If $r < n$ it is convenient introduce the following change of coordinates,

$$\begin{bmatrix} \dot{\mathbf{x}}^{(0)} \\ \mathbf{x}^{(0)} \end{bmatrix} = \begin{bmatrix} \mathbf{U}^{(0)} & 0 \\ 0 & \mathbf{I} \end{bmatrix}^{-1} \begin{bmatrix} \dot{\mathbf{x}} \\ \mathbf{x} \end{bmatrix}, \quad (2.51)$$

and also the matrix $\mathbf{A}^{(0)}$ and vector $\mathbf{f}^{(0)}$ as follows,

$$\begin{aligned} \mathbf{A}^{(0)} &:= -\mathbf{L}^{(0)} \mathbf{A}, \\ \mathbf{f}^{(0)} &:= \mathbf{L}^{(0)} \mathbf{f}. \end{aligned}$$

Let us consider to multiply on the left-hand side eq. (3.49) by $\mathbf{L}^{(0)}$. Then we obtain the following equivalent form,

$$\begin{aligned} \mathbf{L}^{(0)} \begin{bmatrix} \mathbf{E} & -\mathbf{A} \end{bmatrix} \begin{bmatrix} \mathbf{U}^{(0)} & 0 \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{U}^{(0)} & 0 \\ 0 & \mathbf{I} \end{bmatrix}^{-1} \begin{bmatrix} \dot{\mathbf{x}} \\ \mathbf{x} \end{bmatrix} &= \begin{bmatrix} \mathbf{L}^{(0)} \mathbf{E} \mathbf{U}^{(0)} & \mathbf{L}^{(0)} \mathbf{A} \end{bmatrix} \begin{bmatrix} \dot{\mathbf{x}}^{(0)} \\ \mathbf{x}^{(0)} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{E}^{(0)} & \mathbf{A}^{(0)} \end{bmatrix} \begin{bmatrix} \dot{\mathbf{x}}^{(0)} \\ \mathbf{x}^{(0)} \end{bmatrix} = \mathbf{f}^{(0)}(t) \end{aligned}$$

Now notice that matrix $\begin{bmatrix} \mathbf{E}^{(0)} & \mathbf{A}^{(0)} \end{bmatrix} \in \mathbb{R}^{n \times 2n}$ has the following special structure,

$$\begin{bmatrix} \mathbf{\Lambda}^{(0)} & 0 & \mathbf{A}_{11}^{(0)} & \mathbf{A}_{12}^{(0)} \\ 0 & 0 & \mathbf{A}_{21}^{(0)} & \mathbf{A}_{22}^{(0)} \end{bmatrix} \begin{bmatrix} \dot{\mathbf{x}}_1^{(0)} \\ \dot{\mathbf{x}}_2^{(0)} \\ \mathbf{x}_1^{(0)} \\ \mathbf{x}_2^{(0)} \end{bmatrix} = \begin{bmatrix} \mathbf{f}_1^{(0)} \\ \mathbf{f}_2^{(0)} \end{bmatrix}. \quad (2.53)$$

Obviously eq. (3.53) can be also written as a couple of matrix equations as follows,

$$\mathbf{\Lambda}^{(0)} \dot{\mathbf{x}}_1^{(0)} + \mathbf{A}_{11}^{(0)} \mathbf{x}_1^{(0)} + \mathbf{A}_{12}^{(0)} \mathbf{x}_2^{(0)} = \mathbf{f}_1^{(0)}, \quad (2.54a)$$

$$\mathbf{A}_{21}^{(0)} \mathbf{x}_1^{(0)} + \mathbf{A}_{22}^{(0)} \mathbf{x}_2^{(0)} = \mathbf{f}_2^{(0)}, \quad (2.54b)$$

where $\dot{\mathbf{x}}_1^{(0)}, \mathbf{x}_1^{(0)} \in \mathbb{R}^r$, $\dot{\mathbf{x}}_2^{(0)}, \mathbf{x}_2^{(0)} \in \mathbb{R}^{n-r}$ and $\mathbf{x}^{(0)} = [\mathbf{x}_1^{(0)}; \mathbf{x}_2^{(0)}]$. Notice that eq. (3.54a) is a ODE since $\mathbf{\Lambda}^{(0)}$ is invertible, on the contrary eq. (3.54b) is a pure algebraic equation. Therefore in these coordinates the differential and the algebraic components of the DAE are decoupled. We now differentiate with respect to time eq. (3.54b) and we obtain the following,

$$\mathbf{A}_{21}^{(0)} \dot{\mathbf{x}}_1^{(0)} + \mathbf{A}_{22}^{(0)} \dot{\mathbf{x}}_2^{(0)} = \dot{\mathbf{f}}_2^{(0)}. \quad (2.55)$$

Writing now eq. (3.55) with eq. (3.54a) in matrix form we obtain the following,

$$\begin{bmatrix} \mathbf{\Lambda}^{(0)} & 0 & \mathbf{A}_{11}^{(0)} & \mathbf{A}_{12}^{(0)} \\ \mathbf{A}_{21}^{(0)} & \mathbf{A}_{22}^{(0)} & 0 & 0 \end{bmatrix} \begin{bmatrix} \dot{\mathbf{x}}_1^{(0)} \\ \dot{\mathbf{x}}_2^{(0)} \\ \mathbf{x}_1^{(0)} \\ \mathbf{x}_2^{(0)} \end{bmatrix} = \begin{bmatrix} \mathbf{v}_1^{(0)} \\ \dot{\mathbf{f}}_2^{(0)} \end{bmatrix}. \quad (2.56)$$

As we can observe the differentiation w.r.t time of the algebraic constraints in the matrix representation produces just a shift to the left of the two blocks $\mathbf{A}_{21}^{(0)}$, $\mathbf{A}_{22}^{(0)}$. Notice also this shift increases the rank of the left block as it is clearly stated below,

$$\text{rank} \begin{bmatrix} \mathbf{\Lambda}^{(0)} & 0 \\ 0 & 0 \end{bmatrix} < \text{rank} \begin{bmatrix} \mathbf{\Lambda}^{(0)} & 0 \\ \mathbf{A}_{21}^{(0)} & \mathbf{A}_{22}^{(0)} \end{bmatrix} \quad (2.57)$$

Now we can apply again theorem 9 and find a pair of matrices $\mathbf{L}^{(1)} \in \mathbb{R}^{n \times n}$ and $\mathbf{U}^{(1)} \in \mathbb{R}^{n \times n}$ such that,

$$\mathbf{U}^{(1)} \begin{bmatrix} \mathbf{\Lambda}^{(0)} & 0 \\ \mathbf{A}_{21}^{(0)} & \mathbf{A}_{22}^{(0)} \end{bmatrix} \mathbf{L}^{(1)} = \mathbf{E}^{(1)} = \begin{bmatrix} \mathbf{\Lambda}^{(1)} & 0 \\ 0 & 0 \end{bmatrix},$$

and perform another change of coordinates as follows,

$$\begin{aligned} \begin{bmatrix} \dot{\mathbf{x}}^{(1)} \\ \mathbf{x}^{(1)} \end{bmatrix} &= \begin{bmatrix} \mathbf{U}^{(1)} & 0 \\ 0 & \mathbf{I} \end{bmatrix}^{-1} \begin{bmatrix} \dot{\mathbf{x}}^{(0)} \\ \mathbf{x}^{(0)} \end{bmatrix}, \\ \mathbf{A}^{(1)} &= \mathbf{L}^{(1)} \mathbf{A}^{(0)}, \\ \mathbf{f}^{(1)} &= \mathbf{L}^{(0)} \begin{bmatrix} \mathbf{f}_1^{(0)} \\ \mathbf{f}_2^{(0)} \end{bmatrix}. \end{aligned}$$

Please notice that the definition of $\mathbf{f}^{(1)}$ contains part of the field $\mathbf{f}^{(0)}$ but also some derivatives of it. Now we can left multiply eq. (3.56) by $\mathbf{L}^{(1)}$ and use the change of coordinates defined in eq. (3.58). So finally we reach the following form,

$$\begin{bmatrix} \mathbf{E}^{(1)} & \mathbf{A}^{(1)} \end{bmatrix} \begin{bmatrix} \dot{\mathbf{x}}^{(1)} \\ \mathbf{x}^{(1)} \end{bmatrix} = \mathbf{f}^{(1)}. \quad (2.59)$$

Here matrix $\begin{bmatrix} \mathbf{E}^{(1)} & \mathbf{A}^{(1)} \end{bmatrix}$ as seen before has the following structure,

$$\begin{bmatrix} \mathbf{E}^{(1)} & \mathbf{A}^{(1)} \end{bmatrix} = \begin{bmatrix} \mathbf{\Lambda}^{(1)} & 0 & \mathbf{A}_{11}^{(1)} & \mathbf{A}_{12}^{(1)} \\ 0 & 0 & \mathbf{A}_{21}^{(1)} & \mathbf{A}_{22}^{(1)} \end{bmatrix}. \quad (2.60)$$

It is important to point out that the dimensions of matrix $\mathbf{\Lambda}^{(1)}$ are greater than the ones of matrix $\mathbf{\Lambda}^{(0)}$, formally

$$\dim \mathbf{\Lambda}^{(1)} > \dim \mathbf{\Lambda}^{(0)}. \quad (2.61)$$

It is therefore clear that at each iteration (derivation) of the constraints the dimensions of $\mathbf{\Lambda}^{(k)}$, for $k = 1, \dots, \nu$ increases until $\dim \mathbf{\Lambda}^{(\nu)} = n$ and the matrix $\begin{bmatrix} \mathbf{E}^{(\nu)} & \mathbf{A}^{(\nu)} \end{bmatrix}$ coincide with $\begin{bmatrix} \mathbf{\Lambda}^{(\nu)} & \mathbf{A}^{(\nu)} \end{bmatrix}$. At this point the matrix $\mathbf{\Lambda}^{(\nu)} \in \mathbb{R}^{n \times n}$ is full rank and then the DAE can be easily transformed into an ODE.

$$\mathbf{\Lambda}^{(\nu)} \dot{\mathbf{x}}^{(\nu)} + \mathbf{A}^{(\nu)} \mathbf{x}^{(\nu)} = \mathbf{f}^{(\nu)}. \quad (2.62)$$

The minimum integer such that $\dim \mathbf{\Lambda}^{(\nu)} = n$ is the Kronecker index (or differential index) of the DAE.

Example 6. Consider the following set of linear differential equation,

$$\dot{x}_1(t) + \dot{x}_2(t) + \dot{x}_3(t) = 2x_1(t) + x_3(t) + t, \quad (2.63a)$$

$$\dot{x}_1(t) + \dot{x}_2(t) + \dot{x}_3(t) = x_1(t) + x_2(t) + t^3, \quad (2.63b)$$

$$0 = x_1(t) + x_2(t) + x_3(t) + \sin t, \quad (2.63c)$$

that we can write in a more compact matrix form as follows,

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \dot{x}_1(t) \\ \dot{x}_2(t) \\ \dot{x}_3(t) \end{bmatrix} = \begin{bmatrix} 2 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} x_1(t) \\ x_2(t) \\ x_3(t) \end{bmatrix} + \begin{bmatrix} t \\ t^3 \\ \sin t \end{bmatrix}. \quad (2.64)$$

Which corresponds to the usual form $E\dot{x}(t) = Ax(t) + f(t)$ obvious meaning for the symbols. Now consider to form the matrix $\begin{bmatrix} E & A & f(t) \end{bmatrix} \in \mathbb{R}^{n \times 2n+1}$ in such a way that,

$$\begin{bmatrix} E & A & f(t) \end{bmatrix} \begin{bmatrix} \dot{x}(t) \\ -x(t) \\ -1 \end{bmatrix} = 0 \quad (2.65)$$

First it is necessary to check if the pair (E, A) is a regular pencil,

$$\det(A - \lambda_0 E) = \det \begin{bmatrix} 2 - \lambda_0 & -\lambda_0 & 1 - \lambda_0 \\ 1 - \lambda_0 & 1 - \lambda_0 & -\lambda_0 \\ 1 & 1 & 1 \end{bmatrix} = 2 \quad (2.66)$$

the determinant is different from zero, so the pair (E, A) is a regular pencil.

Now it is easy to manipulate the following equations

$$\begin{bmatrix} 1 & 1 & 1 & 2 & 0 & 1 & t \\ 1 & 1 & 1 & 1 & 1 & 0 & t^3 \\ 0 & 0 & 0 & 1 & 1 & 1 & \sin t \end{bmatrix} \quad (2.67)$$

◦

Differential Algebraic Equations for Mechanics

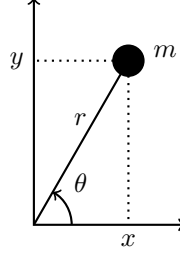
3.1 Overview of Lagrangian Mechanics

Most of the following section is inspired by [1] and [7], which are advanced classical textbooks on classical mechanics. Mechanics deals with the dynamics of particles, rigid bodies, continuous media and field theories such as electromagnetism, gravity, etc. The foundations were laid by Newton in the 17th century and since then many generalizations and extensions have been proposed. To date Mechanics has two main point of view, *Lagrangian mechanics* and *Hamiltonian mechanics*. In these notes we will briefly recall some concepts of Lagrangian Mechanics, for insights see [7]. The Lagrangian formulation of mechanics is based on the observation that there are *variational principles* behind the fundamental laws of force balance, as given by Newton's law $\mathbf{F} = m\mathbf{a}$. This principle is called *least action*, or more accurately, the principle of *stationary action*. This principle can be successfully applied to every mechanical system to obtain the equations of motion. The idea is fairly general, since it can be used to derive Newtonian, Lagrangian, and Hamiltonian equations of motion. However, here we will focus only on mechanical systems with a finite number of *degrees of freedom*. For those systems the *configuration* can be identified through a vector of *generalized coordinates* $\mathbf{q} := (q_1; q_2; \dots; q_n) \in \mathbb{R}^n$, where n is the *minimum* number of *degrees of freedom*. Notice that the choice of vector \mathbf{q} is not unique. After choosing a ordered set of coordinate \mathbf{q} , it is always possible to express *kinetic* and the *potential* energy as a function of it. Thus let us define $T(\mathbf{q}, \dot{\mathbf{q}})$ as the *kinetic energy* and $U(\mathbf{q}, t)$ as the *potential energy*. Roughly we can think T, U as functions whose domain of definition is simply the direct product of Euclidean spaces, i.e., $T : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$ and $U : \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}$. However this is not completely true, since vector \mathbf{q} may represent a coordinatization of more abstract surfaces, namely *manifolds*; however (locally) we can think at Euclidean spaces. Once the *energy* of the system has been defined we can introduce the so called *Lagrangian function* $L(\mathbf{q}, \dot{\mathbf{q}}, t)$, whose domain is $L : \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}$. In classical Mechanics the Lagrangian is defined as the difference between *kinetic* and *potential* energy, formally

$$L(\mathbf{q}, \dot{\mathbf{q}}, t) := T(\mathbf{q}, \dot{\mathbf{q}}) - U(\mathbf{q}, t). \quad (3.1)$$

Example 7. Let us consider a very simple example, a free mass point in the plane subject to gravity in the y downward direction, see example 1. The kinetic and the potential energy for the system are,

$$T = \frac{1}{2}m(\dot{x}^2 + \dot{y}^2), \quad U = mgy,$$

Figure 3.1: A free point of mass m in the plane with gravity.

thus taking $\mathbf{q} = (x; y) \in \mathbb{R}^2$ as a vector of generalized coordinates the Lagrangian function is of the form

$$L = \frac{1}{2}m\dot{\mathbf{q}}^\top \dot{\mathbf{q}} - mge^\top \mathbf{q},$$

where $e_2 = (0; 1) \in \mathbb{R}^2$. Nevertheless Cartesian coordinates is not the only possible choice, thus considering *polar* coordinates the following relations hold

$$\begin{aligned} x &= r \cos \theta, & y &= r \sin \theta, \\ \dot{x} &= \dot{r} \cos \theta - r \sin \theta \dot{\theta}, & \dot{y} &= \dot{r} \sin \theta + r \cos \theta \dot{\theta}, \end{aligned}$$

where $r \in \mathbb{R}_{\geq 0}$ is the radius and $\theta \in [0, 2\pi)$ the angular coordinate. Then taking some steps the kinetic and potential energies can be expressed as,

$$T = \frac{1}{2}m(\dot{r}^2 + r^2\dot{\theta}^2), \quad U = mgr \sin \theta.$$

Thus taking $\tilde{\mathbf{q}} := (\tilde{q}_1; \tilde{q}_2) = (r; \theta) \in \mathbb{R}^2$ as generalized coordinates the Lagrangian is

$$L = \frac{1}{2}m(\dot{\tilde{q}}_1^2 + \tilde{q}_1^2 \dot{\tilde{q}}_2^2) - mg\tilde{q}_1 \sin \tilde{q}_2. \quad (3.3)$$

This simple example shows that the choice of generalized coordinates is not unique but it can be performed according to a criterion of simplicity or utility. As we will see shortly the Lagrangian formulation provides the correct equation of motion despite the chosen reference system. \circ

A very related concept in classical mechanics is the *action* $S \in \mathbb{R}$, this is defined as the integral of the Lagrangian function performed between two instants of time $t_1, t_2 \in \mathbb{R}$, i.e.,

$$S := \int_{t_1}^{t_2} L(\mathbf{q}, \dot{\mathbf{q}}, t) dt. \quad (3.4)$$

Technically speaking the action is a *functional*, i.e., a map that associate to each element of a vector space (usually a space of curves) a scalar. Thus a functional takes curves as input and returns a scalar.

Remark 5. Notice that $S = S(\mathbf{q})$ should be thought as a function of $\mathbf{q}(t)$ over the whole time interval $[t_1, t_2]$. From this point onwards we will use the notation $\mathbf{q}(\cdot)$ to denote the function \mathbf{q} , while $\mathbf{q}(t)$ will be denote the *value of function* \mathbf{q} at the point t . Clarified this it is very important to stress that S does not depends on $q(t)$ but from $q(\cdot)$ and it is just a number. Thus for every choice of the “shape” of curve $\mathbf{q}(\cdot)$, the action S assumes a different (scalar) value. \lrcorner

Example 8. Suppose that $\mathbf{q} = (q_1; q_2) \in \mathbb{R}^2$, $t_1 = 0$, $t_2 = 1$ and $L = q_1^2 + q_2^2$, thus the action is the following,

$$S = \int_0^1 q_1^2 + q_2^2 dt.$$

Take now $\mathbf{q} = (2; t)$, calculating the value for the action results

$$S = \int_0^1 2^2 + t^2 dt = \left[4t + \frac{t^3}{3} \right]_0^1 = \frac{13}{3}.$$

◦

Similarly as seen in analysis courses you may ask yourself if there exist a curve that minimizes or maximizes the action. The answer has created a new branch of mathematics, i.e., the *calculus of variations*. Calculus of variations concerns the *extremal* of functionals whose domain is an infinite-dimensional space: the space of curves. As we will see shortly the *stationary action principle* simply says that every mechanical system evolves along a curve $\mathbf{q}(\cdot)$ that is an extremal of the action functional S . In order to state more formally the stationary action principle, let us consider an initial time t_1 and an initial configuration associated to it, $\mathbf{q}_1 = \mathbf{q}(t_1)$, consider also a final time t_2 and a final configuration $\mathbf{q}_2 = \mathbf{q}(t_2)$.

Theorem 7 (Stationary action). *Motions of mechanical systems from configuration \mathbf{q}_1 to configuration \mathbf{q}_2 , with $t \in [t_1, t_2]$, coincide with extremals curves of the functional in eq. (2.4). Or equivalently, for small variations of the curve $\mathbf{q}(\cdot)$ the action S is stationary, i.e.,*

$$\delta S = 0. \quad (3.5)$$

Where δ denotes a small perturbation of the functional S and sometimes it is also called *variational operator*.

Example 9. The δ operator can be successfully applied also to ordinary calculus. Let us consider the function $f(x, y) = \cos(x) + y^2$, then a variation of f results

$$\delta f = \frac{\partial f}{\partial x} \delta x + \frac{\partial f}{\partial y} \delta y = \begin{pmatrix} \cos x & y \end{pmatrix} \begin{pmatrix} \delta x \\ \delta y \end{pmatrix}.$$

In this setting stationary points are the ones that satisfies $\delta f = 0$ since for *first order* variation $(\delta x, \delta y)$ the corresponding variation of f is zero. (The surface is locally flat and the variation on f is visible only at the second order). Therefore the stationary points are $(\frac{\pi}{2} + k\pi, 0)$ for $k = 1, 2, \dots$ ◦

Now let us develop the consequences of this principle. Assuming L regular enough is possible to move the δ operator under the integral,

$$\delta S = \int_{t_1}^{t_2} \delta L(\mathbf{q}, \dot{\mathbf{q}}, t) dt = \int_{t_1}^{t_2} \frac{\partial L}{\partial \mathbf{q}} \delta \mathbf{q} + \frac{\partial L}{\partial \dot{\mathbf{q}}} \delta \dot{\mathbf{q}} dt. \quad (3.6)$$

Remark 6. Notice that variations are performed at *fixed time*, therefore the generic variation $\delta L(\mathbf{q}(t), \dot{\mathbf{q}}(t), t) = \frac{\partial L}{\partial \mathbf{q}} \delta \mathbf{q} + \frac{\partial L}{\partial \dot{\mathbf{q}}} \delta \dot{\mathbf{q}}$ and not $\delta L(\mathbf{q}(t), \dot{\mathbf{q}}(t), t) \neq \frac{\partial L}{\partial \mathbf{q}} \delta \mathbf{q} + \frac{\partial L}{\partial \dot{\mathbf{q}}} \delta \dot{\mathbf{q}} + \frac{\partial L}{\partial t} \delta t$. Moreover we must think at the $\dot{\mathbf{q}}$ variable as independent of \mathbf{q} , since it can be proven that for small variations those quantities are uncorrelated each other. ┘

Now it is useful to consider the following equality that directly follows from integration by part,

$$\frac{\partial L}{\partial \dot{\mathbf{q}}} \delta \dot{\mathbf{q}} = \frac{d}{dt} \left(\frac{\partial L}{\partial \dot{\mathbf{q}}} \delta \mathbf{q} \right) - \frac{d}{dt} \frac{\partial L}{\partial \dot{\mathbf{q}}} \delta \mathbf{q}, \quad (3.7)$$

then substituting eq. (2.7) into eq. (2.6) results,

$$\delta S = \int_{t_1}^{t_2} \frac{d}{dt} \left(\frac{\partial L}{\partial \dot{\mathbf{q}}} \delta \mathbf{q} \right) - \frac{d}{dt} \frac{\partial L}{\partial \dot{\mathbf{q}}} \delta \mathbf{q} + \frac{\partial L}{\partial \mathbf{q}} \delta \mathbf{q} dt \quad (3.8a)$$

$$= \left[\frac{\partial L}{\partial \dot{\mathbf{q}}} \delta \mathbf{q} \right]_{t_1}^{t_2} + \int_{t_1}^{t_2} \left(\frac{\partial L}{\partial \mathbf{q}} - \frac{d}{dt} \frac{\partial L}{\partial \dot{\mathbf{q}}} \right) \delta \mathbf{q} dt \quad (3.8b)$$

$$= \int_{t_1}^{t_2} \left(\frac{\partial L}{\partial \mathbf{q}} - \frac{d}{dt} \frac{\partial L}{\partial \dot{\mathbf{q}}} \right) \delta \mathbf{q} dt = 0 \quad (3.8c)$$

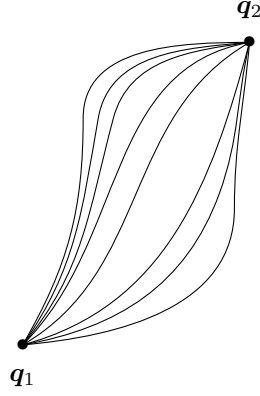


Figure 3.2: A graphical illustration of the possible variations with fixed endpoints $\mathbf{q}_1, \mathbf{q}_2$.

Remark 7. In developing passage (2.8b) we used $\left[\frac{\partial L}{\partial \dot{\mathbf{q}}} \delta \mathbf{q}\right]_{t_1}^{t_2} = 0$, since variations $\delta \mathbf{q}(t_1) = \delta \mathbf{q}(t_2)$ are assumed to be zero because we fixed *a priori* initial and final conditions $\mathbf{q}_1, \mathbf{q}_2$. The concept is graphically illustrated in fig. (2.2), where a set of possible variations is considered but all are vanishing at the endpoints. \lrcorner

We are now ready to further simplify eq. (2.8b) using the *fundamental lemma of calculus of variations*, which adapted for our purposes states the following:

Lemma 7.3 (Fundamental lemma of calculus of variations). *Let S be a differentiable (in variational sense) functional, then if*

$$\int_{t_1}^{t_2} \left(\frac{\partial L}{\partial \mathbf{q}} - \frac{d}{dt} \frac{\partial L}{\partial \dot{\mathbf{q}}} \right) \delta \mathbf{q} dt = 0, \quad (3.9)$$

for any smooth variation $\delta \mathbf{q}$ with $\delta \mathbf{q}_1 = \delta \mathbf{q}_2 = 0$, then,

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{\mathbf{q}}} - \frac{\partial L}{\partial \mathbf{q}} = 0. \quad (3.10)$$

The eq. (2.10) are called *Euler-Lagrange* equations and they provide the equations of motions for mechanical systems. Roughly speaking lemma 4.3 states that since the variation $\delta \mathbf{q}$ is arbitrary, the only way to zeroing the action's variation δS is to satisfy Euler-Lagrange equations. Thus we can collect our results in the following theorem.

Theorem 8. *The curve $q : t \mapsto q(t)$, for $t \in [t_1, t_2]$, is an extremal of the functional S on the space of curves passing through the points $\mathbf{q}(t_1) = \mathbf{q}_1$ and $\mathbf{q}(t_2) = \mathbf{q}_2$, if and only if*

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{\mathbf{q}}} - \frac{\partial L}{\partial \mathbf{q}} = 0. \quad (3.11)$$

The key idea of theorem 5 is that a curve $\mathbf{q}(\cdot)$ in order to be a stationary point for the action S , must satisfy the Euler-Lagrange equations, and therefore according to theorem 4 Euler-Lagrange equations provide the equations of motion for mechanical systems. Notice also that the variational approach is *coordinate free*, meaning that, Euler-Lagrange equations describe the motion regardless the chosen set of coordinates. Indeed the condition for a curve to be an extremal of a functional does not depend on the choice of coordinate systems.

Example 10. Considering the free mass point in example 1 the equations of motions in the Cartesian coordinates are

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{\mathbf{q}}} - \frac{\partial L}{\partial \mathbf{q}} = m\ddot{\mathbf{q}} + mg\mathbf{e}_2 = 0, \quad (3.12)$$

while in polar coordinates,

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{\tilde{q}}_1} - \frac{\partial L}{\partial \tilde{q}_1} = m\ddot{\tilde{q}}_1 + mg \sin \tilde{q}_2, \quad (3.13a)$$

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{\tilde{q}}_2} - \frac{\partial L}{\partial \tilde{q}_2} = m\dot{\tilde{q}}_1^2 \tilde{q}_2 + 2m\tilde{q}_1 \dot{\tilde{q}}_1 \dot{\tilde{q}}_2 + mg\tilde{q}_1 \cos \tilde{q}_2 = 0. \quad (3.13b)$$

Solutions $q(\cdot)$ and $\tilde{q}(\cdot)$ respectively solutions of eq. (2.12) and eq. (2.13) describe the same motion in the plane, just in different coordinate systems. \circ

Nevertheless as presented so far the Euler-Lagrange approach has a big limitation, since L depends only on kinetic and potential energy, we can describe only systems subject to forces that admit a potential function U , i.e., *conservative systems*. Luckily the Euler-Lagrange approach can be extended to *nonconservative forces* using the *D'Alembert's principle*, thus eq. (2.10) is modified as follows

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{\mathbf{q}}}^\top - \frac{\partial L}{\partial \mathbf{q}}^\top = \boldsymbol{\tau}, \quad (3.14)$$

where $\boldsymbol{\tau} \in \mathbb{R}^n$ is a vector of *generalized* non conservative forces.

3.2 Mechanical equations of motion

Although Euler-Lagrange equations are very general is often useful to consider an explicit representation for mechanical systems with some special properties. For this purpose we will assume that the kinetic energy $T(\mathbf{q}, \dot{\mathbf{q}})$ has a very special structure, i.e., is a quadratic form in the generalized velocities $\dot{\mathbf{q}}$. This is not a very restrictive assumption since every mechanical system with a finite number of degrees of freedom owns a quadratic kinetic energy. From this simple assumption a lot of interesting properties can be derived; for a good survey see [10]. An explicit expression for kinetic energy is

$$T(\mathbf{q}, \dot{\mathbf{q}}) = \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n m_{jk}(\mathbf{q}) \dot{q}_j \dot{q}_k = \frac{1}{2} \dot{\mathbf{q}}^\top \mathbf{M}(\mathbf{q}) \dot{\mathbf{q}}, \quad (3.15)$$

where $\mathbf{M}(\mathbf{q}) \in \mathbb{R}^{n \times n}$ is a symmetric positive definite matrix, usually called *mass matrix*. The entries $m_{jk}(\mathbf{q}) \in \mathbf{M}(\mathbf{q})$, for $j, k = 1, \dots, n$ are *generalized masses* and *generalized inertias*.

Remark 8. The matrix $\mathbf{M}(\mathbf{q})$ is symmetric since the quadratic form associated to a matrix only depends on symmetric part. \lrcorner

We are now interested in developing an explicit representation for the equations of motion associated to mechanical systems with a *quadratic* kinetic energy. In order to do so we need to perform some tedious calculations. Let us start from Euler-Lagrange equations in eq. (2.10), the term $\frac{d}{dt} \frac{\partial L}{\partial \dot{q}_i}$ appears as

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{q}_i} = \frac{d}{dt} \frac{\partial T}{\partial \dot{q}_i} = \frac{d}{dt} \frac{\partial}{\partial \dot{q}_i} \left(\frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n m_{jk}(\mathbf{q}) \dot{q}_j \dot{q}_k \right).$$

for $i = 1, \dots, n$. Then exploiting the linearity of summation we can exchange the order between partial derivative and the sum, so that

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{q}_i} = \frac{d}{dt} \left(\frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n m_{jk}(\mathbf{q}) \frac{\partial \dot{q}_j}{\partial \dot{q}_i} \dot{q}_k + \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n m_{jk}(\mathbf{q}) \dot{q}_j \frac{\partial \dot{q}_k}{\partial \dot{q}_i} \right) \quad (3.16a)$$

$$= \frac{d}{dt} \left(\frac{1}{2} \sum_{k=1}^n m_{ik}(\mathbf{q}) \dot{q}_k + \frac{1}{2} \sum_{j=1}^n m_{ji}(\mathbf{q}) \dot{q}_j \right) = \frac{d}{dt} \left(\sum_{j=1}^n m_{ij}(\mathbf{q}) \dot{q}_j \right). \quad (3.16b)$$

In last passage we used the fact that $\frac{\partial \dot{q}_i}{\partial \dot{q}_i} = 1$ for $i = j$ and 0 otherwise, and the same holds for $\frac{\partial \dot{q}_k}{\partial \dot{q}_i}$, moreover we also used the symmetry $m_{ij}(\mathbf{q}) = m_{ji}(\mathbf{q})$ for $i, j = 1, \dots, n$. Finally deriving with respect to time results,

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{q}_i} = \sum_{j=1}^n m_{ij}(\mathbf{q}) \ddot{q}_j + \sum_{j=1}^n \sum_{k=1}^n \frac{\partial m_{ij}(\mathbf{q})}{\partial q_k} \dot{q}_k \dot{q}_j. \quad (3.17)$$

The obtained expression in eq. (2.17) is just a part of the Euler-Lagrange equations; considering then the term $\frac{\partial L}{\partial q_i}$, for $i = 1, \dots, n$, results

$$\frac{\partial L}{\partial q_i} = \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n \frac{\partial m_{jk}(\mathbf{q})}{\partial q_i} \dot{q}_k \dot{q}_j - \frac{\partial U(\mathbf{q}, t)}{\partial q_i}. \quad (3.18)$$

In order to obtain a more compact representation let us define

$$g_i(\mathbf{q}) := \frac{\partial U(\mathbf{q}, t)}{\partial q_i},$$

for $i = 1, \dots, n$ and the associate vector $\mathbf{g}(\mathbf{q}) := (g_1(\mathbf{q}); \dots; g_n(\mathbf{q})) \in \mathbb{R}^n$. Combining together eq. (2.17) and eq. (2.18), the overall equations of motion written componentwise result,

$$\sum_{j=1}^n m_{ij}(\mathbf{q}) \ddot{q}_j + \sum_{j=1}^n \sum_{k=1}^n \frac{\partial m_{ij}(\mathbf{q})}{\partial q_k} \dot{q}_j \dot{q}_k - \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n \frac{\partial m_{jk}(\mathbf{q})}{\partial q_i} \dot{q}_j \dot{q}_k + g_i(\mathbf{q}) = \tau_i, \quad (3.19)$$

for $i = 1, \dots, n$. As it is clear from eq. (2.19) there is a quadratic term in the generalized velocities $\dot{q}_j \dot{q}_k$, therefore it is convenient to introduce the symbol $h_{ijk}(\mathbf{q}) \in \mathbb{R}$, with $i, j, k = 1, \dots, n$, defined as

$$h_{ijk}(\mathbf{q}) := \frac{\partial m_{ij}(\mathbf{q})}{\partial q_k} - \frac{1}{2} \frac{\partial m_{jk}(\mathbf{q})}{\partial q_i}. \quad (3.20)$$

Thus substituting eq. (2.20) into eq. (2.19) the equations of motion result in the more compact representations

$$\sum_{j=1}^n m_{ij}(\mathbf{q}) \ddot{q}_j + \sum_{j=1}^n \sum_{k=1}^n h_{ijk}(\mathbf{q}) \dot{q}_j \dot{q}_k + g_i(\mathbf{q}) = \tau_i, \quad (3.21)$$

for $i = 1, \dots, n$. Now let us analyse eq. (2.21) to provide a physical interpretation of the different terms:

- The coefficients $m_{ii}(\mathbf{q})$, with $i = 1, \dots, n$, represent the generalized masses associated to the i -th degree of freedom.
- The coefficients $m_{ij}(\mathbf{q})$, with $i, j = 1, \dots, n$ but $i \neq j$, represent the *inertia forces* induced by the j -th degree of freedom over the i -th.
- The quadratic terms $h_{ijj} \dot{q}_j^2$, with $j = 1, \dots, n$, represent *centrifugal effect* induced on the i -th degree of freedom by the j -th. Notice that the i -th degree of freedom does not induce any centrifugal effect on itself, i.e., $h_{iii} = 0$.
- The terms $h_{ijj} \dot{q}_j \dot{q}_k$, with $j, k = 1, \dots, n$ but $i \neq j$ represents the *Coriolis* effects induced on the degree of freedom i -th by velocities of degrees of freedom j -th and k -th.
- The terms $g_i(\mathbf{q})$, with $i = 1, \dots, n$, represent the *conservative generalized forces* acting on the i -th degree of freedom.

At this point it is very useful to rewrite the equations of motion (2.21) in a compact matrix notation as follows:

$$\mathbf{M}(\mathbf{q})\ddot{\mathbf{q}} + \mathbf{C}(\mathbf{q}, \dot{\mathbf{q}})\dot{\mathbf{q}} + \mathbf{g}(\mathbf{q}) = \boldsymbol{\tau}, \quad (3.22)$$

where $\mathbf{C}(\mathbf{q}, \dot{\mathbf{q}})$ is a suitable matrix that takes into account the centrifugal and the Coriolis effects. The choice of this matrix is not unique, since it is sufficient that the elements c_{ij} satisfies the relation

$$\sum_{j=1}^n c_{ij}(\mathbf{q})\dot{q}_j = \sum_{j=1}^n \sum_{k=1}^n h_{ijk}(\mathbf{q})\dot{q}_j\dot{q}_k. \quad (3.23)$$

Nevertheless a very common choice is to set the elements $c_{ij}(\mathbf{q}) \in \mathbf{C}(\mathbf{q}, \dot{\mathbf{q}})$ as

$$c_{ij}(\mathbf{q}) = \sum_{k=1}^n c_{ijk}(\mathbf{q})\dot{q}_k, \quad (3.24)$$

with $i, j = 1, \dots, n$. Here the symbol $c_{ijk}(\mathbf{q}) \in \mathbb{R}$ is called *Christoffel symbols of the first type* and it is defined as

$$c_{ijk}(\mathbf{q}) = \frac{1}{2} \left(\frac{\partial m_{ij}(\mathbf{q})}{\partial q_k} + \frac{\partial m_{ik}(\mathbf{q})}{\partial q_j} - \frac{\partial m_{jk}(\mathbf{q})}{\partial q_i} \right). \quad (3.25)$$

Notice that thanks to the symmetry of $\mathbf{M}(\mathbf{q})$ the following remarkable property of indices exchange holds,

$$c_{ijk}(\mathbf{q}) = c_{ikj}(\mathbf{q}). \quad (3.26)$$

Moreover, if the Coriolis-centrifugal matrix $\mathbf{C}(\mathbf{q}, \dot{\mathbf{q}})$ is chosen according to eq. (3.24), then the following matrix

$$\dot{\mathbf{M}}(\mathbf{q}) - 2\mathbf{C}(\mathbf{q}, \dot{\mathbf{q}}), \quad (3.27)$$

is skew symmetric. This property is sometimes useful to synthesize control algorithms for mechanical systems.

3.3 Euler-Lagrange equations with redundant coordinates

What we developed until now holds for systems with a *minimum set of coordinates*. Nevertheless for many real systems is quite hard to write directly kinetic and potential energy using a small set of coordinates. A more common approach is to introduce more generalized coordinates than necessary and then add constraints. This fact provides the motivation to investigate the equations of motions for mechanical systems using a non-minimum set of generalized coordinates. For this purpose let us consider that vector $\mathbf{q} \in \mathbb{R}^n$ *does not* represent a minimum set of coordinates. Therefore in order to correctly represent the system we need to impose some constraints. A set of constrain equations can be formally thought as a vectorial function of the form $\boldsymbol{\phi} : \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}^m$ where m is the number of constrains. The type of constraints equations that we consider are *holomic*, meaning that they depend only on generalized coordinates \mathbf{q} and possibly time t . Holomic constrains can be expressed in the following general form,

$$\boldsymbol{\phi}(\mathbf{q}, t) = 0. \quad (3.28)$$

Remark 9. The total number of degrees of freedom for the system is not $n - m$, unless some regularity for $\boldsymbol{\phi}(\mathbf{q}, t)$ is assumed. Technically if $\forall \mathbf{q} \in \mathbb{R}^n, \forall t \in \mathbb{R}$ the rank $\frac{\partial \boldsymbol{\phi}(\mathbf{q}, t)}{\partial \mathbf{q}} = m$, then constraints are globally linearly independent each other and the system owns exactly $n - m$ degrees of freedom. Trivially we implicitly assumed that $\boldsymbol{\phi}$ is differentiable. \square

Now our aim is to derive the equations of motions for a constrained mechanical systems described by a Lagrangian function $L(\mathbf{q}, \dot{\mathbf{q}}, t)$ and subject to constrain equation $\phi(\mathbf{q}, t) = 0$. In order to do so, we can use the technique of Lagrange multipliers on functional spaces. So let us define the *modified action* $\bar{S} \in \mathbb{R}$ as

$$\bar{S} := \int_{t_1}^{t_2} L(\mathbf{q}, \dot{\mathbf{q}}, t) + \boldsymbol{\lambda}^\top \phi(\mathbf{q}, t) dt. \quad (3.29)$$

Here $\boldsymbol{\lambda} \in \mathbb{R}^m$ is a time dependent vector of *Lagrange multipliers*. The following result is fundamental for our purpose.

Theorem 9 (Stationary action for constrained systems). *Motion of mechanical systems from configuration \mathbf{q}_1 to configuration \mathbf{q}_2 , with $t \in [t_1; t_2]$, coincide with extremals curves of the functional (2.29). Or equivalently, for small variations of the curve $\mathbf{q}(\cdot)$ the modified action \bar{S} is stationary, i.e.,*

$$\delta \bar{S} = 0.$$

Remark 10. Notice that we implicitly assumed that initial and final configuration $\mathbf{q}_1, \mathbf{q}_2$ are consistent with constrain equations, i.e., $\phi(\mathbf{q}_1, t_1) = \phi(\mathbf{q}_2, t_2) = 0$. \lrcorner

Thus we can think to perform variation as done before for systems with a minimal set of coordinates, but there is a difference, since now \bar{S} also depends on the Lagrange multipliers. Therefore the variational operator δ performs variations on $\mathbf{q}, \dot{\mathbf{q}}$ and $\boldsymbol{\lambda}$. The result is the following,

$$\delta \bar{S} = \int_{t_1}^{t_2} \delta L(\mathbf{q}, \dot{\mathbf{q}}, t) + \delta \boldsymbol{\lambda}^\top \phi(\mathbf{q}, t) + \boldsymbol{\lambda}^\top \delta \phi(\mathbf{q}, t) dt \quad (3.30a)$$

$$= \int_{t_1}^{t_2} \left(\frac{\partial L}{\partial \mathbf{q}} - \frac{d}{dt} \frac{\partial L}{\partial \dot{\mathbf{q}}} + \boldsymbol{\lambda}^\top \frac{\partial \phi}{\partial \mathbf{q}} \right) \delta \mathbf{q} + \delta \boldsymbol{\lambda}^\top \phi(\mathbf{q}, t) dt + \left[\frac{\partial L}{\partial \dot{\mathbf{q}}} \delta \mathbf{q} \right]_{t_1}^{t_2} \quad (3.30b)$$

$$= \int_{t_1}^{t_2} \begin{bmatrix} \frac{\partial L}{\partial \mathbf{q}} - \frac{d}{dt} \frac{\partial L}{\partial \dot{\mathbf{q}}} + \boldsymbol{\lambda}^\top \frac{\partial \phi}{\partial \mathbf{q}} & \phi^\top \end{bmatrix} \begin{bmatrix} \delta \mathbf{q} \\ \delta \boldsymbol{\lambda} \end{bmatrix} dt = 0, \quad (3.30c)$$

again we used that endpoints variations $\delta \mathbf{q}_1 = \delta \mathbf{q}_2 = 0$. Finally applying the fundamental lemma of calculus of variations eq. (2.30a), results

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{\mathbf{q}}} - \frac{\partial L}{\partial \mathbf{q}} = \frac{\partial \phi^\top}{\partial \mathbf{q}} \boldsymbol{\lambda}, \quad (3.31a)$$

$$\phi(\mathbf{q}, t) = 0. \quad (3.31b)$$

The equations (2.31) describe the motion of a mechanical system subject to holomic constraints. Note that the term $\frac{\partial \phi^\top}{\partial \mathbf{q}} \boldsymbol{\lambda}$ plays the role of a constrain generalized force which pushes the system to satisfy the constraints $\phi(\mathbf{q}, t) = 0$. Unfortunately problem in eq. (2.31) is completely different from the one in eq. (2.10), since it is a mixture of *differential* and *algebraic* equations, formally a *Differential Algebraic Equation* (DAE). Moreover, apparently eq. (2.31), does not provide any way to determine the value of the Lagrange multipliers. Similarly in presence of external non-conservative generalized forces the equation of motions for a constrained mechanical system are modified as follows,

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{\mathbf{q}}} - \frac{\partial L}{\partial \mathbf{q}} = \frac{\partial \phi^\top}{\partial \mathbf{q}} \boldsymbol{\lambda} + \boldsymbol{\tau}, \quad (3.32a)$$

$$\phi(\mathbf{q}, t) = 0, \quad (3.32b)$$

again with $\boldsymbol{\tau} \in \mathbb{R}^n$ a vector of generalized non-conservative forces.

Example 11. Let us consider a simple mechanical example, the pendulum, see example 5. Pendulum own one-degree of freedom and taking the angular coordinate $\theta \in [0, 2\pi)$ as

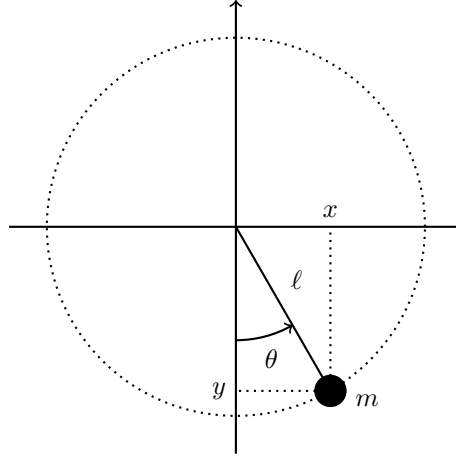


Figure 3.3: A simple pendulum.

generalized coordinate, the following relationships with $(x, y) \in \mathbb{R}^2$ hold,

$$x = \ell \sin \theta, \quad y = -\ell \cos \theta.$$

Thus kinetic and potential energy can be expressed as,

$$T(\theta) = \frac{1}{2}m(\dot{x}^2 + \dot{y}^2) = \frac{1}{2}m\ell^2\dot{\theta}^2, \quad U(\theta) = mgy = -gm\ell \cos \theta,$$

while the Lagrangian for the system results,

$$L = \frac{1}{2}m\ell^2\dot{\theta}^2 + gm\ell \cos \theta.$$

Using the Euler-Lagrange equations in eq. (2.10) yields to

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{\theta}} - \frac{\partial L}{\partial \theta} = m\ell^2\ddot{\theta} + mg\ell \sin \theta = 0,$$

which can be simplified to

$$\ddot{\theta} = -\frac{g}{\ell} \sin \theta.$$

Now suppose we want to describe the pendulum using a set of redundant coordinates, e.g., *natural coordinates*. For this aim let us define the generalized coordinate $\mathbf{q} = (x; y) \in \mathbb{R}^2$. The associated Lagrangian function takes the form,

$$L = \frac{1}{2}m\dot{\mathbf{q}}^\top \dot{\mathbf{q}} - mg\mathbf{e}_2^\top \mathbf{q},$$

where $\mathbf{e}_2 \in \mathbb{R}^2$ is the second vector of the natural basis, roughly speaking the vector $\mathbf{e}_2 = (0; 1) \in \mathbb{R}^2$. Obviously now we need to introduce some constraint equation in order to force mass m to slide over the circle, i.e.,

$$\phi(\mathbf{q}) = x^2 + y^2 - \ell^2 = \mathbf{q}^\top \mathbf{q} - \ell^2 = 0.$$

Using then eq. (2.31) the resultant system of equations is the following,

$$\begin{aligned} m\ddot{\mathbf{q}} + mg\mathbf{e}_2 &= 2\mathbf{q}\lambda, \\ \mathbf{q}^\top \mathbf{q} - \ell^2 &= 0. \end{aligned}$$

◻

Resolution and Numerical Integration of Differential Algebraic Equations

As well pointed out in [8] both from theoretical point of view that the numerical, DAEs are much more challenging than ordinary differential equations (ODEs). Some DAEs can be solved using numerical methods for stiff systems but others cannot and they are very sensitive to the step size and this may cause large errors in the solution and numerical instability. For DAEs coming from constrained mechanical systems a fairly general form is the following,

$$M(\mathbf{q})\ddot{\mathbf{q}} - \frac{\partial \phi^\top}{\partial \dot{\mathbf{q}}} \boldsymbol{\lambda} = \mathbf{n}(\mathbf{q}, \dot{\mathbf{q}}, t), \quad (4.1a)$$

$$\phi(\mathbf{q}, t) = 0 \quad (4.1b)$$

where the term $\mathbf{n}(\mathbf{q}, \dot{\mathbf{q}}, t) \in \mathbb{R}^n$ collects all the contributions coming from Coriolis and centrifugal effects, conservative and non-conservative forces, e.g., it may be definite as follows,

$$\mathbf{n}(\mathbf{q}, \dot{\mathbf{q}}, t) = \boldsymbol{\tau} - \mathbf{C}(\mathbf{q}, \dot{\mathbf{q}})\dot{\mathbf{q}} - \mathbf{g}(\mathbf{q}).$$

It is well known that the equations of motion for a mechanical system are a system of second order (possibly) non-linear differential equations. However for clarity and for sake of generality it is useful to consider a *first order representation*. For this purpose let us define $\mathbf{p} \in \mathbb{R}^n$ as *generalized velocities*, i.e.,

$$\mathbf{p} := \dot{\mathbf{q}}. \quad (4.2)$$

Thus substituting eq. (2.35) into eq. (2.34) results,

$$\dot{\mathbf{q}} = \mathbf{p}, \quad (4.3a)$$

$$M(\mathbf{q})\dot{\mathbf{p}} - \frac{\partial \phi^\top}{\partial \dot{\mathbf{q}}} \boldsymbol{\lambda} = \mathbf{n}(\mathbf{q}, \mathbf{p}, t), \quad (4.3b)$$

$$\phi(\mathbf{q}, t) = 0. \quad (4.3c)$$

Notice that the vector $(\mathbf{q}; \mathbf{p}) \in \mathbb{R}^{2n}$ represents the *state* of a mechanical system, namely the pair *generalized position* and *generalized velocity*.

Now a good question is how to solve the problem in eq. (2.36), the issue is certainly how to determine the value for the Lagrange multiplier $\boldsymbol{\lambda}$. One possible approach is to transform the problem into a system of pure differential equations. In order to do so we need to differentiate the constraint $\phi(\mathbf{q}, t)$, this yields to

$$\frac{d\phi}{dt} = \frac{\partial \phi}{\partial \mathbf{q}} \mathbf{p} + \frac{\partial \phi}{\partial t} = 0. \quad (4.4)$$

Notice that when constraint does not depend on time eq. (2.37) expresses the intuitive idea that generalized velocity \mathbf{p} must be orthogonal to constraint's gradient. This provide also the opportunity to notice that constraint $\phi(\mathbf{q}, t)$ imposes relationships also at the velocity and acceleration level, e.g., eq. (2.37) restricts the set of feasible velocities. The constraints can we can obtain differentiating the constraint are somehow called *hidden constraints*. Finally notice that eq. (2.37) is still algebraic in the \mathbf{p} coordinate, thus in order to obtain a differential relationship let us differentiate again,

$$\frac{d^2\phi(\mathbf{q}, t)}{dt^2} = \frac{d}{dt} \frac{\partial\phi}{\partial\mathbf{q}} \mathbf{p} + \frac{\partial\phi}{\partial\mathbf{q}} \dot{\mathbf{p}} + \frac{\partial^2\phi}{\partial t^2} = 0. \quad (4.5)$$

Notice that now eq. (2.38) does not impose any algebraic constraint on the state of the mechanical system $(\mathbf{q}; \mathbf{p})$. Thus substituting the constraint $\phi(\mathbf{q}, t) = 0$ with the one in eq. (2.38) results the following modified problem,

$$\dot{\mathbf{q}} = \mathbf{p}, \quad (4.6a)$$

$$\mathbf{M}(\mathbf{q})\dot{\mathbf{p}} - \frac{\partial\phi}{\partial\mathbf{q}}^\top \boldsymbol{\lambda} = \mathbf{n}(\mathbf{q}, \mathbf{p}, t), \quad (4.6b)$$

$$-\frac{\partial\phi}{\partial\mathbf{q}} \dot{\mathbf{p}} = \frac{d}{dt} \frac{\partial\phi}{\partial\mathbf{q}} \mathbf{p} + \frac{\partial^2\phi}{\partial t^2} := \mathbf{m}(\mathbf{q}, \mathbf{p}, t). \quad (4.6c)$$

Here to simplify the notation we introduced the symbol $\mathbf{m}(\mathbf{q}, \mathbf{p}, t) \in \mathbb{R}^m$. However even if the idea of differentiating the constraint seem good, the problem in eq. (2.31) and the one in eq. (2.39) are *not equivalent*. This difference produces the so called *drift effect* that will be investigated in next section. At the moment let us ignore this complication and continue on the taken path. System of equations in eq. (2.39) can be also written in compact matrix notation as,

$$\begin{bmatrix} \mathbf{I} & 0 & 0 \\ 0 & \mathbf{M}(\mathbf{q}) & -\frac{\partial\phi}{\partial\mathbf{q}}^\top \\ 0 & -\frac{\partial\phi}{\partial\mathbf{q}} & 0 \end{bmatrix} \begin{bmatrix} \dot{\mathbf{q}} \\ \dot{\mathbf{p}} \\ \boldsymbol{\lambda} \end{bmatrix} = \begin{bmatrix} \mathbf{p} \\ \mathbf{n}(\mathbf{q}, \mathbf{p}, t) \\ \mathbf{m}(\mathbf{q}, \mathbf{p}, t) \end{bmatrix}. \quad (4.7)$$

Now it's easy to see that if the matrix on the left hand side of eq. (2.40) is invertible, then we can express the system in explicit form and use this representation to perform numerical simulations. It easy to see that the matrix is invertible if and only if the sub-block,

$$\begin{bmatrix} \mathbf{M}(\mathbf{q}) & -\frac{\partial\phi}{\partial\mathbf{q}}^\top \\ -\frac{\partial\phi}{\partial\mathbf{q}} & 0 \end{bmatrix} \in \mathbb{R}^{(n+m) \times (n+m)},$$

is non-singular (full rank). In order to prove this we can use the rank additivity formula which states the following,

$$\text{rank} \begin{bmatrix} \mathbf{M}(\mathbf{q}) & -\frac{\partial\phi}{\partial\mathbf{q}}^\top \\ -\frac{\partial\phi}{\partial\mathbf{q}} & 0 \end{bmatrix} = \text{rank} \mathbf{M}(\mathbf{q}) + \text{rank} \left(\frac{\partial\phi}{\partial\mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} \frac{\partial\phi}{\partial\mathbf{q}}^\top \right).$$

First recall that matrix $\mathbf{M}(\mathbf{q})$ is positive definite (so also non-singular), then $\forall \mathbf{q} \in \mathbb{R}^n, \text{rank} \mathbf{M}(\mathbf{q}) = n$. Thus what is missing is the condition

$$\text{rank} \left(\frac{\partial\phi}{\partial\mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} \frac{\partial\phi}{\partial\mathbf{q}}^\top \right) = m, \quad (4.8)$$

which is not generally true. Fortunately assuming the constraints (locally) linearly independent each other, i.e., $\forall \mathbf{q} \in \mathbb{R}^n, \forall t$

$$\text{rank} \left(\frac{\partial\phi}{\partial\mathbf{q}} \right) = m, \quad (4.9)$$

eq. (2.41) holds, matrix in eq. (2.40) is invertible and the overall system can be expressed in *explicit first order form* as follows,

$$\begin{bmatrix} \dot{\mathbf{q}} \\ \dot{\mathbf{p}} \\ \dot{\lambda} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & 0 & 0 \\ 0 & \mathbf{M}(\mathbf{q}) & -\frac{\partial \phi}{\partial \mathbf{q}}^\top \\ 0 & -\frac{\partial \phi}{\partial \mathbf{q}} & 0 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{p} \\ \mathbf{n}(\mathbf{q}, \mathbf{p}, t) \\ \mathbf{m}(\mathbf{q}, \mathbf{p}, t) \end{bmatrix}.$$

An explicit expression for the inverse matrix above can be obtained through tedious calculations and it is shown in appendix A.1. However regardless the specific expression of the solution we want stress a very important point; differentiating the constraint we were able to find an explicit expression for the Lagrange multipliers. For this purpose we differentiated twice and we obtained eq. (2.38). Then isolating the acceleration $\dot{\mathbf{p}}$ from eq. (2.36) results,

$$\dot{\mathbf{p}} = \mathbf{M}(\mathbf{q}) \frac{\partial \phi}{\partial \mathbf{q}}^\top \lambda + \mathbf{M}(\mathbf{q})^{-1} \mathbf{n}(\mathbf{q}, \mathbf{p}, t) = 0. \quad (4.10)$$

Thus substituting eq. (2.43) into eq. (2.38) results,

$$\frac{d}{dt} \frac{\partial \phi}{\partial \mathbf{q}} \mathbf{p} + \frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} \frac{\partial \phi}{\partial \mathbf{q}}^\top \lambda + \frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} \mathbf{n}(\mathbf{q}, \mathbf{p}, t) + \frac{\partial^2 \phi}{\partial t^2} = 0. \quad (4.11)$$

Here we recognize the nonsingular matrix in eq. (2.41), therefore an explicit expression for the Lagrange multipliers is the following,

$$\lambda = - \left(\frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} \frac{\partial \phi}{\partial \mathbf{q}}^\top \right)^{-1} \left(\frac{d}{dt} \frac{\partial \phi}{\partial \mathbf{q}} \mathbf{p} + \frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} \mathbf{n}(\mathbf{q}, \mathbf{p}, t) + \frac{\partial^2 \phi}{\partial t^2} \right).$$

Technically Mechanical systems subject to holonomic constraints are systems of differential-algebraic equations (DAEs) of *differential index* 3. The differential index is one plus the number of differentiations of the constraint that are needed in order to be able to eliminate the Lagrange multipliers. Otherwise, equivalently, the number of times that we need to differentiate the constraint in order to obtain a differential equation also for the Lagrange multiplier.

4.1 The drift effect

In this section we want investigate the effects of the differentiating process done to transform problem in eq. (2.36) into problem eq. (2.39). Unfortunately this method is *not* priceless, indeed formally we substitute an algebraic constraint with its second derivative w.r.t time. To better understand the consequences let us consider a fictitious constraint $\tilde{\phi}(\mathbf{q}, t)$ of the form,

$$\tilde{\phi}(\mathbf{q}, t) := \phi(\mathbf{q}, t) + \mathbf{a}_0 + \mathbf{a}_1 t. \quad (4.12)$$

It is easy to see that consider a mechanical system subject to constraint $\phi(\mathbf{q}, t)$ or $\tilde{\phi}(\mathbf{q}, t)$ differentiating twice is exactly the same, i.e.,

$$\frac{d^2 \phi}{dt^2} = \frac{d^2 \tilde{\phi}}{dt^2}. \quad (4.13)$$

This fact produces the *drift effect*, i.e., the original constraint $\phi(\mathbf{q}, t)$ is violated with an error growing with time. This error grows in the best case linearly, but it can be faster due to *numerical* errors. The reason of this behaviour is the term $\mathbf{a}_1 t$ which does not have any effect in the formulation (2.39). Note also that this is not an undesired effect coming from not enough precision of the numerical scheme, but it is intrinsic effect due to substitution of the constraint ϕ with its second derivative w.r.t time. A simple solution to this undesired effect is the Baumgarte's stabilization that we will discuss in the next section.

Example 12. Continuing example 5 we provide a first order representation for the pendulum; setting $\mathbf{p} = \dot{\mathbf{q}}$ results,

$$\begin{aligned}\mathbf{p} &= \dot{\mathbf{q}}, \\ m\dot{\mathbf{p}} + mg\mathbf{e}_2 &= 2\mathbf{q}\lambda, \\ \mathbf{q}^\top \mathbf{q} - \ell^2 &= 0.\end{aligned}$$

Then in order to eliminate the Lagrange multiplier let us differentiate the constraint,

$$\frac{d\phi}{dt} = 2\mathbf{q}^\top \mathbf{p} = 0.$$

Notice that as expected velocity \mathbf{p} and position \mathbf{q} are orthogonal each other. Differentiating again we obtain

$$\frac{d^2\phi}{dt^2} = 2\dot{\mathbf{q}}^\top \mathbf{p} + 2\mathbf{q}^\top \dot{\mathbf{p}} = 2\mathbf{p}^\top \mathbf{p} + 2\mathbf{q}^\top \dot{\mathbf{p}} = 0.$$

Now we can write the problem in shape of eq. (2.40),

$$\begin{bmatrix} \mathbf{I}_2 & 0 & 0 \\ 0 & m\mathbf{I}_2 & -2\mathbf{q} \\ 0 & -2\mathbf{q}^\top & 0 \end{bmatrix} \begin{bmatrix} \dot{\mathbf{q}} \\ \dot{\mathbf{p}} \\ \lambda \end{bmatrix} = \begin{bmatrix} \mathbf{p} \\ -mg\mathbf{e}_2 \\ 2\mathbf{p}^\top \mathbf{p} \end{bmatrix}, \quad (4.15)$$

◦

where $\mathbf{I}_2 \in \mathbb{R}^{2 \times 2}$ is the identity matrix of dimension 2. The explicit value for the Lagrange multipliers is the following

$$\lambda = -\frac{m}{2\mathbf{q}^\top \mathbf{q}} (\mathbf{p}^\top \mathbf{p} - g\mathbf{q}^\top \mathbf{e}_2) = -\frac{m}{2\ell^2} (\mathbf{p}^\top \mathbf{p} - gq_2). \quad (4.16)$$

To better appreciate the drift effect we can simulate eq. (2.48) using an ODE solver. The results are showed in fig. 2.4. Simulation has been performed for 50 seconds with the following parameters, $g = 9.81$, $m = 2$, $\ell = 2.5$ and initial conditions $\mathbf{q}_0 = (0.86; 2.35) \in \mathbb{R}^2$, $\dot{\mathbf{q}}_0 = \mathbf{p}_0 = (0, 0) \in \mathbb{R}^2$ and $\lambda_0 = 3.69$, which are all consistent with the constraint and with eq. (2.49).

4.2 Baumgarte's stabilization

A very popular way to stabilize the constraint $\phi(\mathbf{q}, t) = 0$ and thus avoid the drift effect is using Baumgarte's stabilization, see [3]. This idea is well presented with other techniques in the survey paper [2]. The key idea is to substitute the constraint with a *differential equation* that *asymptotically* satisfies $\phi(\mathbf{q}, t) = 0$. Obviously this technique is effective only if you are interested in the *long term* behaviour, since during transients the constraint may be violated. Since constrained mechanical systems require to differentiations to eliminate the Lagrange multipliers, we consider a linear second order stable differential equation of the form,

$$\frac{d^2 z(t)}{dt^2} + 2\xi\omega_n \frac{dz(t)}{dt} + \omega_n^2 z(t) = 0, \quad (4.17)$$

where $z(t) \in \mathbb{R}$. In this form we recognize the well known equation for second order linear mechanical or electrical systems, such as the *RLC circuit* or the *mass spring damper*. In these framework $\xi \in \mathbb{R}_{>0}$ is the *damping ratio* and $\omega_n \in \mathbb{R}_{>0}$ is the *natural frequency*. Depending on the value of ξ different cases can arise:

- $0 < \xi < 1$: solution z is *underdamped*,
- $\xi = 1$: solution z is *critically damped*,

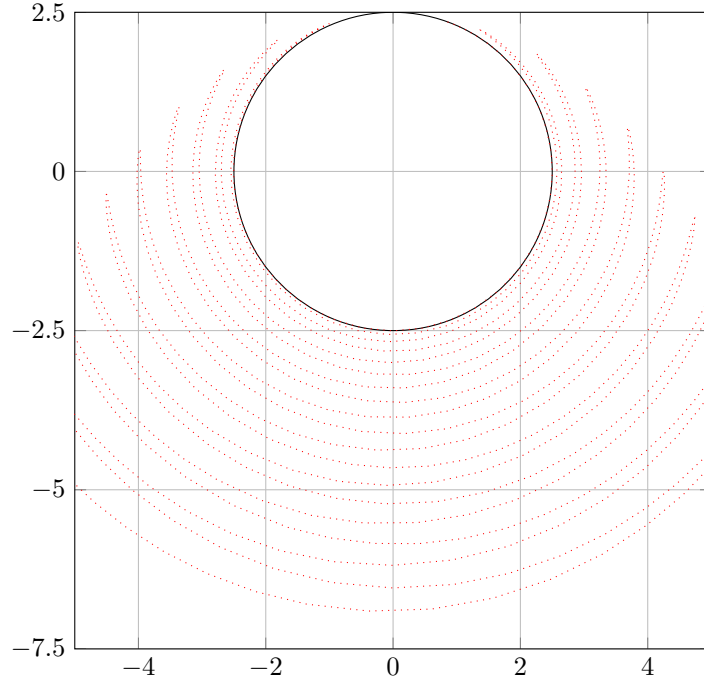


Figure 4.1: Pendulum simulation with drift effect.

- $\xi > 1$: solution z is *overdamped*.

The general form for the solution is the following

$$z(t) = k_1 e^{(-\xi + \sqrt{\xi^2 - 1})\omega_n t} + k_2 e^{(-\xi - \sqrt{\xi^2 - 1})\omega_n t},$$

where $k_1, k_2 \in \mathbb{R}$ are constants determined by the initial conditions $z_0 = z(t=0)$ and $\dot{z}_0 = \dot{z}(t=0)$. Solutions for the different cases can be arranged in different ways, here we show some of them. For a deeper treatment see for example [9]. Let us start from the underdamped case; solution takes the form

$$z(t) = c_1 e^{-\xi\omega_n t} \cos(\sqrt{1 - \xi^2}\omega_n t - c_2), \quad (4.18)$$

where constants $c_1, c_2 \in \mathbb{R}$ are given by,

$$c_1 = \frac{\sqrt{z_0^2 \omega_n^2 + \dot{z}_0^2 + 2z_0 \dot{z}_0 \xi \omega_n}}{\sqrt{1 - \xi^2} \omega_n}$$

$$c_2 = \tan^{-1} \left(\frac{\dot{z}_0 + \xi \omega_n z_0}{x_0 \omega_n \sqrt{1 - \xi^2}} \right).$$

As we can observe in eq. (2.51) solution is a cosine function modulated by a decreasing exponential function. For the critically damped case a linear term appears and the solution takes the form,

$$z(t) = (c_1 + c_2 t) e^{-\omega_n t}, \quad (4.20)$$

where $c_1, c_2 \in \mathbb{R}$ are given by $c_1 = z_0$ and $c_2 = \dot{z}_0 + \omega_n z_0$. Finally for the overdamped case the solution takes the form

$$z(t) = c_1 e^{(-\xi + \sqrt{\xi^2 - 1})\omega_n t} + c_2 e^{(-\xi - \sqrt{\xi^2 - 1})\omega_n t}, \quad (4.21)$$

with constants,

$$c_1 = \frac{x_0 \omega_n (\xi + \sqrt{\xi^2 - 1}) + \dot{x}_0}{2 \omega_n \sqrt{\xi^2 - 1}},$$

$$c_2 = \frac{-x_0 \omega_n (\xi - \sqrt{\xi^2 - 1}) - \dot{x}_0}{2 \omega_n \sqrt{\xi^2 - 1}}.$$

However, regardless of the specific form of the solution, all the three cases satisfy a *global exponential stability* property, i.e., $\forall (z_0, \dot{z}_0) \in \mathbb{R} \times \mathbb{R}$ the solution $z(t, z_0, \dot{z}_0)$ satisfies

$$\lim_{t \rightarrow \infty} z(t, z_0, \dot{z}_0) = 0.$$

The notation $z(t, z_0, \dot{z}_0)$ is just to stress that the solution depends on time but also from initial conditions (z_0, \dot{z}_0) . This consideration suggests to substitute the constraint $\phi(\mathbf{q}, t)$ with an ODE of the form (2.50). In this way, despite the initial conditions, the solution will asymptotically satisfies $\phi(\mathbf{q}, t) = 0$. Thus the explicit *Baumgarte's stabilization* for DAEs of index three takes the form,

$$\frac{d^2 \phi}{dt^2} + 2\xi \omega_n \frac{d\phi}{dt} + \omega_n^2 \phi = 0.$$

Computing explicitly the terms for mechanical systems the stabilization results,

$$\frac{d}{dt} \frac{\partial \phi}{\partial \dot{\mathbf{q}}} \dot{\mathbf{q}} + \frac{\partial \phi}{\partial \mathbf{q}} \ddot{\mathbf{q}} + \frac{\partial^2 \phi}{\partial t^2} + 2\xi \omega_n \frac{\partial \phi}{\partial \dot{\mathbf{q}}} \dot{\mathbf{q}} + 2\xi \omega_n \frac{\partial \phi}{\partial t} + \omega_n^2 \phi = 0. \quad (4.23)$$

which can be re-arranged in the following form

$$-\frac{\partial \phi}{\partial \mathbf{q}} \dot{\mathbf{p}} = \frac{d}{dt} \frac{\partial \phi}{\partial \dot{\mathbf{q}}} \mathbf{p} + \frac{\partial^2 \phi}{\partial t^2} + 2\xi \omega_n \frac{\partial \phi}{\partial \dot{\mathbf{q}}} \mathbf{p} + 2\xi \omega_n \frac{\partial \phi}{\partial t} + \omega_n^2 \phi := \mathbf{m}(\mathbf{q}, \mathbf{p}, t). \quad (4.24)$$

Now eq. (2.57) can be easily implemented without any modification in eq. (2.40); indeed it is sufficient to redefine $\mathbf{m}(\mathbf{q}, \mathbf{p}, t)$ as in eq. (2.57). However a question is still open, how can we tune the parameters ω_n and ξ ?. Unfortunately there is no general rule, but it is useful to notice that the velocity of convergence to zero of the solution depends on both ξ and ω_n . Therefore it is a good idea to select them in order to obtain a high decay rate compatibly with the *numerical stability*. This means that the dynamics associated to the constraint $\phi(\mathbf{q}, t)$ cannot be too fast compared to other ones, otherwise the system could be very hard (stiff) to simulate numerically. Indeed systems where two or more very different time scales are involved are usually called *stiff problems* or *stiff systems*. These may cause some numerical problems, since stiff equations require an integration step extremely small and are numerically unstable. In general it is hard to provide a unique definition of stiffness, but the main idea is that the differential equations includes some terms that can lead to rapid variation in the solution.

Example 13. Continuing example 6 we want use the Baumgarte's technique to stabilize the pendulum. Therefore the idea is to substitute the constraint

$$2\mathbf{p}^\top \mathbf{p} + 2\mathbf{q}\dot{\mathbf{p}} = 0$$

with

$$-2\mathbf{q}^\top \dot{\mathbf{p}} = 2\mathbf{p}^\top \mathbf{p} + 4\omega_n \xi \mathbf{q}^\top \mathbf{p} + \omega_n^2 (\mathbf{q}^\top \mathbf{q} - \ell^2).$$

Thus final equations appears as,

$$\begin{bmatrix} \mathbf{I}_2 & 0 & 0 \\ 0 & m\mathbf{I}_2 & -2\mathbf{q} \\ 0 & -2\mathbf{q}^\top & 0 \end{bmatrix} \begin{bmatrix} \dot{\mathbf{q}} \\ \dot{\mathbf{p}} \\ \lambda \end{bmatrix} = \begin{bmatrix} \mathbf{p} \\ -m\mathbf{g}\mathbf{e}_2 \\ 2\mathbf{p}^\top \mathbf{p} + 4\omega_n \xi \mathbf{q}^\top \mathbf{p} + \omega_n^2 (\mathbf{q}^\top \mathbf{q} - \ell^2) \end{bmatrix}. \quad (4.25)$$

Again to appreciate the effects of Baumgarte's stabilization we can simulate eq. (2.58) with same parameters as example 6 and stabilization parameters $\xi = 0.8$, $\omega_n = 2$. The result is showed in fig. 2.5. \circ

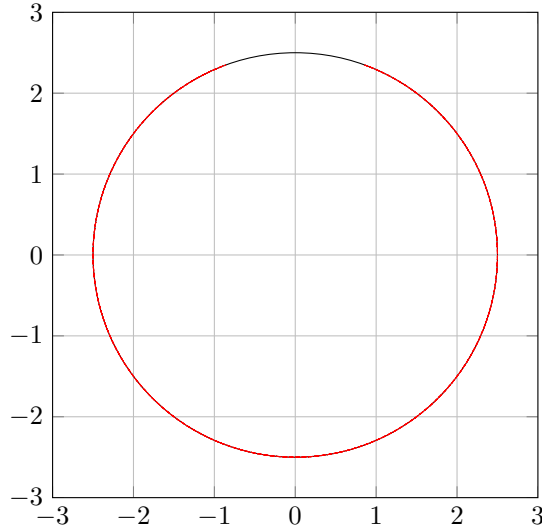


Figure 4.2: Pendulum simulation with Baumgarte's stabilization.

4.3 Projection method

Another way to stabilize the constraint $\phi(\mathbf{q}, t) = 0$ and avoid the drift effect is using the so-called Projection method. The idea is to firstly reduce the index of the DAE, thus obtaining an ODE and the aforementioned *hidden constraints*. Then the resulting ODE can be numerically integrated, leading to a generic intermediate solution $\tilde{\mathbf{q}}(t_n)$. This temporary solution is projected to the hidden constraints in order to minimise the error between the numerical integration solution and the actual system constraints, which is a *nonlinear constrained least squares problem* because of the norm chosen. The projection gives the orthogonal projection to the constraints.

It should be noticed that unlike the previously presented Baumgarte's stabilization method this technique is useful and particularly effective if you are interested in the *short term* behaviour, since even during transients the constraint is always satisfied. Different information from the original and reduced system can be used for various projection methods. For example, we can decide to get an advantage by first projecting the ODE solution $\tilde{\mathbf{q}}(t_n)$ to the position constraints and than to the velocity constraints.

Example 14. Continuing example 6 we want use the Projection technique to stabilize the pendulum. The numerical scheme for integrating the reduced DAE is that in eq. (2.48). We want now to project the results of the generic ODE numerical integrator to the following *hidden constraints*,

$$\frac{d\phi}{dt} = 2\mathbf{q}^T \mathbf{p} = 0.$$

$$\frac{d^2\phi}{dt^2} = 2\dot{\mathbf{q}}^T \mathbf{p} + 2\mathbf{q}^T \dot{\mathbf{p}} = 2\mathbf{p}^T \mathbf{p} + 2\mathbf{q}^T \dot{\mathbf{p}} = 0.$$

At each time step, the resulting orthogonal projection to the constraints is equal to

$$\begin{aligned} \frac{1}{2} \|\tilde{\mathbf{q}} - \mathbf{q}\|_2 &= \min_{\mathbf{q}} \\ \text{subject to: } \phi(\mathbf{q}, t) &= \mathbf{q}^T \mathbf{q} - \ell^2 = 0 \end{aligned} \quad (4.26)$$

for the position coordinates projection, and

$$\begin{aligned} \frac{1}{2} \|\tilde{\mathbf{p}} - \mathbf{p}\|_2 &= \min_{\mathbf{p}} \\ \text{subject to: } \frac{d\phi}{d\mathbf{q}}(\mathbf{q}, t) \mathbf{p} + \frac{d\phi}{dt}(\mathbf{q}, t) &= 0 \end{aligned} \quad (4.27)$$

for the velocity coordinates projection.

◦

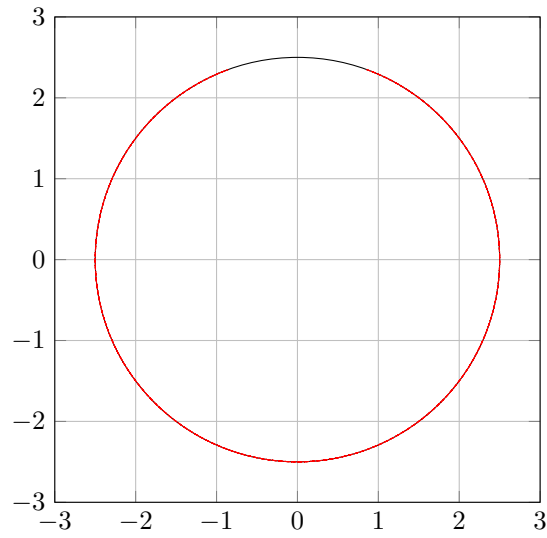


Figure 4.3: Pendulum simulation with Projection method.

Appendix A

A.1 Explicit matrix inverse

Here we provide the explicit expression for the inverse of matrix in eq. (2.40),

$$\begin{bmatrix} \mathbf{I} & 0 & 0 \\ 0 & \mathbf{M}(\mathbf{q}) & -\frac{\partial \phi}{\partial \mathbf{q}}^\top \\ 0 & -\frac{\partial \phi}{\partial \mathbf{q}} & 0 \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{I} & 0 & 0 \\ 0 & \mathbf{X}_{11} & \mathbf{X}_{12} \\ 0 & \mathbf{X}_{21} & \mathbf{X}_{22} \end{bmatrix},$$

where the quantities $\mathbf{X}_{11} \in \mathbb{R}^{n \times n}$, $\mathbf{X}_{12} \in \mathbb{R}^{n \times m}$, $\mathbf{X}_{21} \in \mathbb{R}^{m \times n}$, $\mathbf{X}_{22} \in \mathbb{R}^{m \times m}$ are defined as follows,

$$\begin{aligned} \mathbf{X}_{11} &= \mathbf{M}(\mathbf{q})^{-1} - \mathbf{M}(\mathbf{q})^{-1} \frac{\partial \phi}{\partial \mathbf{q}}^\top \left(\frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} \frac{\partial \phi}{\partial \mathbf{q}}^\top \right)^{-1} \frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1}, \\ \mathbf{X}_{12} &= -\mathbf{M}(\mathbf{q})^{-1} \frac{\partial \phi}{\partial \mathbf{q}}^\top \left(\frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} \frac{\partial \phi}{\partial \mathbf{q}}^\top \right)^{-1}, \\ \mathbf{X}_{21} &= -\left(\frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} \frac{\partial \phi}{\partial \mathbf{q}}^\top \right)^{-1} \frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1}, \\ \mathbf{X}_{22} &= -\left(\frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} \frac{\partial \phi}{\partial \mathbf{q}}^\top \right)^{-1}. \end{aligned}$$

Notice that as expected $\mathbf{X}_{12} = \mathbf{X}_{21}^\top$, moreover in the inversion formula appears the expression,

$$\left(\frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} \frac{\partial \phi}{\partial \mathbf{q}}^\top \right)^{-1}$$

which is exactly the one in eq. (2.41) for which we must guarantee the invertibility. Notice also that in the case of a single constraint, i.e., $\phi(\mathbf{q}) \in \mathbb{R}$ the following expression can be simplified as follows,

$$\mathbf{M}(\mathbf{q})^{-1} \frac{\partial \phi}{\partial \mathbf{q}}^\top \left(\frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} \frac{\partial \phi}{\partial \mathbf{q}}^\top \right)^{-1} \frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} = \frac{\mathbf{M}(\mathbf{q})^{-1} \frac{\partial \phi}{\partial \mathbf{q}}^\top \frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1}}{\frac{\partial \phi}{\partial \mathbf{q}} \mathbf{M}(\mathbf{q})^{-1} \frac{\partial \phi}{\partial \mathbf{q}}^\top}.$$

Bibliography

- [1] Vladimir Igorevich Arnold. *Mathematical methods of classical mechanics*, volume 60. Springer Science & Business Media, 1989.
- [2] Uri M Ascher, Hongsheng Chin, Linda R Petzold, and Sebastian Reich. Stabilization of constrained mechanical systems with daes and invariant manifolds. *Journal of Structural Mechanics*, 23(2):135–157, 1995.
- [3] Joachim Baumgarte. Stabilization of constraints and integrals of motion in dynamical systems. *Computer methods in applied mechanics and engineering*, 1(1):1–16, 1972.
- [4] Earl A Coddington and Norman Levinson. *Theory of ordinary differential equations*. Tata McGraw-Hill Education, 1955.
- [5] Feliks R Gantmacher. *The theory of matrices*. Taylor & Francis, 1964.
- [6] Joao P Hespanha. *Linear systems theory*. Princeton university press, 2009.
- [7] Jerrold E Marsden and Tudor Ratiu. *Introduction to mechanics and symmetry: a basic exposition of classical mechanical systems*, volume 17. Springer Science & Business Media, 2013.
- [8] Linda Petzold. Differential/algebraic equations are not ode’s. *SIAM Journal on Scientific and Statistical Computing*, 3(3):367–384, 1982.
- [9] Singiresu S Rao and Fook Fah Yap. *Mechanical vibrations*, volume 4. Addison-Wesley Reading, 1995.
- [10] Bruno Siciliano, Lorenzo Sciavicco, Luigi Villani, and Giuseppe Oriolo. *Robotics: modelling, planning and control*. Springer Science & Business Media, 2009.