

# Chapter 1. Preliminaries

**Numerical Analysis**, a branch of Applied Mathematics situated at the border between Mathematics and Computer Science, produces methods and procedures for finding numerical solutions of various problems, with a given precision. Standard topics in a Numerical Analysis course are the approximation of problems by simpler problems, the construction of algorithms, iteration methods, error analysis, stability, asymptotic error formulas, the effects of machine arithmetic, etc.

The study will focus on issues such as:

- Problems modeled by functions that do not have an analytical expression, whose values are known only at a discrete set of points. Based on these, we want to approximate values of the function at new points, values of the derivatives or integrals of the function, etc. Such problems lead to finite and divided differences, interpolation formulas, numerical differentiation and integration schemes and many others.
- In many practical situations it is necessary to solve various types of equations or systems of equations, such as: algebraic, transcendental, differential, or integral equations, whose exact solutions cannot be found. They need to be approximated numerically.

An approximating procedure must satisfy several properties:

1. To be *convergent*, meaning the sequence of iterations (successive approximations) must converge to the exact solution, in order to produce “better and better” approximations.
2. To be *stable* (to have stability), meaning that small variations of the input data should lead to small variations in the results (the approximating solutions).
3. From its structure and properties, to be able to estimate the *error* and the *speed of convergence* of the approximating method.

Many times, the conditions that ensure stability of a numerical method coincide with the ones that guarantee its convergence.

# 1 Taylor Polynomials

We start with a very useful tool from Calculus, Taylor's theorem. This will be needed for both the development and understanding of many of the numerical methods discussed later on.

In fact, Taylor polynomials give the very first example of a numerical method, being used as a way to evaluate other functions approximately.

**Theorem 1.1. [Taylor's Theorem]** *Let  $f : [a, b] \rightarrow \mathbb{R}$  be a function with  $n + 1$  continuous derivatives on  $[a, b]$ , for some  $n \geq 0$  and let  $x, x_0 \in [a, b]$ . Then*

$$f(x) = T_n(x) + R_n(x), \quad (1.1)$$

where

$$T_n(x) = f(x_0) + \frac{(x - x_0)}{1!} f'(x_0) + \cdots + \frac{(x - x_0)^n}{n!} f^{(n)}(x_0) \quad (1.2)$$

is **Taylor's polynomial of degree  $n$  of  $f$  at  $x_0$**  and

$$R_n(x) = \frac{1}{n!} \int_{x_0}^x (x - t)^n f^{(n+1)}(t) dt = \frac{(x - x_0)^{n+1}}{(n + 1)!} f^{(n+1)}(\xi), \quad \xi \text{ between } x \text{ and } x_0, \quad (1.3)$$

is the **error** of the Taylor approximation.

**Example 1.2.** Using Taylor's theorem with  $x_0 = 0$ , we obtain the following well-known formulas:

$$\begin{aligned} e^x &= 1 + x + \frac{x^2}{2!} + \cdots + \frac{x^n}{n!} + \frac{x^{n+1}}{(n + 1)!} e^{\xi_x}, \\ \cos x &= 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \cdots + (-1)^n \frac{x^{2n}}{(2n)!} + (-1)^{n+1} \frac{x^{2n+2}}{(2n + 2)!} \cos \xi_x, \\ \sin x &= x - \frac{x^3}{3!} + \frac{x^5}{5!} - \cdots + (-1)^{n-1} \frac{x^{2n-1}}{(2n - 1)!} + (-1)^n \frac{x^{2n+1}}{(2n + 1)!} \sin \xi_x, \\ (1 + x)^a &= 1 + \binom{a}{1} x + \binom{a}{2} x^2 + \cdots + \binom{a}{n} x^n + \binom{a}{n + 1} \frac{x^{n+1}}{(1 + \xi_x)^{n+1-a}}, \end{aligned}$$

with

$$\binom{a}{k} = \frac{a(a - 1) \cdots (a - k + 1)}{k!}, \quad k = 1, 2, 3, \dots, \quad a \in \mathbb{R}.$$

An important special case of the last formula is  $a = -1$ , with  $x$  replaced by  $-x$ :

$$\frac{1}{1-x} = 1 + x + x^2 + \cdots + x^n + \frac{x^{n+1}}{1-x},$$

where the remainder has a simpler form than before. This is easily proved by multiplying both sides by  $1 - x$  and then simplifying. Rearranging the terms, we obtain the familiar formula for a partial sum of the Geometric series:

$$1 + x + x^2 + \cdots + x^n = \frac{1 - x^{n+1}}{1 - x}, \quad x \neq 1.$$

Series representations for the functions in Example 1.2 can be obtained by letting  $n \rightarrow \infty$ . Recall that the series for the first three functions converge for all  $x \in \mathbb{R}$ , while those for the last two converge for  $|x| < 1$ . So, by taking Taylor polynomials of higher and higher degree, we can obtain better and better approximations of the functions above.

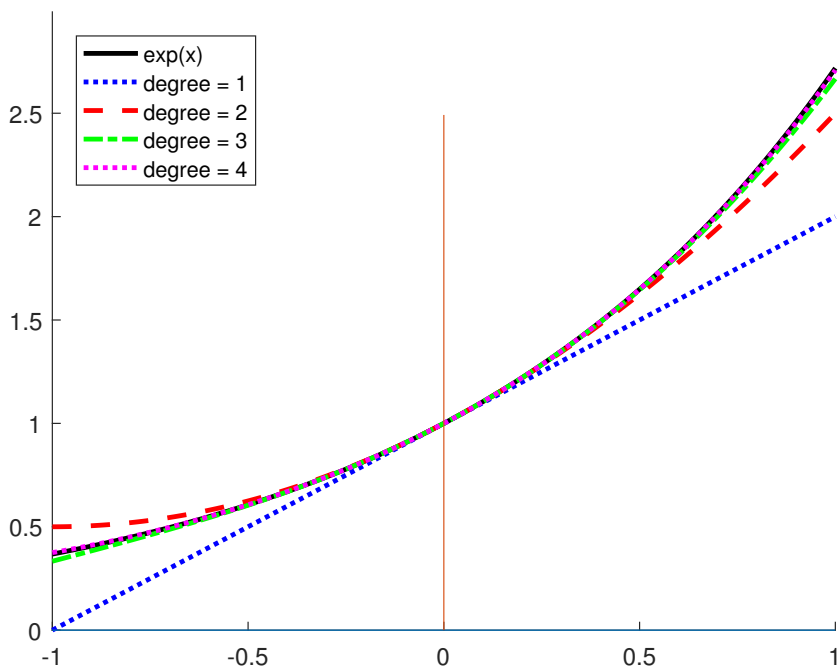


Fig. 1: Taylor approximations of  $e^x$

**Example 1.3.** Consider the function  $f(x) = e^x$ . Figure 1 shows the approximations of  $f$  with

Taylor polynomials of degree 1, 2, 3 and 4, for  $x \in [-1, 1]$ . The errors of these approximations,  $\max_{x \in [-1, 1]} \{e^x - p_n(x)\}$ , are graphed in Figure 2.

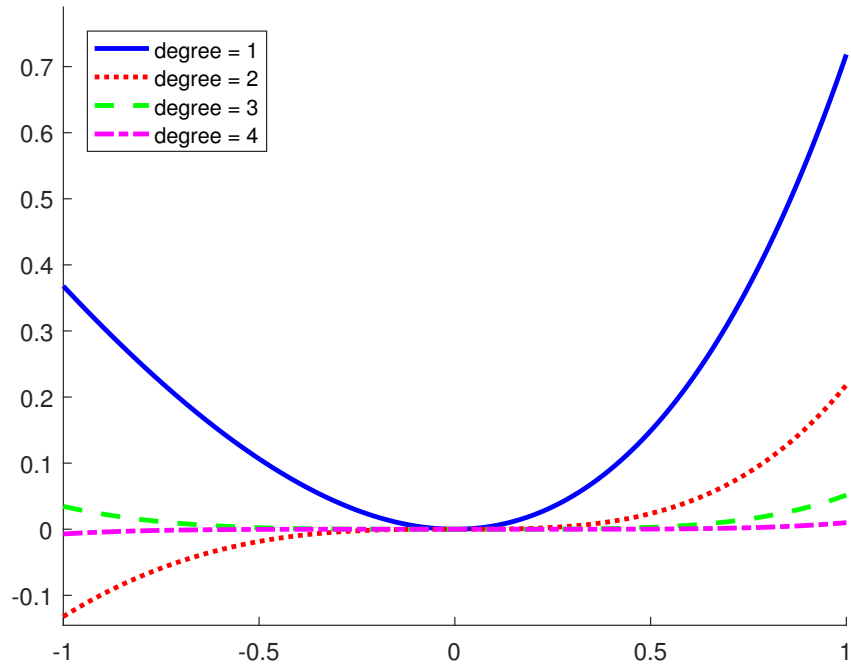


Fig. 2: Errors in Taylor approximations of  $e^x$

### Taylor's formula in two dimensions

**Theorem 1.4.** *Let  $f : D \subset \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be a function which is  $n + 1$  times continuously differentiable on  $D$ , for some  $n \geq 0$  and let  $(x, y), (x_0, y_0) \in D$ . Then*

$$f(x, y) = T_n(x, y) + R_n(x, y), \quad (1.4)$$

where

$$\begin{aligned}
T_n(x, y) &= f(x_0, y_0) + \frac{x - x_0}{1!} f'_x(x_0, y_0) + \frac{y - y_0}{1!} f'_y(x_0, y_0) \\
&+ \frac{1}{2!} \left[ (x - x_0)^2 f''_{x^2}(x_0, y_0) + 2(x - x_0)(y - y_0) f''_{xy}(x_0, y_0) + (y - y_0)^2 f''_{y^2}(x_0, y_0) \right] \\
&+ \sum_{j=1}^n \frac{1}{j!} \left[ (x - x_0) \frac{\partial}{\partial x} + (y - y_0) \frac{\partial}{\partial y} \right]^j f(x, y) \Bigg|_{\substack{x=x_0 \\ y=y_0}}
\end{aligned} \tag{1.5}$$

and

$$R_n(x, y) = \frac{1}{(n+1)!} \left[ (x - x_0) \frac{\partial}{\partial x} + (y - y_0) \frac{\partial}{\partial y} \right]^{n+1} f(x, y) \Bigg|_{\substack{x=\xi \\ y=\eta}}, \tag{1.6}$$

with  $(\xi, \eta)$  a point on the line segment determined by the points  $(x, y)$  and  $(x_0, y_0)$ .

## 2 Errors: Sources, Propagation, Analysis

### 2.1 Types of Numerical Problems

Let us first consider the following simple examples:

**Example 2.1.** Compute

$$\int_1^3 \frac{1}{x} dx.$$

**Solution.** Its exact value is

$$\int_1^3 \frac{1}{x} dx = \ln x \Big|_1^3 = \ln 3 - \ln 1 = \ln 3$$

and its approximation is

$$\int_1^3 \frac{1}{x} dx = 1.0986...$$

■

**Example 2.2.** Solve the equation

$$x^2 = 3, \quad x > 0.$$

**Solution.** The true solution is

$$x = \sqrt{3},$$

with approximating value

$$x = 1.732\dots$$

■

**Example 2.3.** Consider the data

$x_i$	1	2	3	4	5	6	7
$y_i$	5	13	16	23	33	38	40

Discuss the nature of the relationship between  $x$  and  $y$ .

**Solution.** We plot the data (see Figure 3).

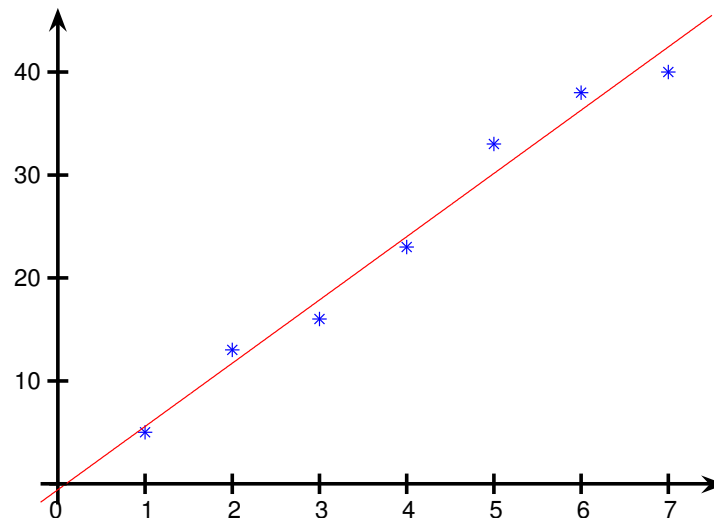


Fig. 3: Scatterplot and linear function

Notice that a straight line “best” approximates the data. So based on these values, we seek a relationship between  $x$  and  $y$  of the form

$$f(x) = ax + b.$$

In the next chapters, we will analyze several rigorous procedures for approximating scattered data. ■

From these examples, we see that a numerical problem is, in general, of the form

$$f(x) = y.$$

- If  $x$  and  $f$  are given, and we seek  $y$ , then this is called a **direct problem** or an **evaluation problem** (Example 2.1).
- If  $y$  and  $f$  are given, and we want to find  $x$ , then it is called an **inverse problem** (Example 2.2).
- If  $x$  and  $y$  are given, and  $f$  must be determined, then we have an **identifying problem** (Example 2.3).

For each type of problem, we obtain an *approximating* value, which is affected by *errors*, perturbations from the true value. The study of errors and their propagation is an important task in Numerical Analysis. In practical problems, it is important that the approximations obtained have *small (negligible)* errors, that do not affect the overall precision of the numerical procedure.

## 2.2 Sources and Propagation of Errors

Let  $\mathcal{A} : \mathbb{R} \rightarrow \mathcal{P}(\mathbb{R})$  be a mapping that assigns a set  $\mathcal{A}(x) \subseteq \mathbb{R}$  to each real number  $x \in \mathbb{R}$ .

**Definition 2.4.** Let  $x \in \mathbb{R}$ . The number  $x^*$  is called an **approximation** of  $x \in \mathbb{R}$ , if  $x^* \in \mathcal{A}(x)$  (notation  $x \approx x^*$ ). The mapping  $\mathcal{A}$  is called an **approximating procedure (method)**.

In practice,  $\mathcal{A}(x)$  consists of numbers that vary “little” from  $x$ .

Let  $x \in \mathbb{R}$  and  $x^* \in \mathcal{A}(x)$  be an approximation of the true (exact) value  $x$ .

**Definition 2.5.** The number

$$\Delta x = x - x^* \tag{2.1}$$

is called the **error** of the approximation  $x^*$ .

If  $\Delta x > 0$ , then  $x^*$  is an **under-approximation**, and if  $\Delta x < 0$ ,  $x^*$  is an **over-approximation**.

For example, for the number  $\pi = 3.141592\dots$ , the number  $x^* = 3.14$  is an under-approximation, while  $x^* = 3.142$  is an over-approximation.

**Definition 2.6.** *The value*

$$|\Delta x| = |x - x^*| \quad (2.2)$$

*is called **absolute error** and the quantity*

$$\delta x = \frac{|\Delta x|}{|x|}, \quad x \neq 0 \quad (2.3)$$

*is called **relative error** .*

Since in practice  $x$  is unknown, we use instead

$$\delta x = \frac{|\Delta x|}{|x^*|}. \quad (2.4)$$

For the relative error in  $\mathbb{R}$ , one can use

$$\delta x = \frac{\Delta x}{x}, \quad (2.5)$$

from which we get

$$\Delta x = x\delta x \implies x^* - x = x\delta x,$$

or

$$x^* = (1 + \delta x)x, \quad (2.6)$$

which is widely used in applications.

For the absolute and relative error, one can also use the notations  $\Delta x^*$  and  $\delta x^*$ .

### Sources of error

- *Mathematical modeling of a physical problem.* A mathematical model for a physical situation is an attempt to give mathematical relationships between certain quantities of physical interest. Because of the complexity of physical reality, a variety of simplifying assumptions are used to construct a more tractable mathematical model. The resulting model has limitations on its accuracy as a consequence of these assumptions, and these limitations may or may not be troublesome, depending on the uses of the model.
- *Blunders (human errors).* In pre-computer times, chance arithmetic errors were always a serious problem. With the introduction of digital computers, the type of blunder has changed. Chance arithmetic errors (e.g., computer malfunctioning) are now relatively rare, and programming errors are currently the main difficulty. Often a program error will be repeated



many times in the course of executing the program, and its existence becomes obvious because of absurd numerical output (although the source of the error may still be difficult to find). This makes good program debugging very important, even though it may not seem very rewarding immediately.

- *Uncertainty in physical data.* Most data from a physical experiment will contain error or uncertainty within it. This must affect the accuracy of any calculations based on the data, limiting the accuracy of the answers.
- *Machine errors.* This means the errors inherent in using the floating-point representation of numbers. Specifically, we mean the rounding/chopping errors and the underflow/overflow errors. The rounding/chopping errors are due to the finite length of the floating-point mantissa; and these errors occur with all of the computer arithmetic operations.
- *Mathematical truncation error.* This name refers to the error of approximation in numerically solving a mathematical problem, and it is the error generally associated with the subject of Numerical Analysis. It involves the approximation of infinite processes by finite ones, replacing noncomputable problems with computable ones, etc.

## 2.3 Propagation of Errors

We distinguish two types of problems:

**Case 1.** Given the errors in the approximations of input data, find the errors in the output.

We start with the case of a function of two variables  $f : D \rightarrow \mathbb{R}$ ,  $D = \{(x, y) | x, y \in \mathbb{R}\}$ . Let  $(x^*, y^*)$  be the approximating values of  $(x, y)$ . Their absolute errors are then

$$\begin{aligned}\Delta x &= x - x^* \Rightarrow x = x^* + \Delta x, \\ \Delta y &= y - y^* \Rightarrow y = y^* + \Delta y.\end{aligned}$$

We want to compute the absolute error

$$\Delta f = f(x, y) - f(x^*, y^*) \tag{2.7}$$

and the relative error  $\delta f$ .

We use Taylor's formula in two variables (1.5), at  $x^*$  and  $y^*$ . We have:

$$\begin{aligned} f(x^* + \Delta x, y^* + \Delta y) &= f(x^*, y^*) + \Delta x f'_x(x^*, y^*) + \Delta y f'_y(x^*, y^*) \\ &+ \frac{(\Delta x)^2}{2!} f''_{x^2}(x^*, y^*) + 2 \frac{\Delta x \Delta y}{2!} f''_{xy}(x^*, y^*) \\ &+ \frac{(\Delta y)^2}{2!} f''_{y^2}(x^*, y^*) + \dots, \end{aligned} \quad (2.8)$$

or

$$\begin{aligned} \Delta f &= \Delta x f'_x(x^*, y^*) + \Delta y f'_y(x^*, y^*) \\ &+ \frac{1}{2!} \left[ (\Delta x)^2 f''_{x^2}(x^*, y^*) + 2 \Delta x \Delta y f''_{xy}(x^*, y^*) + (\Delta y)^2 f''_{y^2}(x^*, y^*) \right] + \dots \end{aligned} \quad (2.9)$$

If  $\Delta x$  and  $\Delta y$  are small, then  $(\Delta x)^2$ ,  $\Delta x \Delta y$  and  $(\Delta y)^2$  can be neglected. Then we find the **absolute error of function  $f$**  or the **maximum absolute error**:

$$\Delta f \simeq \Delta x f'_x(x^*, y^*) + \Delta y f'_y(x^*, y^*). \quad (2.10)$$

In general, for a function of  $n$  variables, we have

$$\Delta f \simeq \Delta x_1 f'_{x_1}() + \Delta x_2 f'_{x_2}() + \dots + \Delta x_n f'_{x_n}() = \sum_{i=1}^n \Delta x_i f'_{x_i}(). \quad (2.11)$$

This is the **propagated error**. Then the relative error  $\delta f$  is

$$\begin{aligned} \delta f &= \frac{\Delta f}{f} \simeq \sum_{i=1}^n \Delta x_i \frac{f'_{x_i}()} {f} = \sum_{i=1}^n \Delta x_i \frac{d}{dx_i} \ln f() \\ &= \sum_{i=1}^n x_i \delta x_i \frac{d}{dx_i} \ln f() = \sum_{i=1}^n x_i \left( \frac{d}{dx_i} \ln f() \right) \delta x_i. \end{aligned} \quad (2.12)$$

**Case 2.** The inverse problem is to determine the precision needed in the input data that will guarantee a (given, preset) accuracy in the output data.

To do this, we use the so-called *principle of equal effects*. This assumes that all terms  $f'_{x_i} \Delta x_i$  in (2.11) have the same effect, i. e.

$$f'_{x_1} \Delta x_1 = f'_{x_2} \Delta x_2 = \dots = f'_{x_n} \Delta x_n.$$

Then (2.11) becomes

$$\Delta f \simeq n \Delta x_i f'_{x_i}(),$$

from which it follows that the absolute error is

$$|\Delta x_i| \simeq \frac{|\Delta f|}{n |f'_{x_i}|} \quad (2.13)$$

and the *relative error* is given by

$$\delta x_i \simeq \frac{\delta f}{n \left| x_i \frac{d}{dx_i} \ln f \right|} \quad (2.14)$$

### 3 Floating-Point Representation of Numbers. Significant Digits

**Definition 3.1.** A number  $r$  written in base  $b$  (an even number) has the *floating-point representation* as

$$r = \pm r_0 r_1 \dots r_p \cdot r_{p+1} \dots r_k \times b^e,$$

where  $r_0, r_1, \dots, r_k$  form the ***mantissa (significand)***,  $e$  is the ***exponent*** and  $b$  is the ***base***.

In order to have uniqueness of the representation, the floating-point numbers are *normalized*, that is, we change the representation, not the value, such that  $r_0 \neq 0$ . The term *floating-point number* will be used to mean a real number that can be exactly represented in this format.

**Definition 3.2.** The ***significant digits*** of a floating-point number written in base  $b$  are any of the digits  $1, 2, \dots, b - 1$  in its representation that are non-zero, or located between non-zero digits, or preceded by at least one non-zero digit.

So 0 can be a significant digit if it does more than just indicate the floating decimal point or fills the places of unknown or omitted digits.

The leftmost significant digit is called the *most significant digit*.

For example,

- the number 7063 has all significant digits;
- the number 0.02340 — the last 4 digits are significant (the zeros in front of 2 merely indicate the decimal point, and are, therefore, *not* significant);
- $1.230 \times 10^3 = 1230$  — zero is a significant digit, as it is preceded by a nonzero digit.

Consider the number  $r > 0$  with the following representation in base 10:

$$r = r_0 10^k + r_1 10^{k-1} + \dots + r_{n-1} 10^{k-n+1} + r_n 10^{k-n} + \dots$$

**Definition 3.3.** *The number*

$$r^* = r_0^* 10^k + r_1^* 10^{k-1} + \dots + r_{n-1}^* 10^{k-n+1}$$

*approximates  $r$  correctly with  $n$  significant digits if*

$$|\Delta r^*| \leq \frac{1}{2} \times 10^{k-n+1} \quad (3.1)$$

**Example 3.4.** Find the number of significant digits with which  $e^* = 2.718282$  approximates correctly the number  $e = 2.71828182\dots$

**Solution.** Let  $e^* = 2.718282$ . This can be written as

$$e^* = 2 \cdot 10^0 + 7 \cdot 10^{-1} + 1 \cdot 10^{-2} + \dots$$

So  $k = 0$ . Then

$$\Delta e^* = e^* - e = 0.00000018 = 0.18 \times 10^{-6} \leq \frac{1}{2} \times 10^{-6}.$$

Comparing it to (3.1), we get

$$10^{k-n+1} = 10^{-6} \Rightarrow 0 - n + 1 = -6 \Rightarrow n = 7.$$

Thus, the number  $e^* = 2.718282$  approximates correctly  $e = 2.71828182\dots$  with 7 significant digits. ■

The following relation exists between the number of significant digits and the relative error:

$$\delta r^* \leq \frac{1}{r_0^* \times 10^{n-1}}, \quad (3.2)$$

where  $r^*$  is the correct approximation with  $n$  significant digits, with the normalized representation

$$r^* = r_0^* 10^k + r_1^* 10^{k-1} + \dots + r_{n-1}^* 10^{k-n+1}. \quad (3.3)$$

This can be considered the *maximum relative error*.

## 4 Stability and Conditioning

A number of mathematical problems have solutions that are quite sensitive to small computational errors, for example rounding errors. To deal with this phenomenon, we use the concepts of *stability* and *condition number*.

Consider a general problem of the type

$$y = f(x), f : \mathbb{R}^m \rightarrow \mathbb{R}^n.$$

We are interested in how “sensitive” is  $f$  to small perturbations of  $x$ , that is, do small perturbations of  $x$  lead to small perturbations in  $y$  – which means the problem is *stable*, or not – *unstable* problem. In the latter case, approximation is not very helpful. We would like to “measure” the degree of sensitivity (how stable or unstable a problem is) by a single number, called the **condition number** of  $f$  at  $x$ .

The function  $f$  is assumed to be evaluated exactly, with infinite precision, as we perturb  $x$ . The condition number of  $f$ , therefore, is an inherent property of the map  $f$  and does not depend on any algorithmic considerations concerning its implementation.

Recall from the propagation of error formula (2.12) that when  $f : \mathbb{R}^m \rightarrow \mathbb{R}$ , we found

$$\delta f \approx \sum_{i=1}^m x_i \frac{f'_{x_i}}{f} \delta x_i.$$

Now, for the general case  $f : \mathbb{R}^m \rightarrow \mathbb{R}^n, x = [x_1, \dots, x_m]^T \in \mathbb{R}^m, f(x) = (f_1(x), \dots, f_n(x))^T \in \mathbb{R}^n$ , let

$$\gamma_{ij} = (\text{cond}_{ij} f)(x) = \frac{x_i \frac{\partial f_j}{\partial x_i}}{f_j(x)}, i = \overline{1, m}, j = \overline{1, n}$$

and

$$\Gamma(x) = [\gamma_{ij}] = \begin{bmatrix} x_1 \frac{\partial f_1}{\partial x_1} & \cdots & x_m \frac{\partial f_1}{\partial x_m} \\ \frac{f_1(x)}{f_1(x)} & \cdots & \frac{f_1(x)}{f_1(x)} \\ \vdots & & \vdots \\ x_1 \frac{\partial f_n}{\partial x_1} & \cdots & x_m \frac{\partial f_n}{\partial x_m} \\ \frac{f_n(x)}{f_n(x)} & \cdots & \frac{f_n(x)}{f_n(x)} \end{bmatrix}, \quad (4.1)$$

called the **conditioning matrix**. Then, the **condition number** of  $f$  at  $x$  is defined by

$$(\text{cond } f)(x) = \|\Gamma(x)\|, \quad (4.2)$$

for a matrix norm  $\|\cdot\|$ . If  $f$  is a linear function, then

$$(\text{cond } f)(x) = \frac{\|x\| \left\| \frac{\partial f}{\partial x} \right\|}{\|f(x)\|}. \quad (4.3)$$

**Particular cases for  $m = n = 1$**

1. If  $x \neq 0, f(x) \neq 0$ , then

$$(\text{cond } f)(x) = \left| \frac{x f'(x)}{f(x)} \right|.$$

2. If  $x = 0, f(x) \neq 0$ , then we take only absolute error for  $x$ ,

$$(\text{cond } f)(x) = \left| \frac{f'(x)}{f(x)} \right|.$$

3. If  $x \neq 0, f(x) = 0$ , then we take only absolute error for  $f(x)$ ,

$$(\text{cond } f)(x) = |x f'(x)|.$$

4. If  $x = f(x) = 0$ , then

$$(\text{cond } f)(x) = |f'(x)|.$$

## Ill-conditioned and ill-posed problems

If the condition number of a problem is large ( $(\text{cond } f)(x) \gg 1$ ), then even for small (relative) errors in the input data, we can expect large errors in the output data. Such problems are called **ill-conditioned** problems. There is no clear line between ill- and well-conditioned problems, it all depends on precision specifications.

If the result of a mathematical problem depends in a discontinuous way on data that varies continuously, then it is impossible to give an accurate numerical solution in a neighborhood of the discontinuity. In such cases, the result can be significantly perturbed, even if the input data is exact and we use high precision procedures. These problems are called **ill-posed** problems. An ill-posed problem can appear if, for example, an integer result is computed from real data (which varies continuously).

### Example 4.1.

#### 1. Number of distinct real roots of a polynomial. Let

$$p(x, c) = x^3 - 2x^2 + x + c, \quad c > -\frac{4}{27}.$$

The polynomial  $p$  can have one real root (if  $c > 0$ ), two distinct real roots (for  $c = 0$ ), or three (if  $c < 0$ ). Therefore, for values of  $c$  close to zero, the number of distinct real zeros of  $p$  is an *ill-posed* problem.

#### 2. Recurrence relations. Let

$$I_n = \int_0^1 \frac{t^n}{t+5} dt, \quad n \in \mathbb{N}.$$

Individual values of this integral can be computed exactly or numerically, e.g.  $I_0 = \ln \frac{6}{5}$ ,  $I_{100} \approx 0.0017$ . What about a recurrence relation?

Writing

$$\frac{t^n}{t+5} = t^{n-1} - 5 \frac{t^{n-1}}{t+5} \quad \text{or} \quad \frac{t^n}{t+5} = \frac{1}{5} \left( t^n - 5 \frac{t^{n+1}}{t+5} \right),$$

we can find *direct recurrence*

$$I_n = \frac{1}{n} - 5I_{n-1} = f_n(I_0),$$

or *inverse recurrence*

$$I_n = \frac{1}{5} \left( \frac{1}{n+1} - 5I_{n+1} \right) = g_n(I_{100}).$$

Then it can be shown that

$$(\text{cond } f_n)(I_0) = \mathcal{O}(5^n) \quad \text{and} \quad (\text{cond } g_n)(I_{100}) = \mathcal{O}(5^{100-n}),$$

which means that direct recurrence is an *ill-conditioned* problem, whereas inverse recurrence is *well-conditioned*.

**Remark 4.2.** In the next chapter, we will consider in greater detail the conditioning of algebraic linear systems and matrices.

## 5 Divided and Finite Differences

These are expressions that are helpful in writing, computing and implementing various iterative numerical procedures.

### 5.1 Divided Differences

**Definition 5.1.** Let  $f : [a, b] \rightarrow \mathbb{R}$  be a differentiable function on  $[a, b]$ , and  $x_i \in [a, b]$ ,  $i = \overline{0, n}$ , be  $n + 1$  distinct nodes. The quantity

$$f[x_0, x_1] = \frac{f(x_1) - f(x_0)}{x_1 - x_0}. \quad (5.1)$$

is called the **first-order divided difference** of  $f$  at the nodes  $x_0$  and  $x_1$ .

**Remark 5.2.**

1. An alternative notation is  $[x_0, x_1; f]$ .
2. The first-order divided difference of a function can be thought of as a *discrete* version of the derivative.
3. If we consider  $f[x_0] = f(x_0)$  the *divided difference of order 0* at  $x_0$ , then (5.1) can be written as

$$f[x_0, x_1] = \frac{f[x_1] - f[x_0]}{x_1 - x_0}. \quad (5.2)$$

We define higher-order divided differences recursively using lower-order ones.



**Definition 5.3.** The *divided difference of order  $n$*  of  $f$  at the distinct nodes  $x_0, x_1, \dots, x_n$  is the quantity

$$f[x_0, x_1, \dots, x_n] = \frac{f[x_1, x_2, \dots, x_n] - f[x_0, x_1, \dots, x_{n-1}]}{x_n - x_0}. \quad (5.3)$$

**Remark 5.4.**

1. The denominator in (5.3) is the difference between the nodes that are *not* common to the differences at the numerator.
2. For easy computation (and implementation) of divided differences, we generate the *table of divided differences*, illustrated below for 4 nodes. The divided differences are obtained on the first row.

$$\begin{array}{ccccccc}
 x_0 & f[x_0] & \longrightarrow & f[x_0, x_1] & \longrightarrow & f[x_0, x_1, x_2] & \longrightarrow & f[x_0, x_1, x_2, x_3] \\
 & & \nearrow & & \nearrow & & \nearrow & \\
 x_1 & f[x_1] & \longrightarrow & f[x_1, x_2] & \longrightarrow & f[x_1, x_2, x_3] & & \\
 & & \nearrow & & \nearrow & & & \\
 x_2 & f[x_2] & \longrightarrow & f[x_2, x_3] & & & & \\
 & & \nearrow & & & & & \\
 x_3 & f[x_3] & & & & & & 
 \end{array}$$

**Example 5.5.** Let  $f(x) = \sin \pi x$ , and the nodes  $x_0 = 0$ ,  $x_1 = \frac{1}{6}$ ,  $x_2 = \frac{1}{2}$ . Let us construct the divided difference table.

**Solution.**

$$\begin{array}{ccccccc}
 x_0 = 0 & f[x_0] = 0 & \longrightarrow & f[x_0, x_1] = \frac{1/2 - 0}{1/6 - 0} = 3 & \longrightarrow & f[x_0, x_1, x_2] = \frac{3/2 - 3}{1/2 - 0} = -3 \\
 & & \nearrow & & \nearrow & & \\
 x_1 = 1/6 & f[x_1] = 1/2 & \longrightarrow & f[x_1, x_2] = \frac{1 - 1/2}{1/2 - 1/6} = 3/2 & & & \\
 & & \nearrow & & & & \\
 x_2 = 1/2 & f[x_2] = 1 & & & & & 
 \end{array}$$

■

## Divided differences with multiple nodes

Divided differences with multiple nodes can be expressed in terms of the derivatives of the function  $f$ , as follows

$$f[x_0, x_0] = \lim_{x_1 \rightarrow x_0} f[x_0, x_1] = \lim_{x_1 \rightarrow x_0} \frac{f(x_1) - f(x_0)}{x_1 - x_0} = f'(x_0),$$

In general, the **divided difference of order  $n$**  at the node  $x_0$ , of multiplicity  $n + 1$ , is defined as

$$f[x_0, x_0, \dots, x_0] = \frac{f^{(n)}(x_0)}{n!}, \quad (5.4)$$

and further, for mixed nodes (some simple, some multiple), we use definition 5.3.

**Example 5.6.** For a function  $f$ , construct the divided differences table for the double node  $x_0$  and the simple node  $x_1$ .

**Solution.** The new nodes are

$$z_0 = x_0, \quad z_1 = x_0, \quad \text{and} \quad z_2 = x_1.$$

The divided differences table:

$$\begin{array}{ccccccc} z_0 = x_0 & f[z_0] = f(x_0) & \longrightarrow & f'(x_0) & \longrightarrow & \frac{\frac{f(x_1) - f(x_0)}{x_1 - x_0} - f'(x_0)}{x_1 - x_0} \\ & & \nearrow & & \nearrow & \\ z_1 = x_0 & f[z_1] = f(x_0) & \longrightarrow & \frac{f(x_1) - f(x_0)}{x_1 - x_0} & & \\ & & \nearrow & & & \\ z_2 = x_1 & f[z_2] = f(x_1) & & & & \end{array}$$

■

**Example 5.7.** Let us see the divided differences table for 3 double nodes.

**Solution.** We have the nodes  $x_0, x_1, x_2$  and the values  $f(x_i), f'(x_i)$ ,  $i = 0, 1, 2$ . We define the sequence of nodes  $z_0, z_1, \dots, z_5$  by

$$z_{2i} = z_{2i+1} = x_i, \quad i = 0, 1, 2.$$

We build the divided difference table relative to the nodes  $z_i, i = \overline{0, 5}$ .

Since  $z_{2i} = z_{2i+1} = x_i$  for every  $i = 0, 1, 2$ ,  $f[z_{2i}, z_{2i+1}] = f[x_i, x_i]$  is a divided difference with a double node and it is equal to  $f'(x_i)$ ; therefore we will use  $f'(x_0), f'(x_1), f'(x_2)$  instead of first order divided differences  $f[z_0, z_1], f[z_2, z_3], f[z_4, z_5]$ .

$$\begin{array}{ccccccc}
 z_0 = x_0 & f[z_0] & \longrightarrow & f[z_0, z_1] = f'(x_0) & \longrightarrow & f[z_0, z_1, z_2] & \dots \\
 & & \nearrow & & \nearrow & & \\
 z_1 = x_0 & f[z_1] & \longrightarrow & f[z_1, z_2] = \frac{f(x_1) - f(x_0)}{x_1 - x_0} & \longrightarrow & f[z_1, z_2, z_3] & \dots \\
 & & \nearrow & & \nearrow & & \\
 z_2 = x_1 & f[z_2] & \longrightarrow & f[z_2, z_3] = f'(x_1) & \longrightarrow & f[z_2, z_3, z_4] & \dots \\
 & & \nearrow & & \nearrow & & \\
 z_3 = x_1 & f[z_3] & \longrightarrow & f[z_3, z_4] = \frac{f(x_2) - f(x_1)}{x_2 - x_1} & \longrightarrow & f[z_3, z_4, z_5] & \dots \\
 & & \nearrow & & \nearrow & & \\
 z_4 = x_2 & f[z_4] & \longrightarrow & f[z_4, z_5] = f'(x_2) & & & \\
 & & \nearrow & & & & \\
 z_5 = x_2 & f[z_5] & & & & & 
 \end{array}$$

■

## 5 Divided and Finite Differences - Continued

### Properties of divided differences

**Theorem 5.1.** *Divided differences have a number of special properties that can simplify work with them:*

a)

$$f[x_0, x_1, \dots, x_n] = \sum_{i=0}^n \frac{f(x_i)}{u'(x_i)} = \sum_{i=0}^n \frac{f(x_i)}{u_i(x_i)}, \quad (5.1)$$

where  $u(x) = (x - x_0)(x - x_1) \dots (x - x_n)$  and  $u_i(x) = \frac{u(x)}{x - x_i}$ .

b) For any permutation  $\{i_0, i_1, \dots, i_n\}$  of the integers  $\{0, 1, \dots, n\}$ ,

$$f[x_{i_0}, x_{i_1}, \dots, x_{i_n}] = f[x_0, x_1, \dots, x_n]. \quad (5.2)$$

c)

$$f[x_0, x_1, \dots, x_n] = \frac{(Wf)(x_0, x_1, \dots, x_n)}{V(x_0, x_1, \dots, x_n)}, \quad (5.3)$$

where

$$Wf(x_0, x_1, \dots, x_n) = \begin{vmatrix} 1 & x_0 & x_0^2 & \dots & x_0^{n-1} & f(x_0) \\ 1 & x_1 & x_1^2 & \dots & x_1^{n-1} & f(x_1) \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^{n-1} & f(x_n) \end{vmatrix} \quad \text{and}$$

$$V(x_0, x_1, \dots, x_n) = \begin{vmatrix} 1 & x_0 & x_0^2 & \dots & x_0^n \\ 1 & x_1 & x_1^2 & \dots & x_1^n \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^n \end{vmatrix} = \prod_{0 \leq j < i \leq n} (x_i - x_j) \text{ is the Vandermonde determinant.}$$

d) Let  $e_k(x) = x^k, k \geq 0$ . Then

$$e_k[x_0, x_1, \dots, x_n] = \begin{cases} 0, & k < n \\ 1, & k = n \end{cases}.$$

For a polynomial of degree  $k$ ,  $P_k = a_0 + a_1x + \cdots + a_kx^k$ ,

$$P_k[x_0, x_1, \dots, x_n] = \begin{cases} 0, & k < n \\ a_n, & k = n \end{cases}. \quad (5.4)$$

e) If  $f \in C^n[a, b]$ , where  $[a, b]$  is the smallest interval containing the distinct nodes  $\{x_0, \dots, x_n\}$ , then there exists  $\xi_n \in (a, b)$  such that

$$f[x_0, x_1, \dots, x_n] = \frac{1}{n!} f^{(n)}(\xi_n), \quad (5.5)$$

(the mean-value formula for divided differences).

**Remark 5.2.** As a consequence of part e), if  $f \in C^n[a, b]$  and  $\alpha \in [a, b]$ , then

$$\lim_{x_0, \dots, x_n \rightarrow \alpha} f[x_0, x_1, \dots, x_n] = \lim_{\xi_n \rightarrow \alpha} \frac{f^{(n)}(\xi_n)}{n!} = \frac{1}{n!} f^{(n)}(\alpha),$$

the computational formula for divided differences with multiple nodes.

## 5.2 Finite Differences

**Definition 5.3.** Consider the equidistant nodes  $x_i = x_0 + ih, i = 0, 1, \dots, n, h > 0$ . The quantity

$$\Delta^1 f(x_i) = f(x_{i+1}) - f(x_i) = f_{i+1} - f_i \quad (5.6)$$

is called the **first-order forward difference** of  $f$  with step  $h$  at  $x_i$ , and

$$\Delta^k f(x_i) = \Delta^{k-1} f(x_{i+1}) - \Delta^{k-1} f(x_i) = \Delta^{k-1} f_{i+1} - \Delta^{k-1} f_i \quad (5.7)$$

is the  **$k$ th-order forward difference** of  $f$  with step  $h$ , at  $x_i$ .

**Remark 5.4.**

1. As a convention, we use  $\Delta^0 f(x_i) = f(x_i) = f_i$ .
2. For easy computation (and implementation) of forward differences, we construct a *table of forward differences*, similar to the one used for divided differences, illustrated below for 4 nodes.

$$\begin{array}{c|ccccccc}
x_0 & f_0 & \longrightarrow & \Delta f_0 & \longrightarrow & \Delta^2 f_0 & \longrightarrow & \Delta^3 f_0 \\
& & \nearrow & & \nearrow & & \nearrow & \\
x_1 & f_1 & \longrightarrow & \Delta f_1 & \longrightarrow & \Delta^2 f_1 & & \\
& & \nearrow & & \nearrow & & & \\
x_2 & f_2 & \longrightarrow & \Delta f_2 & & & & \\
& & \nearrow & & & & & \\
x_3 & f_3 & & & & & & 
\end{array}$$

In a similar way, we define the **backward difference**  $\nabla$  by

$$\begin{aligned}
\nabla^0 f_i &= f_i, \\
\nabla^1 f_i &= f_i - f_{i-1}, \\
\nabla^k f_i &= \nabla^{k-1} f_i - \nabla^{k-1} f_{i-1},
\end{aligned} \tag{5.8}$$

and they can also be easily computed in a table.

$$\begin{array}{c|ccccccc}
x_0 & f_0 & & & & & & \\
& & \searrow & & & & & \\
x_1 & f_1 & \longrightarrow & \nabla f_1 & & & & \\
& & \searrow & & \searrow & & & \\
x_2 & f_2 & \longrightarrow & \nabla f_2 & \longrightarrow & \nabla^2 f_2 & & \\
& & \searrow & & \searrow & & \searrow & \\
x_3 & f_3 & \longrightarrow & \nabla f_3 & \longrightarrow & \nabla^2 f_3 & \longrightarrow & \nabla^3 f_3
\end{array}$$

**Remark 5.5.** These differences are referred to collectively as *finite differences*. Usually, if nothing is specified, by “finite” differences we mean “forward” differences.

Denote by

$$X = \{x_i \mid x_i = x_0 + ih, \ i = \overline{0, n}, x_0, h \in \mathbb{R}\}$$

and for  $f : X \rightarrow \mathbb{R}$ , by  $f_i = f(x_i)$ .

Let us write a few finite differences and notice a pattern.

$$\begin{aligned}
\Delta^1 f(x_i) &= f_{i+1} - f_i, \\
\Delta^2 f(x_i) &= \Delta^1 f_{i+1} - \Delta^1 f_i = f_{i+2} - f_{i+1} - (f_{i+1} - f_i) = f_{i+2} - 2f_{i+1} + f_i, \\
\Delta^3 f(x_i) &= \Delta^2 f_{i+1} - \Delta^2 f_i = f_{i+3} - 2f_{i+2} + f_{i+1} - (f_{i+2} - 2f_{i+1} + f_i) \\
&= f_{i+3} - 3f_{i+2} + 3f_{i+1} - f_i.
\end{aligned}$$

It can easily be proved (by induction) that

$$\Delta^n f(x_i) = \sum_{k=0}^n (-1)^k \binom{n}{k} f_{n-k+i},$$

or, equivalently, by the symmetry of combinations,

$$\Delta^n f(x_i) = \sum_{k=0}^n (-1)^{n-k} \binom{n}{k} f_{k+i}.$$

In particular, we have

$$\Delta^n f(x_0) = \sum_{k=0}^n (-1)^{n-k} \binom{n}{k} f_k. \quad (5.9)$$

Finite and divided differences for equally spaced nodes are closely related.

**Proposition 5.6.** *Let  $f : X \rightarrow \mathbb{R}$ . Then*

$$f[a, a+h, \dots, a+nh] = \frac{1}{n!h^n} \Delta^n f(a). \quad (5.10)$$

# Chapter 2. Numerical Solution of Systems of Linear Algebraic Equations

Systems of simultaneous linear equations occur in solving problems in a wide variety of disciplines, including Mathematics, Statistics, physical, biological and social sciences, engineering, business and many more. They arise directly in solving real-world problems, and they also occur as part of the solution process for other problems. Numerical solutions of boundary value problems and initial boundary value problems for differential equations are a rich source of linear systems, especially large-size ones.

In this chapter, we will examine the following problem: given a matrix  $A \in \mathbb{R}^{n \times n}$  and a vector  $b \in \mathbb{R}^n$ , find  $x \in \mathbb{R}^n$  such that

$$Ax = b.$$

There are two types of methods for the solution of algebraic linear systems:

- *direct (exact)* methods, that provide a solution in a finite number of steps (e.g., Cramer, Gaussian elimination, factorizations);
- *iterative* methods, which approximate the solution by a sequence converging to it (e.g., Jacobi, Gauss-Seidel, SOR).

## 1 Direct Methods

### 1.1 Gaussian Elimination

A linear system is easy to solve when the matrix of the system is *triangular*:

**Definition 1.1.** A matrix  $A = [a_{ij}]_{i,j=\overline{1,n}}$  is called

- *upper triangular*, if  $a_{ij} = 0, \forall i > j$ ,

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ & a_{22} & \dots & a_{2n} \\ & & \ddots & \vdots \\ 0 & & & a_{nn} \end{bmatrix}, \quad (1.1)$$



- **lower triangular**, if  $a_{ij} = 0, \forall i < j$ ,

$$A = \begin{bmatrix} a_{11} & & & 0 \\ a_{21} & a_{22} & & \\ \vdots & & \ddots & \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}, \quad (1.2)$$

- **diagonal**, if it is both upper and lower triangular,  $a_{ij} = 0, \forall i \neq j$ ,

$$A = \begin{bmatrix} a_{11} & & & 0 \\ & a_{22} & & \\ & & \ddots & \\ 0 & & & a_{nn} \end{bmatrix}. \quad (1.3)$$

**Remark 1.2.** The determinant of an upper or lower triangular matrix is equal to the product of its diagonal elements

$$\det(A) = a_{11}a_{22} \dots a_{nn}.$$

So, an upper or lower triangular matrix is nonsingular if and only if all of its diagonal entries are nonzero.

**Example 1.3.** Solve the triangular systems

**a)**

$$\begin{bmatrix} 2 & 4 & 2 \\ 0 & -1 & 1 \\ 0 & 0 & -1 \end{bmatrix} x = \begin{bmatrix} 8 \\ 0 \\ -1 \end{bmatrix}, \quad (1.4)$$

**b)**

$$\begin{bmatrix} 1 & 0 & 0 \\ 1/2 & 1 & 0 \\ 1/2 & 1 & 1 \end{bmatrix} x = \begin{bmatrix} 8 \\ 4 \\ 3 \end{bmatrix}. \quad (1.5)$$

**Solution.**

**a)** The upper triangular system is

$$\begin{cases} 2x_1 + 4x_2 + 2x_3 = 8 \\ -x_2 + x_3 = 0 \\ -x_3 = -1 \end{cases}$$

We start from the bottom (the last equation) and solve recursively for each unknown:

$$\begin{aligned} x_3 &= \frac{-1}{-1} = 1, \\ x_2 &= \frac{1}{-1}(0 - x_3) = 1, \\ x_1 &= \frac{1}{2}(8 - 4x_2 - 2x_3) = 1. \end{aligned}$$

We found the solution

$$x = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = [1 \ 1 \ 1]^T.$$

**b)** For the lower triangular system:

$$\begin{cases} x_1 = 8 \\ 1/2x_1 + x_2 = 4 \\ 1/2x_1 + x_2 + x_3 = 3 \end{cases}$$

we start from the top and solve each equation going down:

$$\begin{aligned} x_1 &= \frac{8}{1} = 8, \\ x_2 &= \frac{1}{1}\left(4 - \frac{1}{2}x_1\right) = 0, \\ x_3 &= \frac{1}{1}\left(3 - \frac{1}{2}x_1 - x_2\right) = -1. \end{aligned}$$

So the solution is

$$x = \begin{bmatrix} 8 \\ 0 \\ -1 \end{bmatrix} = [8 \ 0 \ -1]^T.$$

■

So, in general, for a nonsingular upper triangular matrix  $U$ , the system  $Ux = b$  is easily solved by **backward substitution**:

$$\begin{aligned} x_n &= \frac{b_n}{u_{nn}}, \\ x_i &= \frac{1}{u_{ii}} \left( b_i - \sum_{j=i+1}^n u_{ij} x_j \right), \quad i = \overline{n-1, 1} \end{aligned} \quad (1.6)$$

and if the nonsingular matrix  $L$  is lower triangular, then the system  $Lx = b$  is solved by **forward substitution**:

$$\begin{aligned} x_1 &= \frac{b_1}{l_{11}}, \\ x_i &= \frac{1}{l_{ii}} \left( b_i - \sum_{j=1}^{i-1} l_{ij} x_j \right), \quad i = \overline{2, n}. \end{aligned} \quad (1.7)$$

**Gaussian elimination** is a procedure for transforming a system into an equivalent (upper) triangular one, by doing the following elementary row operations:

- multiplying a row (equation) by a constant  $\lambda \neq 0$ ,

$$(\lambda R_i) \rightarrow (R_i),$$

- multiplying a row by a constant  $\lambda \neq 0$  and adding it to another row,

$$(R_i + \lambda R_j) \rightarrow (R_i),$$

- interchanging (permuting) two rows,

$$(R_i) \longleftrightarrow (R_j).$$

All these elementary operations are performed on the *augmented (extended)* matrix of the system

$$\tilde{A} = [A \mid b] = \left[ \begin{array}{cccc|c} a_{11} & a_{12} & \dots & a_{1n} & a_{1,n+1} \\ a_{21} & a_{22} & \dots & a_{2n} & a_{2,n+1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} & a_{n,n+1} \end{array} \right], \quad (1.8)$$

where  $a_{i,n+1} = b_i$ ,  $i = \overline{1, n}$ .

Gaussian elimination goes as follows:

Assuming  $a_{11} \neq 0$ , at the first step, we eliminate (make 0) the coefficients of  $x_1$  from every row below, i.e., every  $R_j, j = \overline{2, n}$ , using  $a_{11}$ , i.e. by

$$(R_j - \frac{a_{j1}}{a_{11}} R_1) \rightarrow (R_j).$$

Then we proceed the same for the coefficients of each  $x_i, i = \overline{2, n-1}, j = \overline{i+1, n}$ . This way we obtain a finite sequence of augmented matrices

$$\tilde{A}^{(1)}, \tilde{A}^{(2)}, \dots, \tilde{A}^{(n)},$$

where  $\tilde{A}^{(1)} = \tilde{A}$  and (at step  $k$ )  $\tilde{A}^{(k)} = [a_{ij}^{(k)}]$  obtained by

$$\left( R_i - \frac{a_{i,k-1}^{(k-1)}}{a_{k-1,k-1}^{(k-1)}} R_{k-1} \right) \rightarrow (R_i).$$

Here, we denoted by  $a_{ij}^{(l)}$  the  $(i, j)$  entry at step  $l$ . . The quantities

$$m_{i,k-1} = \frac{a_{i,k-1}^{(k-1)}}{a_{k-1,k-1}^{(k-1)}}, \quad i = k, \dots, n \quad (1.9)$$

are called *multipliers*. For equations  $i = k, \dots, n$ , we subtract  $m_{i,k-1}$  times  $E_{k-1}$  from  $E_i$ , eliminating  $x_{k-1}$  from  $E_i$ . The new coefficients and the right-hand sides in equations  $E_k$  through  $E_n$  are defined by

$$\begin{aligned} a_{ij}^{(k)} &= a_{ij}^{(k-1)} - m_{i,k-1} a_{kj}^{(k-1)}, \quad i, j = k, \dots, n, \\ b_i^{(k)} &= b_i^{(k-1)} - m_{i,k-1} b_k^{(k-1)}, \quad i = k, \dots, n. \end{aligned}$$

At every step  $k$ , the system corresponding to the augmented matrix  $\tilde{A}^{(k)}$  is equivalent to the original linear system (meaning, it has the same solution) and in it, the variable  $x_{k-1}$  was eliminated from the equations  $E_k, E_{k+1}, \dots, E_n$ . Then the system corresponding to  $\tilde{A}^{(n)}$  is an equivalent upper triangular one:

$$\begin{cases} a_{11}^{(1)}x_1 + a_{12}^{(1)}x_2 + \dots + a_{1n}^{(1)}x_n = a_{1,n+1}^{(1)} \\ \quad + a_{22}^{(2)}x_2 + \dots + a_{2n}^{(2)}x_n = a_{2,n+1}^{(2)} \\ \quad \quad \quad \ddots \quad \quad \quad \vdots \\ \quad \quad \quad \quad \quad \quad a_{nn}^{(n)}x_n = a_{n,n+1}^{(n)} \end{cases}, \quad (1.10)$$

which is solved by backward substitution (1.6).

Of course, in all this we need  $a_{ii}^{(i)} \neq 0$ . The element  $a_{ii}^{(i)}$  is called **pivot**. If at any time during the elimination process, we find  $a_{kk}^{(k)} = 0$ , then we look further down in that column for a pivot, i.e., we interchange rows

$$(R_k) \longleftrightarrow (R_p),$$

where  $p$  is the smallest integer  $k+1 \leq p \leq n$  with  $a_{pk}^{(k)} \neq 0$ . If the original system is nonsingular, it can be shown that one of the equations following  $E_k$  *must* contain a term involving  $x_k$  with a nonzero coefficient, so it is always possible to find a pivot.

In fact, in practice, when implementing Gaussian elimination, pivoting is necessary even if the pivot is *not* zero, but small, compared to the rest of the elements in that column. We should avoid using coefficients that are nearly zero as pivot elements, because such a pivot can produce substantial rounding errors and even cancellations. Instead, we can use several types of *pivoting*.

- We can choose the pivot to be the largest element (in absolute value) in that column, below the main diagonal, i.e.

$$|a_{pk}^{(k)}| = \max_{k \leq l \leq n} |a_{lk}^{(k)}|. \quad (1.11)$$

In most instances, it decreases the propagated effects of rounding errors. With partial pivoting, the multipliers  $m_{i,k}$  in (1.9) will satisfy

$$|m_{ik}| \leq 1, \quad 1 \leq k < i \leq n.$$

This will help reduce loss-of-significance errors, because multiplications by  $m_{ik}$  will not lead to much larger numbers.

This is called **partial pivoting (maximal pivoting on columns)** and it is the most popular one in practice.

- We can do **scaled pivoting on columns**: First, we define a scaling factor for each row

$$s_i = \max_{j=\overline{1,n}} |a_{ij}| \text{ or } s_i = \sum_{j=1}^n |a_{ij}|, \quad i = \overline{1,n}.$$

If there exists an  $i$  such that  $s_i = 0$ , then the matrix is singular. For a nonsingular matrix, we use the scaling factor to choose the pivot. At each step  $i$ , we find the smallest  $p$ ,  $i \leq p \leq n$  such that

$$\frac{|a_{pi}|}{s_i} = \max_{1 \leq j \leq n} \frac{|a_{ji}|}{s_j} \quad (1.12)$$

and then interchange rows  $(R_i) \longleftrightarrow (R_p)$  so the pivot is  $a_{pi}$ . This ensures the fact that the maximal element in each column has the relative size 1, before we compare and interchange rows. Also, dividing by the scaling factor does not produce any extra rounding errors.

- The third method is **total (maximal) pivoting**. At each step  $k$ , we find

$$|a_{pq}| = \max\{|a_{ij}|, i, j = \overline{k, n}\} \quad (1.13)$$

and interchange both the rows and the columns,

$$(R_k) \longleftrightarrow (R_p), \quad (C_k) \longleftrightarrow (C_q).$$

But then we have to keep track of the columns (unknowns) interchanges.

**Remark 1.4.** If  $A$  is singular of rank  $p - 1$ , then at step  $p$  we get

$$\tilde{A}^{(p)} = \left[ \begin{array}{cccccc|c} a_{11}^{(1)} & a_{12}^{(1)} & \cdots & a_{1,p-1}^{(1)} & a_{1p}^{(1)} & \cdots & a_{1n}^{(1)} & a_{1,n+1}^{(1)} \\ 0 & a_{22}^{(2)} & \cdots & a_{2,p-1}^{(2)} & a_{2p}^{(2)} & \cdots & a_{2n}^{(2)} & a_{2,n+1}^{(2)} \\ \vdots & & \ddots & \vdots & \vdots & & \vdots & \vdots \\ \vdots & & & a_{p-1,p-1}^{(p-1)} & a_{p-1,p}^{(p-1)} & & a_{p-1,n}^{(p-1)} & a_{p-1,n+1}^{(p-1)} \\ \vdots & & & & 0 & \cdots & 0 & a_{p,n+1}^{(p)} \\ \vdots & & & & & \ddots & \vdots & \vdots \\ 0 & \cdots & \cdots & \cdots & \cdots & \cdots & 0 & a_{n,n+1}^{(n)} \end{array} \right].$$

So, if  $a_{i,n+1}^{(i)} = b_i^{(i)} = 0$ , for *all*  $i = p, p+1, \dots, n$ , then the system is compatible, but undetermined (i.e., it has an infinite number of solutions), otherwise, the system is incompatible (no solution).

Thus, Gaussian elimination can also be used to discuss the solvability of the linear system.

**Example 1.5.** Solve the system

$$\begin{cases} x_1 - x_2 + x_3 = -1 \\ -2x_1 + 2x_2 + x_3 = 2 \\ -3x_1 - x_2 + 5x_3 = -5 \end{cases},$$

by Gaussian elimination with different types of pivoting.

**Solution.** The augmented matrix of the system is

$$\tilde{A} = \left[ \begin{array}{ccc|c} 1 & -1 & 1 & -1 \\ -2 & 2 & 1 & 2 \\ -3 & -1 & 5 & -5 \end{array} \right],$$

### Partial pivoting

On the first column, the largest element in absolute value is  $-3$ , so that will be the pivot. Thus, first we interchange  $(R_1) \longleftrightarrow (R_3)$ . We get

$$\tilde{A} \sim \left[ \begin{array}{ccc|c} \boxed{-3} & -1 & 5 & -5 \\ -2 & 2 & 1 & 2 \\ 1 & -1 & 1 & -1 \end{array} \right] \sim \left[ \begin{array}{ccc|c} -3 & -1 & 5 & -5 \\ 0 & 8/3 & -7/3 & 16/3 \\ 0 & -4/3 & 8/3 & -8/3 \end{array} \right] \begin{array}{l} (-2/3 R_1 + R_2) \rightarrow (R_2) \\ (1/3 R_1 + R_3) \rightarrow (R_3) \end{array}$$

Further, we have

$$\tilde{A} \sim \left[ \begin{array}{ccc|c} -3 & -1 & 5 & -5 \\ 0 & \boxed{8/3} & -7/3 & 16/3 \\ 0 & -4/3 & 8/3 & -8/3 \end{array} \right] \sim \left[ \begin{array}{ccc|c} -3 & -1 & 5 & -5 \\ 0 & 8/3 & -7/3 & 16/3 \\ 0 & 0 & 3/2 & 0 \end{array} \right] \begin{array}{l} (\frac{1}{2} R_2 + R_3) \rightarrow (R_3) \end{array}$$

Now we solve by back substitution (1.6), to get

$$x = [1 \ 2 \ 0]^T.$$

### Scaled partial pivoting

We compute the scaling factors using sums on each row. At the first step  $k = 1$ , we get

$$s = [3, 5, 9]$$
$$\left[ \frac{|a_{j,1}|}{s_j} \right] = [1/3, 2/5, 3/9] = [5/15, 6/15, 5/15], j = 1, 2, 3.$$

The maximum of the three fractions is the second, so  $p = 2$ . We interchange  $(R_1) \longleftrightarrow (R_2)$  and make zeros below it.

$$\tilde{A} \sim \left[ \begin{array}{ccc|c} \boxed{-2} & 2 & 1 & 2 \\ 1 & -1 & 1 & -1 \\ -3 & -1 & 5 & -5 \end{array} \right] \sim \left[ \begin{array}{ccc|c} -2 & 2 & 1 & 2 \\ 0 & 0 & 3/2 & 0 \\ 0 & -4 & 7/2 & -8 \end{array} \right] \begin{array}{l} (1/2 R_1 + R_2) \rightarrow (R_2) \\ (-3/2 R_1 + R_3) \rightarrow (R_3) \end{array}$$

At step  $k = 2$ , obviously,  $p = 3$ , so we interchange  $(R_2) \longleftrightarrow (R_3)$ ,

$$\tilde{A} \sim \left[ \begin{array}{ccc|c} -2 & 2 & 1 & 2 \\ 0 & -4 & 7/2 & -8 \\ 0 & 0 & 3/2 & 0 \end{array} \right]$$

and we are done. By back substitution we get the (obviously, same) solution

$$x = [1 \ 2 \ 0]^T.$$

### Total pivoting

At step  $k = 1$ , since

$$\max_{i,j=1,3} |a_{ij}| = 5 = |a_{33}|,$$

we interchange both rows and columns,  $(R_1) \longleftrightarrow (R_3)$ ,  $(C_1) \longleftrightarrow (C_3)$ , to get

$$\tilde{A} \sim \left[ \begin{array}{ccc|c} -3 & -1 & 5 & -5 \\ -2 & 2 & 1 & 2 \\ 1 & -1 & 1 & -1 \end{array} \right] \sim \left[ \begin{array}{ccc|c} 5 & -1 & -3 & -5 \\ 1 & 2 & -2 & 2 \\ 1 & -1 & 1 & -1 \end{array} \right],$$

which is now a system for the *new unknown*  $x' = [x_3 \ x_2 \ x_1]^T$ . We proceed to make zeros on the



first column below the diagonal.

$$\tilde{A} \sim \left[ \begin{array}{ccc|c} \boxed{5} & -1 & -3 & -5 \\ 1 & 2 & -2 & 2 \\ 1 & -1 & 1 & -1 \end{array} \right] \sim \left[ \begin{array}{ccc|c} 5 & -1 & -3 & -5 \\ 0 & 11/5 & -7/5 & 3 \\ 0 & -4/5 & 8/5 & 0 \end{array} \right] \begin{array}{l} (-1/5 R_1 + R_2) \rightarrow (R_2) \\ (-1/5 R_1 + R_3) \rightarrow (R_3) \end{array}$$

At step  $k = 2$ ,

$$\max_{i,j=2,3} |a_{ij}| = \frac{11}{5} = |a_{22}|,$$

so no (row or column) interchanges are necessary. We have

$$\tilde{A} \sim \left[ \begin{array}{ccc|c} 5 & -1 & -3 & -5 \\ 0 & \boxed{11/5} & -7/5 & 3 \\ 0 & -4/5 & 8/5 & 0 \end{array} \right] \sim \left[ \begin{array}{ccc|c} 5 & -1 & -3 & -5 \\ 0 & 11/5 & -7/5 & 3 \\ 0 & 0 & 12/11 & 12/11 \end{array} \right] \left( \frac{4}{11} R_2 + R_3 \right) \rightarrow (R_3)$$

By back substitution, we get

$$x' = [0 \ 2 \ 1]^T \text{ and } x = [1 \ 2 \ 0]^T.$$

■

### Remark 1.6.

1. The elements under the main diagonal (which become 0) need not be computed.
2. When pivoting, we do not need to *physically* interchange rows or columns. Just keep one (or two) permutation vector(s)  $p$  ( $q$ ) with  $p[i]$  ( $q[j]$ ) meaning that the row (column)  $p$  ( $q$ ) has been interchanged with row (column)  $i$  ( $j$ ). This is especially a good solution if matrices are stored row by row or column by column.
3. Gaussian elimination can be used to find the inverse  $A^{-1}$  of a nonsingular matrix. For each  $k = \overline{1, n}$ , column  $k$  of  $A^{-1}$  can be found by solving the system  $Ax = e_k$ , where  $\{e_k\}$  is the canonical basis of  $\mathbb{R}^n$ ,  $e_k = [0 \ 0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0]^T$ , with 1 on the  $k$ th slot. Alternatively, the inverse of  $A$  can be found by Gaussian elimination on the matrix

$$[A \mid I] \sim \dots \sim [I \mid A^{-1}].$$

### Computational Complexity

Since the time required for a certain algorithm to run depends on the details of the hardware used, it is more representative to count the number of some elementary operations, such as multiplica-

tions, divisions, additions, and subtractions. Strictly speaking, we should count separately additions/subtractions and multiplications/divisions, since the latter take slightly more time to be performed, but we will count everything together. Let us assess the computational cost of Gaussian elimination and compare it to other methods.

### *Gaussian elimination*

At step  $k = 1$ , we perform  $n - 1$  divisions,  $(n - 1)n$  multiplications and  $(n - 1)n$  additions, so a total of  $2n(n - 1) + (n - 1)$  flops. At step  $k = 2$ , there are  $2(n - 1)(n - 2) + (n - 2)$  flops and so on until step  $k = n - 1$ . So, the actual elimination process requires

$$\sum_{k=1}^{n-1} [2(n - k)(n - k + 1) + (n - k)] = \sum_{i=1}^{n-1} [2i(i + 1) + i] = \frac{n(n - 1)(4n + 7)}{6}$$

flops. Back substitution adds another

$$1 + 3 + \cdots + 2n - 1 = \sum_{i=1}^{2n-1} i - 2 \sum_{i=1}^{n-1} i = n^2$$

flops, for a total of

$$\frac{n(4n^2 + 9n - 7)}{6} = \mathcal{O}\left(\frac{2}{3}n^3\right)$$

operations.

### *Cramer's rule*

Assume the determinants in Cramer's rule are computed using expansion by minors. That means that to solve an  $n \times n$  system, we have to calculate  $n + 1$  determinants. If  $D_n$  denotes the number of elementary operations needed to compute the determinant of an  $n \times n$  matrix, then

$$\begin{aligned} D_2 &= 2 + 1 = 3 \text{ (two multiplications and one subtraction),} \\ D_3 &= 3D_2 + 3 + 2 \text{ (3}D_2\text{, 3 multiplications and 2 additions/subtractions),} \\ &\dots \\ D_n &= nD_{n-1} + n + n - 1. \end{aligned}$$

So,

$$D_n > nD_{n-1} > n(n - 1)D_{n-2} > \dots > n!.$$

Then the operation count for Cramer's rule is

$$\mathcal{O}((n+1)!).$$

For  $n = 10$ , Gaussian elimination uses about 805 operations, while Cramer's rule uses around 3,628,800 operations. This should emphasize the point that Cramer's rule is not a practical computational method, and that it should be considered as just a theoretical mathematics tool.

## 1.2 Factorization Based Methods

These are methods using the fact that the matrix of coefficients of a linear system being solved can be *factored (decomposed)* into the product of two triangular matrices.

### 1.2.1 LU Factorization

**Theorem 1.7.** *If no row interchanges are necessary in the Gaussian elimination process for solving the system  $Ax = b$ , then  $A$  can be factored as*

$$A = LU, \tag{1.14}$$

where  $L$  and  $U$  are lower and upper triangular matrices, respectively. The pair  $(L, U)$  is called an **LU factorization (decomposition)** of the matrix  $A$ .

*Sketch of Proof.* The first step is to partition  $A$  as

$$A = \left[ \begin{array}{c|ccc} a_{11} & a_{12} & \cdots & a_{1n} \\ \hline a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \cdots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{array} \right] = \begin{bmatrix} a_{11} & w^* \\ v & A' \end{bmatrix},$$

where  $v$  is a column vector of length  $n - 1$ ,  $w^*$  is a row vector of length  $n - 1$  and  $A'$  is an  $(n - 1) \times (n - 1)$  matrix. Then we can factor  $A$  as

$$A = \begin{bmatrix} a_{11} & w^* \\ v & A' \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ v/a_{11} & I_{n-1} \end{bmatrix} \begin{bmatrix} a_{11} & w^* \\ 0 & A' - vw^*/a_{11} \end{bmatrix}.$$

The matrix  $A' - vw^*/a_{11}$  is called the **Schur complement** of  $A$  with respect to  $a_{11}$ .

Then, we proceed recursively:

$$A' - vw^*/a_{11} = L'U'.$$

So

$$\begin{aligned} A &= \begin{bmatrix} 1 & 0 \\ v/a_{11} & I_{n-1} \end{bmatrix} \begin{bmatrix} a_{11} & w^* \\ 0 & A' - vw^*/a_{11} \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 \\ v/a_{11} & I_{n-1} \end{bmatrix} \begin{bmatrix} a_{11} & w^* \\ 0 & L'U' \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 \\ v/a_{11} & L' \end{bmatrix} \begin{bmatrix} a_{11} & w^* \\ 0 & U' \end{bmatrix} \end{aligned}$$

until we get a scalar (a  $1 \times 1$  matrix) that can no longer be partitioned.

□

**Remark 1.8.** If  $A = LU$ , then solving the system  $Ax = b$  is reduced to solving two triangular systems

$$\begin{aligned} Ly &= b \text{ and} \\ Ux &= y. \end{aligned} \tag{1.15}$$

**Example 1.9.** Use  $LU$  decomposition to solve the system

$$\begin{cases} 2x_1 + 4x_2 + 2x_3 = 8 \\ x_1 + x_2 + 2x_3 = 4 \\ x_1 + x_2 + x_3 = 3 \end{cases}$$

**Solution.** We have

$$A = \begin{bmatrix} 2 & 4 & 2 \\ 1 & 1 & 2 \\ 1 & 1 & 1 \end{bmatrix} = \left[ \begin{array}{ccc|ccc} 2 & 4 & 2 & & & \\ 1 & 1 & 2 & & & \\ 1 & 1 & 1 & & & \end{array} \right],$$

so, at the first step,

$$a_{11} = 2, \quad v = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad w^* = [4 \ 2], \quad A' = \begin{bmatrix} 1 & 2 \\ 1 & 1 \end{bmatrix}, \quad \left[ \begin{array}{ccc|ccc} 2 & 4 & 2 & & & \\ 1/2 & & & & & \\ 1/2 & & & & & \end{array} \right].$$

The first Schur complement is

$$A' - vw^*/a_{11} = \begin{bmatrix} 1 & 2 \\ 1 & 1 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix} [4 \ 2] = \begin{bmatrix} 1 & 2 \\ 1 & 1 \end{bmatrix} - \begin{bmatrix} 2 & 1 \\ 2 & 1 \end{bmatrix} = \begin{bmatrix} -1 & 1 \\ -1 & 0 \end{bmatrix}$$

and, for now, we have

$$\left[ \begin{array}{c|cc} 2 & 4 & 2 \\ \hline 1/2 & -1 & 1 \\ \hline 1/2 & -1 & 0 \end{array} \right] = \left[ \begin{array}{c|cc} 2 & 4 & 2 \\ \hline 1/2 & -1 & 1 \\ \hline 1/2 & -1 & 0 \end{array} \right] = \left[ \begin{array}{c|cc} 2 & 4 & 2 \\ \hline 1/2 & -1 & 1 \\ \hline 1/2 & \textcolor{red}{1} & \end{array} \right],$$

where  $\textcolor{red}{1}$  was obtained by dividing  $\frac{-1}{-1}$ .

The last Schur complement is

$$0 - (-1)/(-1) \cdot 1 = -1$$

and the final decomposition is

$$\left[ \begin{array}{c|cc} \textcolor{blue}{2} & \textcolor{blue}{4} & \textcolor{blue}{2} \\ \hline \textcolor{violet}{1/2} & -1 & 1 \\ \hline \textcolor{violet}{1/2} & \textcolor{violet}{1} & -1 \end{array} \right].$$

We take the upper triangular part (**including** the main diagonal) for  $U$  and the lower triangular part (**without** the main diagonal) for  $L$  (which will have all  $\mathbf{1}$ 's on the main diagonal), to get

$$L = \begin{bmatrix} \mathbf{1} & 0 & 0 \\ \textcolor{violet}{1/2} & \mathbf{1} & 0 \\ \textcolor{violet}{1/2} & \textcolor{violet}{1} & \mathbf{1} \end{bmatrix}, \quad U = \begin{bmatrix} \textcolor{blue}{2} & \textcolor{blue}{4} & \textcolor{blue}{2} \\ 0 & -1 & 1 \\ 0 & 0 & -1 \end{bmatrix}$$

and check that indeed  $A = LU$ .

Now, solve  $Ly = b = [8 \ 4 \ 3]^T$ , which, from Example 1.3b) has solution  $y = [8 \ 0 \ -1]^T$  and then  $Ux = [8 \ 0 \ -1]^T$ , which, from Example 1.3a), gives the solution

$$x = [1 \ 1 \ 1]^T.$$

Check that  $Ax = b$ . ■

### Remark 1.10.

**1.** If all that is required is that  $L$  be lower and  $U$  be upper triangular, then the  $LU$  decomposition is *not* unique. We can make it unique by imposing more conditions. For instance, if we require

$l_{ii} = 1, i = \overline{1, n}$ , we have *Doolittle* factorization and if we impose  $u_{ii} = 1, i = \overline{1, n}$ , we get the *Crout* factorization. The procedure described in the proof of Theorem 1.7 leads to Doolittle factorization.

2. The matrix  $U = [u_{ij}]$  in the Doolittle factorization is the upper triangular matrix obtained by Gaussian elimination,

$$u_{ij} = a_{ij}^{(i)}, \quad i \leq j, \quad (1.16)$$

while  $L = [l_{ij}]$  is the matrix of the *multipliers*

$$l_{ij} = m_{ij} = \frac{a_{i,j}^{(j)}}{a_{j,j}^{(j)}}, \quad i \geq j. \quad (1.17)$$

3. Examples of cases when no row interchanges are necessary:

-  $A$  is **diagonally dominant on rows**,

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = \overline{1, n}. \quad (1.18)$$

-  $A$  is **positive definite**,

$$x^T A x > 0, \quad \forall x \neq 0. \quad (1.19)$$

## 1.2 Factorization Based Methods - Continued

### 1.2.2 LUP Factorization

So, we can find an  $LU$  factorization for a matrix  $A$ , whenever row swaps are *not* necessary. What if row interchanges (pivoting) *are* necessary? A row interchange is a permutation of two rows. We keep track of those in a *permutation* matrix, which is simply a matrix obtained from the corresponding identity matrix  $I$  by permuting rows. So, for a matrix  $A$  we find its  **$LUP$  factorization (decomposition)**, i.e., a triplet  $(L, U, P)$ , with  $L$  a lower triangular,  $U$  an upper triangular and  $P$  a permutation matrix, such that

$$PA = LU. \quad (1.1)$$

**Remark 1.1.**

1. Multiplication of a matrix  $A$  to the *left* by a permutation matrix  $P$  will yield the same *row* interchanges on the matrix  $A$  as in  $P$ , while multiplication on the *right* will result in the same *column* interchanges in  $A$  as in  $P$ .
2. Solving the system  $Ax = b$  is now equivalent to solving two triangular systems

$$\begin{aligned} Ly &= Pb \text{ and} \\ Ux &= y. \end{aligned} \quad (1.2)$$

3. The procedure for obtaining an  $LUP$  factorization is similar to the previous one, while keeping track of the row interchanges in a permutation matrix  $P$ .

**Example 1.2.** Find an  $LUP$  factorization for the matrix

$$A = \begin{bmatrix} 2 & 1 & -2 \\ 1 & 1 & -1 \\ 3 & -1 & 1 \end{bmatrix}.$$

**Solution.** At the first step, we do partial pivoting and interchange  $(R_1) \longleftrightarrow (R_3)$ .

At each row interchange, instead of writing the entire matrix  $P$ , we only emphasize which rows are

permuted. Other than that, we proceed as before. We have

$$A \sim \left[ \begin{array}{ccc|ccc} 3 & -1 & 1 & & & \\ 1 & 1 & -1 & & & \\ 2 & 1 & -2 & & & \end{array} \right] = \left[ \begin{array}{ccc|ccc} 3 & -1 & 1 & & & \\ 1 & 1 & -1 & & & \\ 2 & 1 & -2 & & & \end{array} \right] \sim \left[ \begin{array}{ccc|ccc} 3 & -1 & 1 & & & \\ 1/3 & & & & & \\ 2/3 & & & & & \end{array} \right], \quad \left[ \begin{array}{c} 3 \\ 2 \\ 1 \end{array} \right].$$

Schur complement

$$\left[ \begin{array}{cc} 1 & -1 \\ 1 & -2 \end{array} \right] - \frac{1}{3} \left[ \begin{array}{c} 1 \\ 2 \end{array} \right] [-1 \ 1] = \left[ \begin{array}{cc} 1 & -1 \\ 1 & -2 \end{array} \right] - \left[ \begin{array}{cc} -1/3 & 1/3 \\ -2/3 & 2/3 \end{array} \right] = \left[ \begin{array}{cc} 4/3 & -4/3 \\ 5/3 & -8/3 \end{array} \right],$$

so, at this point we have

$$A \sim \left[ \begin{array}{ccc|ccc} 3 & -1 & 1 & & & \\ 1/3 & 4/3 & -4/3 & & & \\ 2/3 & 5/3 & -8/3 & & & \end{array} \right] \sim \left[ \begin{array}{ccc|ccc} 3 & -1 & 1 & & & \\ 2/3 & 5/3 & -8/3 & & & \\ 1/3 & 4/3 & -4/3 & & & \end{array} \right], \quad \left[ \begin{array}{c} 3 \\ 1 \\ 2 \end{array} \right],$$

because we interchanged  $(R_2) \longleftrightarrow (R_3)$ . Further, we have

$$A \sim \left[ \begin{array}{ccc|ccc} 3 & -1 & 1 & & & \\ 2/3 & 5/3 & -8/3 & & & \\ 1/3 & 4/5 & & & & \end{array} \right] \sim \left[ \begin{array}{ccc|ccc} 3 & -1 & 1 & & & \\ 2/3 & 5/3 & -8/3 & & & \\ 1/3 & 4/5 & & 4/5 & & \end{array} \right],$$

the last Schur complement being

$$-\frac{4}{3} - \frac{4}{5} \cdot \left(-\frac{8}{3}\right) = \frac{4}{5}.$$

So, we obtained

$$L = \left[ \begin{array}{ccc} 1 & 0 & 0 \\ 2/3 & 1 & 0 \\ 1/3 & 4/5 & 1 \end{array} \right], \quad U = \left[ \begin{array}{ccc} 3 & -1 & 1 \\ 0 & 5/3 & -8/3 \\ 0 & 0 & 4/5 \end{array} \right], \quad P = \left[ \begin{array}{ccc} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{array} \right].$$

Check that  $PA = LU$ .

■

### Remark 1.3.

1. The computational cost for  $LU$  (and  $LUP$ ) factorization is about the same as for Gaussian



elimination,  $O(n^3)$  flops. However, for *tridiagonal* matrices, that cost drops to  $O(n)$  operations. The *Thomas algorithm*, based on *LUP* decomposition is an efficient way of solving tridiagonal matrix systems. In addition, only three one-dimensional arrays for the three diagonals are needed to store the matrix. This means that very large systems can be solved rapidly and efficiently, and systems of order over  $n = 10,000$  are not unusual in some applications, for example, in solving boundary value problems for differential equations.

**2.** More generally, a *band* or *banded* matrix is a sparse matrix whose non-zero entries are confined to a diagonal band, comprising of the main diagonal and zero or more diagonals on either side. If all matrix elements are zero outside a diagonally bordered band whose range is determined by constants  $k_1, k_2 \geq 0$ ,

$$a_{ij} = 0, \text{ if } j < i - k_1 \text{ or } j > i + k_2$$

then the quantities  $k_1$  and  $k_2$  are called the *lower bandwidth* and *upper bandwidth*, respectively. The *bandwidth* of the matrix is then defined as

$$w = \max \{k_1, k_2\},$$

i.e., it is the number  $w$  such that

$$a_{ij} = 0, \text{ if } |i - j| > w.$$

It can be shown that *LU* factorization with partial pivoting for  $n \times n$  banded matrices with bandwidth  $w$  requires  $O(w^2n)$  flops, while triangular solvers require  $O(wn)$  flops.

### 1.2.3 QR Factorization

**Definition 1.4.** A real square matrix  $Q$  is called *orthogonal* if

$$Q \cdot Q^T = Q^T \cdot Q = I. \quad (1.3)$$

**Theorem 1.5.** Let  $A$  be a real square matrix. Then there exist unique matrices  $Q$  and  $R$  such that

$$A = QR, \quad (1.4)$$

with  $Q$  orthogonal and  $R$  upper triangular with positive elements on the main diagonal,  $r_{ii} > 0, \forall i$ . The pair  $(Q, R)$  is called the **QR factorization** of  $A$ .

**Remark 1.6.**

1. If  $A = QR$ , then solving the system  $Ax = b$  is equivalent to solving the upper triangular systems

$$Rx = Q^T b. \quad (1.5)$$

2. Relation (1.3) automatically implies that any orthogonal matrix is nonsingular with  $Q^{-1} = Q^T$ . Orthogonal matrices are very useful in Numerical Analysis, as they preserve lengths, angles, and do not magnify errors.

**1.2.4 Cholesky Factorization**

**Definition 1.7.** A real square matrix  $A$  is called **positive definite**, if

$$x^T A x = \sum_{i,j=1}^n a_{ij} x_i x_j > 0, \forall x \in \mathbb{R}^n, x \neq 0. \quad (1.6)$$

Symmetric positive definite matrices can be decomposed into triangular factors twice as fast as general matrices. The standard algorithm for this, *Cholesky factorization*, is a variant of Gaussian elimination, which operates both on the left and the right of the matrix at once, preserving and exploiting the symmetry. These matrices have many interesting properties. Among them, the fact that a symmetric matrix is positive definite if and only if all its e-values are real and positive. Also, the e-vectors corresponding to distinct e-values of such a matrix, are orthogonal. Systems having symmetric positive definite matrices play an important role in Numerical Linear Algebra and its applications. Many matrices that arise in physical systems are symmetric and positive definite because of the fundamental physical laws.

**Theorem 1.8.** Let  $A$  be a symmetric positive definite matrix. Then  $A$  has a unique **Cholesky factorization**

$$A = R^T R, \quad (1.7)$$

where  $R$  is an upper triangular matrix with positive elements on the main diagonal,  $r_{ii} > 0, \forall i$ .

*Sketch of Proof.* First off, let us use (1.6) for

$$x = e_1 = [1 \ 0 \ \dots \ 0]^T.$$

We get

$$\begin{aligned}
x^T A x &= [1 \ 0 \ \dots \ 0] \begin{bmatrix} a_{11} & \dots & a_{1n} \\ a_{21} & \dots & a_{2n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \\
&= [1 \ 0 \ \dots \ 0] \begin{bmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{n1} \end{bmatrix} = a_{11}.
\end{aligned}$$

So, any positive definite matrix  $A$  has  $a_{11} > 0$  and we can set  $\alpha = \sqrt{a_{11}}$ . Then we proceed in a similar way as with LU factorization, keeping in mind that  $A$  is also symmetric, so we work on the left and on the right at the same time.

$$\begin{aligned}
A &= \begin{bmatrix} a_{11} & w^T \\ w & A' \end{bmatrix} \\
&= \begin{bmatrix} \alpha & 0 \\ w/\alpha & I_{n-1} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & A' - ww^T/a_{11} \end{bmatrix} \begin{bmatrix} \alpha & w^T/\alpha \\ 0 & I_{n-1} \end{bmatrix} = R_1^T A_1 R_1.
\end{aligned}$$

By induction, all matrices that appear during the factorization are positive definite and so, the process cannot break down. This procedure is repeated until

$$A = \underbrace{R_1^T R_2^T \dots R_n^T}_{R^T} \underbrace{R_n R_{n-1} \dots R_1}_R = R^T R.$$

The uniqueness follows from the fact that at each step, the value  $\alpha = \sqrt{a_{11}}$  is uniquely determined from the factorization and once  $\alpha$  is determined, all the rest of the  $R_i$ 's are also uniquely determined.  $\square$

**Remark 1.9.** This method requires only  $n(n+1)/2$  storage locations for  $R$ , rather than the usual  $n^2$  locations. Since only half the matrix needs to be stored, it follows that half of the arithmetic operations can be avoided and the number of operations is about  $O\left(\frac{1}{3}n^3\right)$ , rather than the number  $O\left(\frac{2}{3}n^3\right)$  required for the usual LU decomposition.

**Example 1.10.** Find the Cholesky factorization (if it exists) of the matrix

$$A = \begin{bmatrix} 4 & 12 & -16 \\ 12 & 37 & -43 \\ -16 & -43 & 98 \end{bmatrix}.$$

**Solution.** The matrix is symmetric and its e-values are

$$0.0188, \quad 15.5040, \quad 123.4772,$$

real and positive. Therefore,  $A$  is positive definite and has a Cholesky decomposition. We will only work on the lower triangular part, the other will follow by symmetry. We have

$$A = \left[ \begin{array}{c|cc} 4 & & \\ \hline 12 & 37 & \\ -16 & -43 & 98 \end{array} \right] \sim \left[ \begin{array}{c|cc} \sqrt{4} & & \\ \hline 6 & & \\ -8 & & \end{array} \right].$$

The first Schur complement is

$$\begin{aligned} A' - ww^T/a_{11} &= \begin{bmatrix} 37 & \\ -43 & 98 \end{bmatrix} - \begin{bmatrix} 6 \\ -8 \end{bmatrix} [6 \quad -8] \\ &= \begin{bmatrix} 37 & \\ -43 & 98 \end{bmatrix} - \begin{bmatrix} 36 & \\ -48 & 64 \end{bmatrix} = \begin{bmatrix} 1 & \\ 5 & 34 \end{bmatrix} \end{aligned}$$

and

$$A \sim \left[ \begin{array}{c|cc} 2 & & \\ \hline 6 & 1 & \\ -8 & 5 & 34 \end{array} \right] \sim \left[ \begin{array}{c|cc} 2 & & \\ \hline 6 & \sqrt{1} & \\ -8 & 5 & \end{array} \right] \sim \left[ \begin{array}{c|cc} 2 & & \\ \hline 6 & 1 & \\ -8 & 5 & \sqrt{9} \end{array} \right] = \begin{bmatrix} 2 & & \\ 6 & 1 & \\ -8 & 5 & 3 \end{bmatrix},$$

with the last Schur complement being  $34 - 5 \cdot 5 = 9$  and its square root 3. Then

$$R^T = \begin{bmatrix} 2 & 0 & 0 \\ 6 & 1 & 0 \\ -8 & 5 & 3 \end{bmatrix}, \quad R = \begin{bmatrix} 2 & 6 & -8 \\ 0 & 1 & 5 \\ 0 & 0 & 3 \end{bmatrix} \quad \text{and} \quad A = R^T R.$$

■

## 2 Iterative Methods

The linear systems  $Ax = b$  that occur in many applications can have very large orders ( $10^3, 10^5, 10^6$ ). For such systems, the Gaussian elimination method (and consequent factorization methods) of the last section is often too expensive in either computation time or computer memory requirements, or possibly both. Moreover, the accumulation of round-off errors can sometimes prevent the numerical solution from being accurate. As an alternative, such linear systems are usually solved with *iteration methods*. In an iterative method, a sequence of progressively accurate iterates is produced to approximate the solution. Thus, in general, we do not expect to get the *exact* solution in a finite number of iteration steps, even if the round-off error effect is not taken into account. In the study of iteration methods, a most important issue is the *convergence property*. We will provide a framework for the convergence analysis of a general iteration method.

### 2.1 Jacobi and Gauss-Seidel Methods

We begin with some numerical examples that illustrate two popular iteration methods. Following that, we give a more general discussion of iteration methods.

Consider the linear system

$$\begin{aligned} 9x_1 + x_2 + x_3 &= b_1 \\ 2x_1 + 10x_2 + 3x_3 &= b_2 \\ 3x_1 + 4x_2 + 11x_3 &= b_3 \end{aligned} \tag{2.1}$$

We proceed as follows: in the equation numbered  $k$ , solve for  $x_k$  in terms of the remaining unknowns. In the above case,

$$\begin{aligned} x_1 &= \frac{1}{9}[b_1 - x_2 - x_3] \\ x_2 &= \frac{1}{10}[b_2 - 2x_1 - 3x_3] \\ x_3 &= \frac{1}{11}[b_3 - 3x_1 - 4x_2] \end{aligned} \tag{2.2}$$

Let

$$x^{(0)} = [x_1^{(0)}, x_2^{(0)}, x_3^{(0)}]^T$$

be an initial guess of the true solution  $x$ . Then define an iteration sequence:

$$\begin{aligned}x_1^{(k+1)} &= \frac{1}{9} \left[ b_1 - x_2^{(k)} - x_3^{(k)} \right] \\x_2^{(k+1)} &= \frac{1}{10} \left[ b_2 - 2x_1^{(k)} - 3x_3^{(k)} \right] \\x_3^{(k+1)} &= \frac{1}{11} \left[ b_3 - 3x_1^{(k)} - 4x_2^{(k)} \right]\end{aligned}\tag{2.3}$$

for  $k = 0, 1, \dots$ . This is called the **Jacobi iteration method** or the *method of simultaneous replacements (substitution)*.

$k$	$x_1^{(k)}$	$x_2^{(k)}$	$x_3^{(k)}$	Error	Ratio
0	0	0	0	$2.00e + 0$	
1	1.1111	1.9000	0	$1.00e + 0$	0.500
2	0.9000	1.6778	-0.9939	$3.22e - 1$	0.322
3	1.0351	2.0182	-0.8556	$1.44e - 1$	0.448
4	0.9819	1.9496	-1.0162	$5.06e - 2$	0.349
5	1.0074	2.0085	-0.9768	$2.32e - 2$	0.462
6	0.9965	1.9915	-1.0051	$8.45e - 3$	0.364
7	1.0015	2.0022	-0.9960	$4.03e - 3$	0.477
8	0.9993	1.9985	-1.0012	$1.51e - 3$	0.375
9	1.0003	2.0005	-0.9993	$7.40e - 4$	0.489
10	0.9999	1.9997	-1.0003	$2.83e - 4$	0.382
30	1.0000	2.0000	-1.0000	$3.01e - 11$	0.447
31	1.0000	2.0000	-1.0000	$1.35e - 11$	0.447

Table 1: Jacobi iteration for solving system (2.1)

In Table 1, we give a number of the iterations for the case that  $b = [10, 19, 0]^T$ , which yields the true solution

$$x = [1, 2, -1]^T.$$

In the table, the error is computed as

$$||x - x^{(k)}|| = \max_{1 \leq i \leq n} |x_i - x_i^{(k)}|.$$

Notice that the errors decrease as  $k$  increases and the values of the ratio eventually approach a

limiting constant of approximately 0.447 as  $k$  becomes much larger.

As another approach to the iterative solution of system (2.1) through the use of (2.2), we use *all* the information we obtain in the calculation of each new component. Specifically, let us define

$$\begin{aligned}x_1^{(k+1)} &= \frac{1}{9} \left[ b_1 - x_2^{(k)} - x_3^{(k)} \right] \\x_2^{(k+1)} &= \frac{1}{10} \left[ b_2 - 2x_1^{(k+1)} - 3x_3^{(k)} \right] \\x_3^{(k+1)} &= \frac{1}{11} \left[ b_3 - 3x_1^{(k+1)} - 4x_2^{(k+1)} \right]\end{aligned}\tag{2.4}$$

for  $k = 0, 1, \dots$ . This is called the **Gauss-Seidel iteration method** or the *method of successive replacements (substitution)*. This method is usually more rapidly convergent than the Jacobi method.

In Table 2, we give a number of iterations for solving the system (2.1). Compare these results to those in Table 1. The speed of convergence is much higher than with the Jacobi method (2.3). The values of the ratio, however, do not appear to approach a limiting value, even when looking at values of  $k$  larger than those in the table.

$k$	$x_1^{(k)}$	$x_2^{(k)}$	$x_3^{(k)}$	Error	Ratio
0	0	0	0	$2.00e + 0$	
1	1.1111	1.6778	-0.9131	$3.22e - 1$	0.161
2	1.0262	1.9687	-0.9958	$3.13e - 2$	0.097
3	1.0030	1.9981	-1.0001	$3.00e - 3$	0.096
4	1.0002	2.0000	-1.0001	$2.24e - 4$	0.074
5	1.0000	2.0000	-1.0000	$1.65e - 5$	0.074
6	1.0000	2.0000	-1.0000	$2.58e - 6$	0.155

Table 2: Gauss-Seidel iteration for solving system (2.1)

## 2.2 Iterative Methods – General Theory

To understand the behavior of iteration methods, it is best to put them into a vector-matrix format. To this end, we recall some notions and results from Linear Algebra.

**Definition 2.1.** Let  $A \in \mathbb{R}^{n \times n}$ .

– The polynomial  $p(\lambda) = \det(A - \lambda I_n)$  is called the **characteristic polynomial of  $A$**  and the equa-

tion  $p(\lambda) = 0$  the **characteristic equation of  $A$** .

– The roots of  $p(\lambda)$  are called **eigenvalues (e-values) of  $A$** .

– If  $\lambda \in \mathbb{C}$  is an e-value of  $A$ , a vector  $x \in \mathbb{R}^n, x \neq 0$  satisfying  $(A - \lambda I_n)x = 0$  is called an **eigenvector (e-vector) of  $A$** , corresponding to the e-value  $\lambda$ .

– The set of all e-values of  $A$ , denoted by  $\lambda(A)$  is called the **spectrum of  $A$** .

– The value  $\rho(A) = \max\{|\lambda| \mid \lambda \in \lambda(A)\}$  is called the **spectral radius of  $A$** .

– The value  $\text{tr}(A) = a_{11} + \cdots + a_{nn}$  is called the **trace of  $A$** .

**Definition 2.2.** A **matrix norm** is a function  $\|\cdot\| : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$  satisfying the conditions:  $\forall A, B \in \mathbb{R}^{n \times n}, \forall \alpha \in \mathbb{R}$ ,

(i)  $\|A\| \geq 0, \|A\| = 0 \Leftrightarrow A = 0_n$ .

(ii)  $\|\alpha A\| = |\alpha| \cdot \|A\|$ .

(iii)  $\|A + B\| \leq \|A\| + \|B\|$ .

(iv)  $\|AB\| \leq \|A\| \cdot \|B\|$ .

The first three conditions define *any* norm on a vector space. The fourth one is *specific* to matrix norms and it is necessary due to the fact that matrix multiplication is not done component-wise.

The easiest way of obtaining a matrix norm is from a vector one.

**Definition 2.3.** Let  $\|\cdot\|$  be a norm on  $\mathbb{R}^n$ . Then

$$\|A\| = \sup_{v \neq 0} \frac{\|Av\|}{\|v\|} = \sup_{\|v\| \leq 1} \|Av\| = \sup_{\|v\|=1} \|Av\| \quad (2.5)$$

is the **natural (subordinate, induced) matrix norm** associated with the vector norm  $\|\cdot\|$ .

**Remark 2.4.**

1. It can be easily checked that (2.5) satisfies the conditions of Definition 2.2 and is indeed a matrix norm.

2. A subordinate matrix norm is just a particular case for the norm of a linear mapping  $A : \mathbb{R}^n \rightarrow \mathbb{R}^n$ .

3. For any induced norm,

$$\|I\| = 1. \quad (2.6)$$



**Theorem 2.5.** Let  $A \in \mathbb{R}^{n \times n}$ . Then

a)

$$\begin{aligned} \|A\|_1 &= \sup_{v \neq 0} \frac{\|Av\|_1}{\|v\|_1} = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}| \quad (\text{the Minkovski norm}), \\ \|A\|_\infty &= \sup_{v \neq 0} \frac{\|Av\|_\infty}{\|v\|_\infty} = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \quad (\text{the Chebyshev norm}), \\ \|A\|_2 &= \sup_{v \neq 0} \frac{\|Av\|_2}{\|v\|_2} = \sqrt{\rho(A^T A)} \quad (\text{the Euclidean norm}). \end{aligned} \quad (2.7)$$

b) The mapping  $\|\cdot\|_F : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^n$  given by

$$\|A\|_F = \left[ \sum_{i=1}^n \sum_{j=1}^n |a_{ij}|^2 \right]^{1/2} = \sqrt{\text{tr}(A^T A)} \quad (2.8)$$

is a nonsubordinate ( $\|I_n\|_F = \sqrt{n}$ ) matrix norm, called the **Frobenius norm**.

Now, to solve the system  $Ax = b$ , for a nonsingular matrix  $A \in \mathbb{R}^{n \times n}$ , suppose there exist  $T \in \mathbb{R}^{n \times n}$  and  $c \in \mathbb{R}^n$ , such that  $I - T$  is invertible and the solution  $x$  of  $Ax = b$  is the unique fixed point of the equation

$$x = Tx + c. \quad (2.9)$$

Let  $x^*$  be the solution and  $x^{(0)}$  be an arbitrary vector (the initial approximation). Then, we use (2.9) to define an iterative method by

$$x^{(k+1)} = Tx^{(k)} + c, \quad k \in \mathbb{N}. \quad (2.10)$$

The matrix  $T$  should be chosen such that the system  $Tx = f$  is “easily solvable” (diagonal, triangular, tridiagonal, etc.)

Regarding the convergence of such methods, we have the following results from Calculus and Linear Algebra:

**Lemma 2.6** (Geometric Series). Let  $X \in \mathbb{R}^{n \times n}$ . If  $\rho(X) < 1$ , then  $(I - X)^{-1}$  exists and

$$(I - X)^{-1} = I + X + \cdots + X^k + \cdots \quad (2.11)$$

Conversely, if the series in (2.11) is convergent, then  $\rho(X) < 1$ .

**Theorem 2.7.** *The following are equivalent:*

- a) *The iteration method (2.10) is convergent;*
- b)  $\rho(T) < 1$ ;
- c)  $\|T\| < 1$  for some matrix norm  $\|\cdot\|$ .

**Theorem 2.8.** *If  $\|T\| < 1$  for some matrix norm  $\|\cdot\|$ , then the sequence  $\{x^{(k)}\}_{k \in \mathbb{N}}$  defined in (2.10) converges to the unique fixed point  $x^*$ , starting with any  $x^{(0)} \in \mathbb{R}^n$ , and the error bounds*

$$\|x^* - x^{(k)}\| \leq \frac{\|T\|}{1 - \|T\|} \|x^{(k)} - x^{(k-1)}\| \quad (2.12)$$

and

$$\|x^* - x^{(k)}\| \leq \frac{\|T\|^k}{1 - \|T\|} \|x^{(1)} - x^{(0)}\| \leq \frac{\|T\|}{1 - \|T\|} \|x^{(1)} - x^{(0)}\|, \quad (2.13)$$

hold for every  $k \in \mathbb{N}^*$ .

**Remark 2.9.**

1. By Theorem 2.8, for a given error  $\varepsilon$ , we compute iterations until

$$\|x^{(k)} - x^{(k-1)}\| \leq \frac{1 - \|T\|}{\|T\|} \varepsilon. \quad (2.14)$$

2. In particular, if  $\|T\| < 1/2$ , then

$$\|x^* - x^{(k)}\| \leq \|x^{(k)} - x^{(k-1)}\|$$

and the stopping criterion can be

$$\|x^{(k)} - x^{(k-1)}\| \leq \varepsilon.$$

Now, *how* to actually find the matrix  $T$  and the scalar  $c$ , satisfying (2.9)? Suppose we can write  $A$  as

$$A = M - N. \quad (2.15)$$

This is called a *splitting* of  $A$ . If  $M$  is easily invertible (diagonal, triangular, etc.), then we can write

$$Ax = b \iff Mx = Nx + b \iff x = M^{-1}Nx + M^{-1}b,$$

which is of the form (2.9), with

$$\begin{aligned} T &= M^{-1}N = M^{-1}(M - A) = I - M^{-1}A, \\ c &= M^{-1}b. \end{aligned}$$

We then define the iteration method by

$$x^{(k+1)} = M^{-1}Nx^{(k)} + M^{-1}b, \quad k \in \mathbb{N}, \quad (2.16)$$

with  $x^{(0)}$  an arbitrary vector.

Assume  $A$  is nonsingular, with  $a_{ii} \neq 0, i = \overline{1, n}$ . We can write

$$A = D - L - U,$$

with

$$D = \begin{bmatrix} a_{11} & & & 0 \\ & a_{22} & & \\ & & \ddots & \\ 0 & & & a_{nn} \end{bmatrix}, \quad -L = \begin{bmatrix} 0 & & & 0 \\ a_{21} & 0 & & \\ \vdots & & \ddots & \\ a_{n1} & a_{n2} & \dots & 0 \end{bmatrix}, \quad -U = \begin{bmatrix} 0 & a_{12} & \dots & a_{1n} \\ & 0 & \dots & a_{2n} \\ & & \ddots & \vdots \\ 0 & & & 0 \end{bmatrix},$$

the diagonal, the lower triangular (without the diagonal) and the upper triangular (without the diagonal) parts of  $A$ .

For **Jacobi iteration**, take

$$\begin{aligned} M &= D, \quad N = L + U, \quad \text{so} \\ T_J &= D^{-1}(L + U), \quad c_J = D^{-1}b. \end{aligned}$$

The method is defined by

$$x^{(k+1)} = D^{-1}(L + U)x^{(k)} + D^{-1}b, \quad k \in \mathbb{N}, \quad x^{(0)} \in \mathbb{R}^n, \quad (2.17)$$

or, component-wise,

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left[ b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j^{(k)} \right], \quad i = \overline{1, n}. \quad (2.18)$$

What can be said about the convergence of the method? By Theorem 2.7, we need a matrix norm such that  $\|T_J\| < 1$ . Using Theorem 2.5, we want

$$\|T_J\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n \left| \frac{a_{ij}}{a_{ii}} \right| < 1,$$

which means

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = \overline{1, n},$$

so, a *diagonally dominant* matrix  $A$ . Thus, for any diagonally dominant system, the Jacobi iterative method converges and the error estimates from Theorem 2.8 can be used. More generally, a necessary and sufficient condition for the convergence of the Jacobi iteration is

$$\rho(T_J) < 1.$$

For **Gauss-Seidel iteration**, we take

$$\begin{aligned} M &= D - L, \quad N = U, \quad \text{so} \\ T_{GS} &= (D - L)^{-1}U, \quad c_{GS} = (D - L)^{-1}b. \end{aligned}$$

Then the method is defined by

$$x^{(k+1)} = (D - L)^{-1}Ux^{(k)} + (D - L)^{-1}b, \quad k \in \mathbb{N}, \quad x^{(0)} \in \mathbb{R}^n, \quad (2.19)$$

and each component, by

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left[ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right], \quad i = \overline{1, n}. \quad (2.20)$$

Although it is not so trivial, it can be shown that for a diagonally dominant matrix,  $\|T_{GS}\| < 1$  and

so the Gauss-Seidel iterative method converges at least as fast as the Jacobi one.

### Acceleration methods; SOR Method

Most iterative methods have a regular pattern in which the error decreases. This can often be used to *accelerate* the convergence. Rather than giving a general theory for the acceleration of iteration methods for solving  $Ax = b$ , we just describe an acceleration of the Gauss-Seidel method. This is one of the main cases of interest in applications.

We introduce an *acceleration parameter*  $\omega$  and consider the following modification of the method:

$$\begin{aligned} M &= \frac{D}{\omega} - L, \quad N = \left( \frac{1-\omega}{\omega} D + U \right), \text{ so} \\ T_\omega &= \left( \frac{D}{\omega} - L \right)^{-1} \left( \frac{1-\omega}{\omega} D + U \right), \quad c_\omega = \left( \frac{D}{\omega} - L \right)^{-1} b. \end{aligned}$$

The acceleration method is defined by

$$x_i^{(k+1)} = \frac{\omega}{a_{ii}} \left[ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right] + (1-\omega) x_i^{(k)}, \quad i = \overline{1, n}. \quad (2.21)$$

This is called the *relaxation method*. We have the following cases:

- $\omega < 1$  is called *subrelaxation*;
- $\omega = 1$  is the Gauss-Seidel method;
- $\omega > 1$  is called *overrelaxation*, the **SOR method**, an abbreviation for *successive overrelaxation*.

It can be shown that, if  $a_{ii} \neq 0$ ,  $i = \overline{1, n}$ , then  $\rho(T_\omega) \geq |\omega - 1|$ . Thus, by Theorem 2.7, a necessary condition for the convergence of the SOR method is

$$0 < \omega < 2. \quad (2.22)$$

Also, the following holds:

**Theorem 2.10 (Ostrowski-Reich).** *If  $A$  is a positive definite matrix and  $0 < \omega < 2$ , then the SOR iteration method converges for any choice of the initial approximation  $x^{(0)} \in \mathbb{R}^n$ .*

The parameter  $\omega$  is to be chosen to minimize the error, in order to make  $x^{(k)}$  converge to  $x$  as rapidly as possible. It was found that the optimal value for  $\omega$  is

$$\omega^* = \frac{2}{1 + \sqrt{1 - (\rho(T_J))^2}}. \quad (2.23)$$

**Remark 2.11.** Iterative methods are rarely used for systems of small order, because they are inefficient, since the time needed to get the desired precision exceeds the time required for Gaussian elimination. But for large systems ( $n \geq 10^3$ ), especially for sparse matrices, they can really make a huge difference in the implementation and computational cost.

### 3 Conditioning of a Linear System

Recall (from Lecture 1) the issue of *stability* (sensitivity to errors/perturbations) and *conditioning* (a measure of that sensitivity) of a mathematical problem.

For a general problem of the type

$$y = f(x), \quad f : \mathbb{R}^m \rightarrow \mathbb{R}^n,$$

we define

$$\gamma_{ij} = (\text{cond}_{ij} f)(x) = \frac{x_i \frac{\partial f_j}{\partial x_i}}{f_j(x)}, \quad i = \overline{1, m}, \quad j = \overline{1, n}$$

and

$$\Gamma(x) = [\gamma_{ij}] = \begin{bmatrix} \frac{x_1 \frac{\partial f_1}{\partial x_1}}{f_1(x)} & \cdots & \frac{x_m \frac{\partial f_1}{\partial x_m}}{f_1(x)} \\ \vdots & & \vdots \\ \frac{x_1 \frac{\partial f_n}{\partial x_1}}{f_n(x)} & \cdots & \frac{x_m \frac{\partial f_n}{\partial x_m}}{f_n(x)} \end{bmatrix}, \quad (3.1)$$

called the **conditioning matrix**. Then, the **condition number** of  $f$  at  $x$  is defined by

$$(\text{cond } f)(x) = \|\Gamma(x)\|, \quad (3.2)$$

for a matrix norm  $\|\cdot\|$ . If  $f$  is a linear function, then

$$(\text{cond } f)(x) = \frac{\|x\| \left\| \frac{\partial f}{\partial x} \right\|}{\|f(x)\|}. \quad (3.3)$$

Now, for a linear system, we have  $A \in \mathbb{R}^{n \times n}$ , nonsingular and  $b \in \mathbb{R}^n$  given. The problem is finding  $x \in \mathbb{R}^n$  such that

$$Ax = b.$$

So, in this case, the input data consists of  $A$  and  $b$  and the output data is the vector  $x$ . Then, we can regard this problem as

$$x = f(b) = A^{-1}b, \quad f : \mathbb{R}^n \rightarrow \mathbb{R}^n. \quad (3.4)$$

Since  $f$  is linear and  $\frac{\partial f}{\partial b} = A^{-1}$ , the condition number is

$$(\text{cond } f)(b) = \frac{\|b\| \|A^{-1}\|}{\|A^{-1}b\|} = \frac{\|Ax\| \|A^{-1}\|}{\|x\|} = \|A^{-1}\| \frac{\|Ax\|}{\|x\|}.$$

Then

$$\max_{\substack{b \in \mathbb{R}^n \\ b \neq 0}} (\text{cond } f)(b) = \|A^{-1}\| \max_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{\|Ax\|}{\|x\|} = \|A^{-1}\| \|A\|.$$

This is the **conditioning number of the matrix**  $A$  (and of the system):

$$\text{cond}(A) = \|A\| \|A^{-1}\|. \quad (3.5)$$

If the matrix  $A$  is singular, by convention,  $\text{cond}(A) = \infty$ .

The number  $\text{cond}(A)$  will vary with the norm being used, but it is always bounded below by one, since

$$1 = \|I\| = \|AA^{-1}\| \leq \|A\| \|A^{-1}\| = \text{cond}(A).$$

If the condition number is nearly 1, then small relative perturbations in  $b$  will lead to similarly small relative perturbations in the solution  $x$ . But if  $\text{cond}(A)$  is large, then there may be small relative perturbations of  $b$  that will lead to large relative perturbations in  $x$ .

**Example 3.1.** (Ill-conditioned Matrices)**1. Hilbert Matrix**

$$H_n = \left[ \frac{1}{i+j-1} \right]_{i,j=\overline{1,n}} = \begin{bmatrix} 1 & \frac{1}{2} & \cdots & \frac{1}{n} \\ \frac{1}{2} & \frac{1}{3} & \cdots & \frac{1}{n+1} \\ \vdots & & \ddots & \vdots \\ \frac{1}{n} & \frac{1}{n+1} & \cdots & \frac{1}{2n-1} \end{bmatrix}. \quad (3.6)$$

This is a symmetric and positive definite matrix, so it is nonsingular. However, it is very ill-conditioned, and increasingly so as  $n$  increases.

$n$	$\text{cond}_2(H_n)$
10	$1.6e + 13$
20	$2.45e + 28$
40	$7.65e + 58$

Table 3: Condition numbers of Hilbert matrix

**2. Vandermonde Matrix**

$$V_n = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ t_1 & t_2 & \cdots & t_n \\ \vdots & \vdots & \ddots & \vdots \\ t_1^{n-1} & t_2^{n-1} & \cdots & t_n^{n-1} \end{bmatrix}. \quad (3.7)$$

For  $t_i = \frac{1}{i}, i = \overline{1, n}$ , it can be shown that

$$\text{cond}_\infty(V_n) > n^{n+1}.$$



# Chapter 3. Approximation of Functions

Approximation of functions is one of the most important tasks in Numerical Analysis.

Most functions encountered in mathematical problems and applications cannot be evaluated exactly, even though we usually handle them as if they were completely known quantities. The simplest and most important of these are  $\sqrt{x}$ ,  $e^x$ ,  $\log x$ , and the trigonometric functions; and there are many other functions that occur commonly in physics, engineering, and other disciplines. In evaluating functions, by hand or using a computer, we are essentially limited to the elementary arithmetic operations  $+$ ,  $-$ ,  $\times$  and  $\div$ . Combining these operations means that we can evaluate polynomials and rational functions, which are polynomials divided by polynomials. All other functions must be evaluated by using approximations based on polynomials or rational functions, including piecewise variants of them (e.g., spline functions). Although rational functions generally give slightly more efficient approximations, polynomials are adequate for most problems and their theory is much easier to handle.

**Interpolation** is the process of finding and evaluating a function whose graph goes through a set of given points. The points may arise as measurements in a physical problem, or they may be obtained from a known function. The interpolating function is usually chosen from a restricted class of functions and *polynomials* are the most commonly used class.

Interpolation is an important tool in producing computable approximations to commonly used functions. Moreover, to numerically integrate or differentiate a function, we often replace the function with a simpler approximating expression, which is then integrated or differentiated. These simpler expressions are almost always obtained by interpolation. Also, some of the most widely used numerical methods for solving differential equations are obtained from interpolating approximations. Finally, interpolation is widely used in computer graphics, to produce smooth curves and surfaces when the geometric object of interest is given at only a discrete set of data points.

Later on, we will briefly consider other forms of function approximations.

## 1 Polynomial Interpolation

**Interpolation problem.** Given  $n + 1$  distinct points – called *nodes (or knots)* –  $x_i \in [a, b]$ ,  $i = \overline{0, n}$  and the values  $f(x_i) = y_i$  of an unknown function  $f : [a, b] \rightarrow \mathbb{R}$ , find a polynomial  $P(x)$  of minimum degree, satisfying

$$P(x_i) = f(x_i), \quad i = \overline{0, n}, \quad (1.1)$$

called *interpolation conditions*. This polynomial approximates function  $f$ .

## 1.1 Lagrange Interpolation

### Linear interpolation

We start with a simple case: consider two interpolation nodes,  $(x_0, y_0), (x_1, y_1), x_0 \neq x_1$ .

We know that there is a unique *line* passing through these points. That means we can find a polynomial of degree 1 that interpolates the data. Let us find it.

The slope of the line is

$$m = \frac{y_1 - y_0}{x_1 - x_0}$$

and its equation is

$$y - y_0 = \frac{y_1 - y_0}{x_1 - x_0}(x - x_0).$$

We find the linear interpolation polynomial as

$$\begin{aligned} P_1(x) &= y_0 + \frac{y_1 - y_0}{x_1 - x_0}(x - x_0) \\ &= \left(1 - \frac{x - x_0}{x_1 - x_0}\right)y_0 + \frac{x - x_0}{x_1 - x_0}y_1 \\ &= \frac{x - x_1}{x_0 - x_1}y_0 + \frac{x - x_0}{x_1 - x_0}y_1. \end{aligned}$$

**Example 1.1.** Consider the function  $f(x) = \sqrt{x}$  and the nodes  $x_0 = 1, x_1 = 4$ , i.e. the data  $(1, 1), (4, 2)$ .

**Solution.** The linear interpolation polynomial is

$$P_1(x) = \frac{x - 4}{1 - 4} \cdot 1 + \frac{x - 1}{4 - 1} \cdot 2 = \frac{1}{3}x + \frac{2}{3}.$$

The graphs of  $f$  and  $P_1$  on the interval  $[0, 15]$  are shown in Figure 1.

■

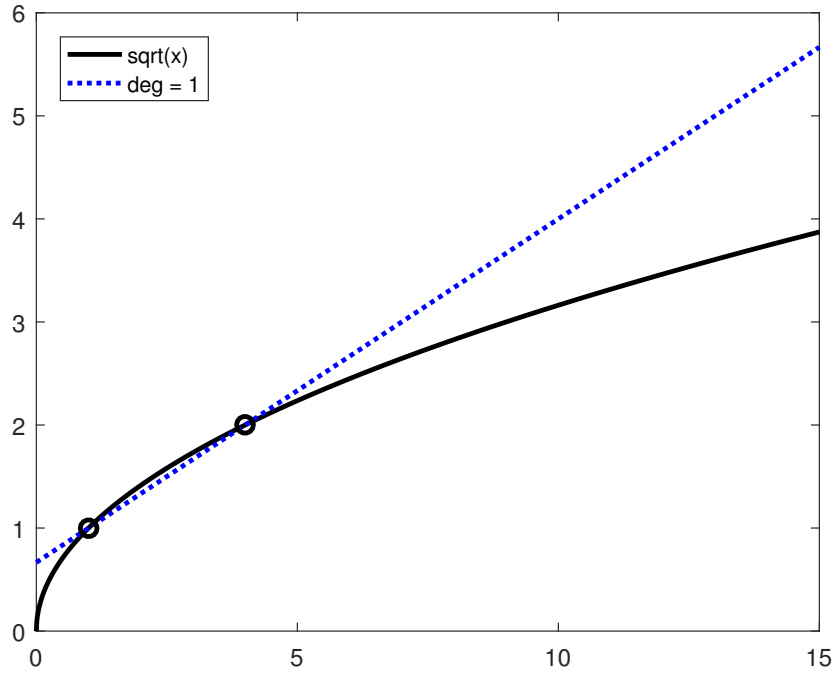


Fig. 1: Linear interpolation of function  $\sqrt{x}$

### Quadratic interpolation

We go further and consider 3 distinct nodes  $(x_0, y_0), (x_1, y_1), (x_2, y_2)$ . It can easily be checked that the quadratic polynomial

$$P_2(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)}y_0 + \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)}y_1 + \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)}y_2$$

interpolates these data.

**Example 1.2.** In Example 1.1 we add the node  $(9, 3)$ .

**Solution.** With nodes  $(1, 1)$ ,  $(4, 2)$  and  $(9, 3)$ , the quadratic interpolation polynomial is given by

$$\begin{aligned} P_2(x) &= \frac{(x - 4)(x - 9)}{(1 - 4)(1 - 9)} \cdot 1 + \frac{(x - 1)(x - 9)}{(4 - 1)(4 - 9)} \cdot 2 + \frac{(x - 1)(x - 4)}{(9 - 1)(9 - 4)} \cdot 3 \\ &= -\frac{1}{60}x^2 + \frac{5}{12}x + \frac{3}{5}. \end{aligned}$$

The graphs of  $f$  and the two interpolation polynomials  $P_1$  and  $P_2$  are plotted in Figure 2. ■

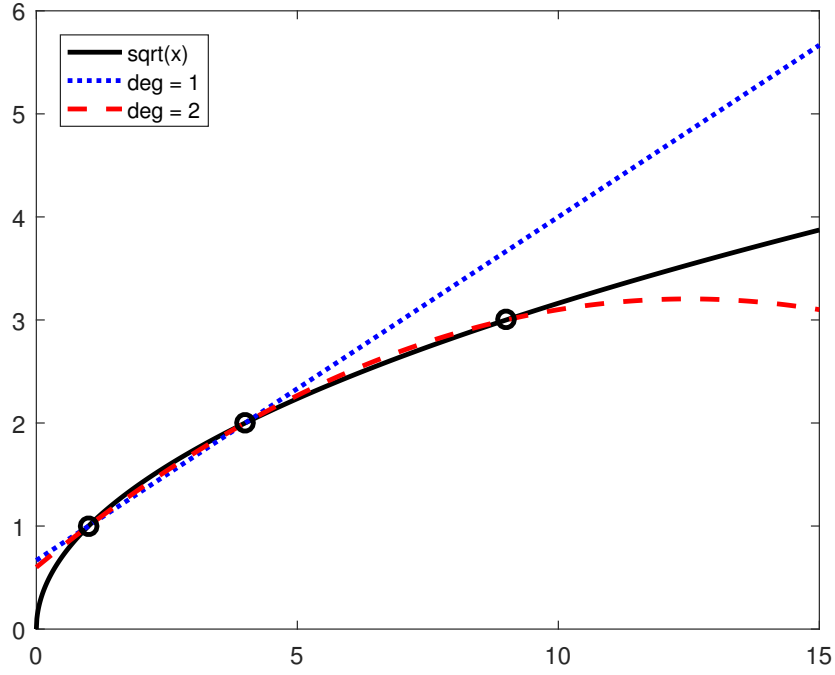


Fig. 2: Linear and quadratic interpolation of function  $\sqrt{x}$

### General case

Consider the interval  $[a, b] \subset \mathbb{R}$ , a function  $f : [a, b] \rightarrow \mathbb{R}$  and a set of  $n + 1$  distinct nodes  $\{x_0, x_1, \dots, x_n\} \subset [a, b]$ .

Recall the notations

$$\begin{aligned} u(x) &= \prod_{j=0}^n (x - x_j), \\ u_j(x) &= \frac{u(x)}{x - x_j}, \quad j = 0, 1, \dots, n. \end{aligned} \tag{1.2}$$

**Theorem 1.3.** *There is a unique polynomial  $L_n f$  of degree at most  $n$ , satisfying the interpolation conditions (1.1). This polynomial can be written as*

$$L_n f(x) = \sum_{i=0}^n l_i(x) f(x_i), \tag{1.3}$$

where

$$l_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j} = \frac{u_i(x)}{u_i(x_i)} = \frac{u_i(x)}{u'(x_i)}. \quad (1.4)$$

$L_n f$  is called the **Lagrange interpolation polynomial** of  $f$  at the nodes  $x_0, x_1, \dots, x_n$ . The functions  $l_i(x), i = \overline{0, n}$  are called **Lagrange fundamental (basis) polynomials** associated with these points.

*Proof.* It can easily be checked that  $l_i$  is a polynomial of degree at most  $n$  and that

$$l_i(x_j) = \delta_{ij} = \begin{cases} 0, & i \neq j \\ 1, & i = j \end{cases}.$$

Hence, the polynomial  $L_n f$  defined in (1.3) is also a polynomial of degree at most  $n$  and it satisfies conditions (1.1).

To prove uniqueness, assume there exists another polynomial  $P_n^*$  (of degree at most  $n$ ) satisfying conditions (1.1) and consider

$$Q_n = L_n - P_n^*.$$

By (1.1),  $Q_n(x_i) = 0, i = 0, \dots, n$ , which means  $Q_n$ , a polynomial of degree at most  $n$ , has  $n + 1$  distinct roots. By the Fundamental Theorem of Algebra,  $Q_n$  must be identically zero, thus proving the uniqueness of  $L_n$ . □

So, given  $n + 1$  distinct points, we can find a unique polynomial of degree *at most*  $n$ , interpolating the data. It is possible that the degree of the interpolation polynomial to be actually *less* than  $n$ .

**Example 1.4.** Find the polynomial of minimum degree that interpolates the data

$$\begin{array}{c|cccccc} x & -2 & -1 & 0 & 1 & 2 & 3 \\ \hline y & -5 & 1 & 1 & 1 & 7 & 25 \end{array}.$$

**Solution.** Given 6 points, we find, by (1.3)-(1.4) (after simplifying), the polynomial

$$L_5 f(x) = \sum_{i=0}^5 l_i(x) y_i = x^3 - x + 1.$$

which, actually, has degree  $3 < 5$ . ■

## Error and convergence

First of all, for any set of distinct nodes, the interpolation problem is *well-posed*, meaning that it has a unique solution that depends continuously on the data. Moreover, it can be expressed in terms of basis polynomials in the form (1.3).

We want to use the approximation

$$f(x) \approx L_n f(x), \quad x \in [a, b].$$

To this end, we must assess (bound) the error (the remainder)

$$(R_n f)(x) = f(x) - (L_n f)(x), \quad x \in [a, b]. \quad (1.5)$$

**Theorem 1.5.** *Let  $[a, b] \subset \mathbb{R}$ ,  $f : [a, b] \rightarrow \mathbb{R}$  a function of class  $C^{n+1}[a, b]$  and consider the distinct nodes  $\{x_0, x_1, \dots, x_n\} \subset [a, b]$ . Then there exists  $\xi \in (a, b)$  such that*

$$(R_n f)(x) = \frac{u(x)}{(n+1)!} f^{(n+1)}(\xi), \quad (1.6)$$

where  $u(x) = (x - x_0)(x - x_1) \dots (x - x_n)$ .

**Example 1.6.** For linear and quadratic interpolation, the remainders are given by

$$\begin{aligned} (R_1 f)(x) &= \frac{(x - x_0)(x - x_1)}{2} f''(\xi), \\ (R_2 f)(x) &= \frac{(x - x_0)(x - x_1)(x - x_2)}{6} f'''(\xi). \end{aligned}$$

For  $f(x) = \sqrt{x} = x^{1/2}$ , the derivatives are

$$\begin{aligned} f'(x) &= \frac{1}{2} x^{-1/2} = \frac{1}{2} \cdot \frac{1}{\sqrt{x}}, \\ f''(x) &= \frac{1}{2} \left( -\frac{1}{2} \right) x^{-3/2} = -\frac{1}{4} \cdot \frac{1}{x\sqrt{x}}, \\ f'''(x) &= -\frac{1}{4} \left( -\frac{3}{2} \right) x^{-5/2} = \frac{3}{8} \cdot \frac{1}{x^2\sqrt{x}}. \end{aligned}$$

So, for the remainders, we have

$$\begin{aligned} |(R_1f)(x)| &= \frac{|(x-x_0)(x-x_1)|}{8} \cdot \frac{1}{\xi\sqrt{\xi}}, \\ |(R_2f)(x)| &= \frac{3}{8} \frac{|(x-x_0)(x-x_1)(x-x_2)|}{6} \cdot \frac{1}{\xi^2\sqrt{\xi}} = \frac{|(x-x_0)(x-x_1)(x-x_2)|}{16} \cdot \frac{1}{\xi^2\sqrt{\xi}}. \end{aligned}$$

**Remark 1.7.** In general, an upper bound of the interpolation error is given by

$$|(R_nf)(x)| \leq \frac{|u(x)|}{(n+1)!} M_{n+1}(f), \quad (1.7)$$

where

$$M_{n+1}(f) = \sup_{t \in [a,b]} |f^{(n+1)}(t)|.$$

Regarding the convergence of the Lagrange polynomial  $L_nf$  to  $f$ , this *does not* happen, in general. The polynomial does converge, if, for instance,  $f \in C^\infty[a, b]$ , with  $|f^{(k)}(x)| \leq M_k$ ,  $\forall x \in [a, b]$ ,  $k = 0, 1, 2, \dots$  and satisfies

$$\lim_{k \rightarrow \infty} \frac{(b-a)^k}{k!} M_k = 0.$$

In the early 1900's, it was proved (by Bernstein and Faber) that for each triangular array of nodes

$$\begin{array}{ccccccc} & & & & & & x_0^{(0)} \\ & & & & & & x_0^{(1)} & x_1^{(1)} \\ & & & & & & x_0^{(2)} & x_1^{(2)} & x_2^{(2)} \\ & & & & & & \vdots & \vdots & \vdots & \ddots \\ & & & & & & x_0^{(n)} & x_1^{(n)} & x_2^{(n)} & \dots & x_n^{(n)} \end{array}$$

in  $[a, b]$  there exists a function  $f \in C[a, b]$ , such that the sequence of Lagrange interpolation polynomials

$$L_nf = L_n(f, x_0^{(n)}, x_1^{(n)}, \dots, x_n^{(n)}; x)$$

*does not converge* uniformly to  $f$  on  $[a, b]$ . Moreover, they proved that for any array of nodes as above, there exists a function  $f \in C[a, b]$ , such that the corresponding sequence  $\{L_nf\}_n$  is *divergent*.

**Example 1.8 (Bernstein's Example).** Let

$$f(x) = |x|, x \in [-1, 1]$$

and consider the equidistant nodes

$$x_k^{(n)} = -1 + \frac{2k}{n}, k = \overline{0, n}.$$

One can show that

$$\lim_{n \rightarrow \infty} |f(x) - L_n f(x)| = \infty,$$

for every  $x \in [-1, 1]$ , except  $x = -1, x = 0$  and  $x = 1$  (see Figure 3). Convergence in  $x = \pm 1$  is trivial, since they are interpolation nodes (where the error is zero). The same is true for  $x = 0$ , when  $n$  is even.

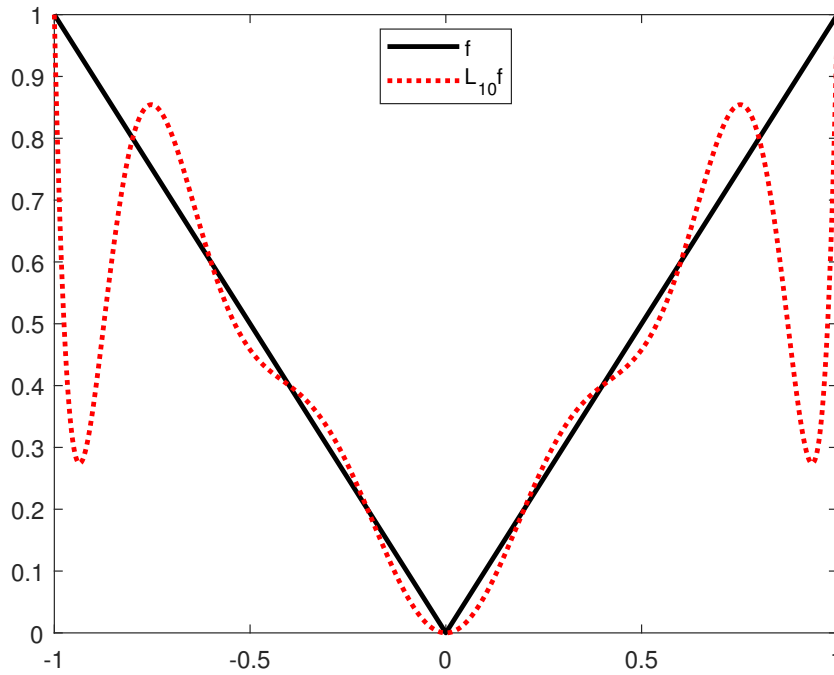


Fig. 3: Bernstein's Example, equidistant nodes,  $n = 10$



**Example 1.9 (Runge's Example).** Consider the function

$$f(x) = \frac{1}{1+x^2}, \quad x \in [-5, 5]$$

and the equally spaced nodes

$$x_k^{(n)} = -5 + 10 \frac{k}{n}, \quad k = \overline{0, n}.$$

It can be shown that

$$\lim_{n \rightarrow \infty} |f(x) - L_n f(x)| = \begin{cases} 0, & \text{if } |x| < 3.633\dots \\ \infty, & \text{if } |x| > 3.633\dots \end{cases}$$

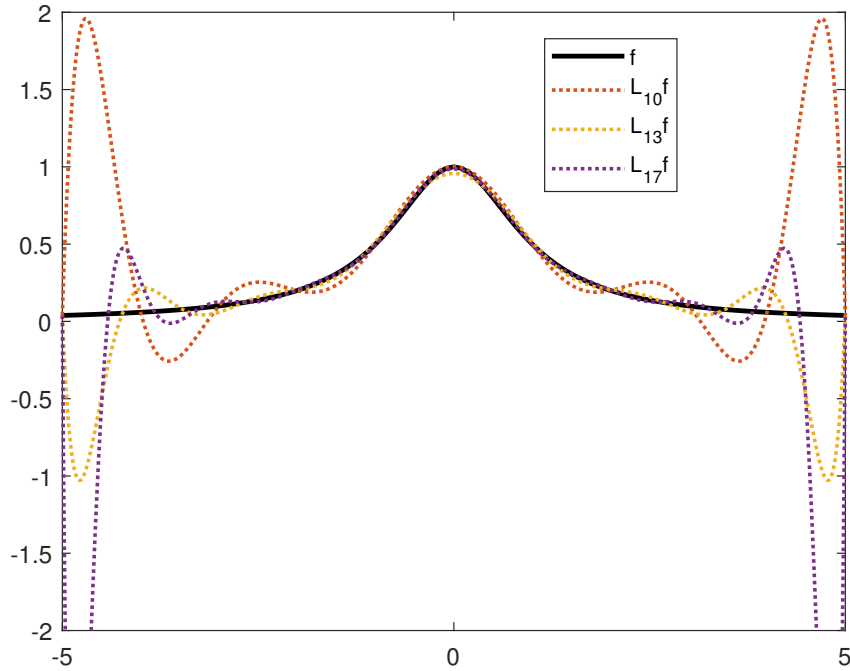


Fig. 4: Runge's Example,  $n = 10, 13, 17$

### Optimal choice of nodes

It may seem that choosing equally spaced nodes is beneficial, because it makes computations easier. However, as the last two examples showed, that is not the case. In fact, polynomial interpolation in equally spaced points is highly *ill-conditioned*: small changes in the data may cause huge changes

in the interpolant. This is known as **Runge's phenomenon**: a problem of oscillation at the edges of an interval, that occurs when using polynomial interpolation with polynomials of high degree over a set of *equidistant* nodes.

For polynomial interpolation to be a well-conditioned process, unless  $n$  is rather small, one must dispense with equally spaced points. The alternative is to use point sets that are clustered at the endpoints of the interval. Such families of nodes will minimize the term  $|u(x)|$  in the error (1.7).

The simplest examples of clustered point sets are the families of *Chebyshev points*, obtained by projecting equally spaced points on the unit circle down to the unit interval  $[-1, 1]$ .

Assume the interval is  $[-1, 1]$ . Then, for a general interval  $[a, b]$ , we use the linear change of variables

$$x = \frac{b-a}{2}t + \frac{b+a}{2}, \quad t \in [-1, 1], \quad x \in [a, b].$$

*Chebyshev points of the first kind*

An optimal choice of nodes are the roots of the **Chebyshev polynomial of the first kind**:

$$T_n(x) = \cos(n \arccos x), \quad x \in [-1, 1]. \quad (1.8)$$

With the change of variables  $x = \cos t, t \in [0, \pi]$ , we get

$$\begin{aligned} T_n(x) &= \cos(nt) = \frac{1}{2} (e^{int} + e^{-int}) \\ &= \frac{1}{2} [(\cos t + i \sin t)^n + (\cos t - i \sin t)^n] \\ &= \frac{1}{2} [(x + i\sqrt{1-x^2})^n + (x - i\sqrt{1-x^2})^n]. \end{aligned}$$

The odd powers of the radical will be canceled, resulting in a polynomial of degree  $n$  in  $x$ , with leading coefficient  $2^{n-1}$ .

Chebyshev polynomials of the first kind have some remarkable properties:

1. Polynomials of degree 0, 1, 2 and 3 are easily computable using trigonometric identities. They

are

$$\begin{aligned} T_0(x) &= 1, \\ T_1(x) &= x, \\ T_2(x) &= 2x^2 - 1, \\ T_3(x) &= 4x^3 - 3x. \end{aligned}$$

2. Higher degree polynomials can be obtained from the recurrence relation

$$T_{k+1}(x) = 2xT_k(x) - T_{k-1}(x), \quad k = 1, 2, \dots$$

For example, the next polynomial is

$$T_4(x) = 8x^4 - 8x^2 + 1.$$

3.  $\{T_n(x)\}_{n \in \mathbb{N}}$  is a sequence of *orthogonal* polynomials on  $(-1, 1)$  with respect to the weight function  $w(x) = \frac{1}{\sqrt{1-x^2}}$ , i.e.

$$\int_{-1}^1 \frac{T_n(x)T_m(x)}{\sqrt{1-x^2}} dx = \begin{cases} 0, & n \neq m \\ \pi, & n = m = 0 \\ \frac{\pi}{2}, & n = m \neq 0 \end{cases}.$$

To minimize the term  $|u(x)|$  in the error (1.7), on the interval  $[-1, 1]$ , we choose

$$u(x) = \tilde{T}_{n+1}(x) = \frac{1}{2^n} T_{n+1}(x),$$

i.e. the nodes

$$x_k = \cos\left(\frac{2k+1}{2n+2}\pi\right), \quad k = 0, \dots, n, \tag{1.9}$$

the roots of the Chebyshev polynomial  $T_{n+1}$ . In this case, we have

$$\|R_n f\| \leq \frac{1}{2^n(n+1)!} \|f^{(n+1)}\|.$$

For the general case, on the interval  $[a, b]$ , we take

$$u(x) = \tilde{T}_{n+1}(x; a, b) = \frac{(b-a)^{n+1}}{2^{2n+1}} \cos \left( (n+1) \arccos \frac{2x-a-b}{b-a} \right)$$

and for the remainder, we have

$$\|R_n f\| \leq \frac{(b-a)^{n+1}}{2^{2n+1}(n+1)!} \|f^{(n+1)}\|.$$

**Example 1.10.** Let us revisit Bernstein's Example, i.e. the function

$$f(x) = |x|, \quad x \in [-1, 1],$$

only with Chebyshev nodes, this time, given by (1.9). Figure 5 shows a much better behaviour of the Lagrange polynomial.

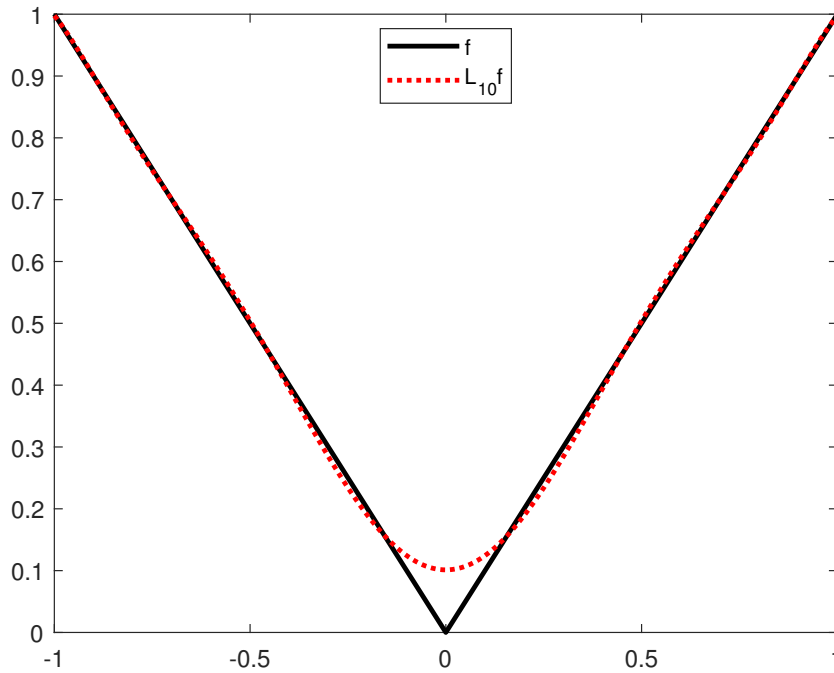


Fig. 5: Bernstein's Example, Chebyshev nodes,  $n = 10$

*Chebyshev points of the second kind*

**Chebyshev polynomials of the second kind** are defined by

$$Q_n(x) = \frac{\sin((n+1)\arccos x)}{\sqrt{1-x^2}}, \quad x \in (-1, 1). \quad (1.10)$$

These polynomials are orthogonal on  $[-1, 1]$  with respect to the weight function  $w(x) = \sqrt{1-x^2}$ :

$$\int_{-1}^1 \sqrt{1-x^2} Q_n(x) Q_m(x) dx = \begin{cases} 0, & n \neq m \\ \frac{\pi}{2}, & n = m \end{cases}.$$

With the same change of variables as before,  $x = \cos t, t \in [0, \pi]$ , we find

$$Q_n(t) = \frac{\sin((n+1)t)}{\sin t}, \quad t \in [0, \pi].$$

The roots of  $Q_n$  are

$$x_k = \cos\left(\frac{k}{n+1}\pi\right), \quad k = 1, \dots, n. \quad (1.11)$$

$Q_n(x)$  can be generated using the recurrence relations

$$\begin{aligned} Q_{k+1}(x) &= 2xQ_k(x) - Q_{k-1}(x), \quad k = 1, 2, \dots \\ Q_0(x) &= 1, \quad Q_1(x) = 2x. \end{aligned}$$

The first few Chebyshev polynomials of the second kind are

$$\begin{aligned} Q_0(x) &= 1, \\ Q_1(x) &= 2x, \\ Q_2(x) &= 4x^2 - 1, \\ Q_3(x) &= 8x^3 - 4x. \end{aligned}$$

**Remark 1.11.** The Chebyshev polynomials of the first and second kinds are closely related. It can

be easily checked that they satisfy a pair of mutual recurrence relations:

$$\begin{aligned} T_{n+1}(x) &= x T_n(x) - (1 - x^2) Q_{n-1}(x), \\ Q_{n+1}(x) &= x Q_n(x) + T_{n+1}(x). \end{aligned}$$

Orthogonal polynomials are used extensively in many problems in Numerical Analysis, so we will revisit them (these and other families) later on.

## 1.2 Efficient Computation of Interpolation Polynomials

Recall the Lagrange interpolation polynomial of a function  $f$ , at the distinct nodes  $(x_i, f_i)$ ,  $f_i = f(x_i)$ ,  $i = \overline{0, n}$ :

$$L_n f(x) = \sum_{i=0}^n l_i(x) f_i, \quad (1.12)$$

where

$$\begin{aligned} l_i(x) &= \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j} = \frac{u_i(x)}{u_i(x_i)} = \frac{u_i(x)}{u'(x_i)}, \\ u(x) &= \prod_{j=0}^n (x - x_j), \quad u_j(x) = \frac{u(x)}{x - x_j}, \quad j = 0, 1, \dots, n. \end{aligned}$$

This formula is well-suited for many theoretical uses of interpolation, but it is less desirable for practical computations. Among its shortcomings:

- for each value of  $x$ , all of the basis functions  $l_i$  must be evaluated at  $x$ , which requires a product of  $n$  terms; thus, the total work is  $O(n^2)$  flops (additions and multiplications) for every value of  $x$ ;
- adding a new node  $(x_{n+1}, f_{n+1})$  requires a new computation *from scratch* and knowing  $L_n f(x)$  does not lead to a less expensive way to evaluate  $L_{n+1} f(x)$ ;
- the computation is numerically unstable.

For these reasons, we need alternative and more easily computable formulations and expressions for interpolation polynomials.

### 1.2.1 Barycentric interpolation

The Lagrange formula (1.12) can be rewritten in such a way that it can be evaluated and updated in  $O(n)$  flops.

$$L_n f(x) = \sum_{i=0}^n \frac{u_i(x)}{u'(x_i)} f_i = \sum_{i=0}^n \frac{u(x)}{(x - x_i)u'(x_i)} f_i = u(x) \sum_{i=0}^n \frac{1}{x - x_i} \frac{u'(x_i)}{u'(x_i)} f_i.$$

Let

$$w_i = \frac{1}{u'(x_i)} = \frac{1}{\prod_{\substack{j=0 \\ j \neq i}}^n (x_i - x_j)}, \quad i = 0, 1, \dots, n. \quad (1.13)$$

These are called **barycentric weights**. With these, the Lagrange interpolation polynomial can be written as

$$L_n f(x) = u(x) \sum_{i=0}^n \frac{w_i}{x - x_i} f_i. \quad (1.14)$$

Formula (1.14) is called the **first barycentric formula** (also known as the *modified Lagrange interpolation formula*).

Now, Lagrange interpolation is a formula requiring  $O(n^2)$  operations for calculating some quantities independent of  $x$ , the weights  $w_i$ , followed by  $O(n)$  flops for evaluating  $L_n f(x)$ , once these numbers are known. Incorporating a new node  $x_{n+1}$  entails two calculations:

- dividing each  $w_i, i = 0, \dots, n$  by  $x_i - x_{n+1}$ , for a cost of  $n + 1$  flops,
- computing a new weight  $w_{n+1}$  using formula (1.13), for another  $n + 1$  flops.

Formula (1.14) can be improved even further. Notice that for the constant function  $f \equiv 1$ , the Lagrange polynomial is  $f$  itself

$$L_n f \equiv 1,$$

by the uniqueness of the interpolation polynomial. Substituting in (1.14), we find

$$1 = u(x) \sum_{i=0}^n \frac{w_i}{x - x_i}$$

and further

$$u(x) = \frac{1}{\sum_{i=0}^n \frac{w_i}{x - x_i}}.$$

Then the Lagrange polynomial can be written as

$$L_n f(x) = \frac{\sum_{i=0}^n \frac{w_i}{x - x_i} f_i}{\sum_{i=0}^n \frac{w_i}{x - x_i}}, \quad (1.15)$$

called the **second barycentric formula** (or, simply, the *barycentric formula*).

**Example 1.12.** Consider our previous example, the function  $f(x) = \sqrt{x}$  and the nodes  $x_0 = 1$  and  $x_1 = 4$ .

**Solution.** We have

$$w_0 = \frac{1}{x_0 - x_1} = \frac{1}{1 - 4} = -\frac{1}{3}, \quad w_1 = \frac{1}{x_1 - x_0} = \frac{1}{4 - 1} = \frac{1}{3}$$

and

$$\begin{aligned} L_1 f(x) &= \frac{\frac{w_0}{x - x_0} f(x_0) + \frac{w_1}{x - x_1} f(x_1)}{\frac{w_0}{x - x_0} + \frac{w_1}{x - x_1}} = \frac{-\frac{1/3}{x - 1} \cdot 1 + \frac{1/3}{x - 4} \cdot 2}{-\frac{1/3}{x - 1} + \frac{1/3}{x - 4}} \\ &= \frac{-\frac{1}{x - 1} + \frac{2}{x - 4}}{-\frac{1}{x - 1} + \frac{1}{x - 4}} = \frac{-x + 4 + 2x - 2}{-x + 4 + x - 1} = \frac{1}{3}x + \frac{2}{3}, \end{aligned}$$

as before. ■

**Remark 1.13.** Even though formula (1.15) *hardly* looks like a polynomial, it *actually* is, as seen in the example above. However, there is one troublesome aspect: the interpolation polynomial should agree with the function value at the nodes, but, it is technically undefined when  $x$  equals one of the



nodes. In fact, it can easily be shown (using L'Hôpital's rule) that

$$\lim_{x \rightarrow x_k} \frac{\sum_{i=0}^n \frac{w_i}{x - x_i} f_i}{\sum_{i=0}^n \frac{w_i}{x - x_i}} = f_k, \quad k = 0, 1, \dots, n,$$

so a continuous extension to the nodes is justified. This aspect is particularly important in the implementation of the barycentric formula.

We see that the barycentric formula is a Lagrange formula, but one with a special and beautiful symmetry. The weights  $w_i$  appear in the denominator exactly as in the numerator, except without the function values  $f_i$ . This means that *any common factor in all the weights  $w_i$  may be canceled without affecting the value of  $L_n f$* , something that will be very useful next.

### Computation of the barycentric weights

For some special sets of nodes  $x_j$ , one can give explicit formulas for the barycentric weights  $w_j$ , using the identity

$$w_j = \frac{1}{u'(x_j)}.$$

- The obvious place to start is *equidistant nodes* with spacing  $h = 2/n$  on the interval  $[-1, 1]$ . In this case,

$$w_j = (-1)^{n-j} \binom{n}{j} / (h^n n!)$$

Since any common factor that does not depend on  $j$  will be canceled in the numerator and denominator of (1.15), this can be simplified to

$$w_j = (-1)^j \binom{n}{j}, \quad j = 0, \dots, n. \quad (1.16)$$

For a general interval  $[a, b]$  we would multiply this formula by  $2^n(b-a)^{-n}$ , but this constant factor too can be dropped, so we end up with (1.16) again, regardless of  $a$  and  $b$ .

- For *Chebyshev points of the first kind*, after canceling factors independent of  $j$ , we find

$$w_j = (-1)^j \sin \frac{(2j+1)\pi}{2n+2}, \quad j = 0, \dots, n. \quad (1.17)$$

If  $n$  is large, the weights  $w_j$  in (1.16) for equispaced barycentric interpolation vary by exponentially large factors, of order approximately  $2^n$ . The effect will be that even small data near the center of the interval are associated with large oscillations in the interpolant, of the order of  $2^n$  times bigger, near the edge of the interval, i.e. Runge's phenomenon.

In contrast, the weights in (1.17) vary by factors  $\mathcal{O}(n)$ , not exponentially, reflecting the good distribution of the points and making polynomial interpolation a well-conditioned problem.

- If *Chebyshev points of the second kind* are used, then the weights are given by

$$w_j = \begin{cases} (-1)^j \frac{1}{2}, & \text{if } j = 0 \text{ or } j = n, \\ (-1)^j, & \text{otherwise.} \end{cases} \quad (1.18)$$

### 1.2.2 Newton-type methods

The next procedures gives an alternative form for the interpolation polynomial, as well as for the remainder.

#### Newton's divided difference formula

To better understand (and write) the transition from  $n$  to  $n + 1$  nodes, we slightly change the notations. For the monic polynomial of the nodes (“monic” means the leading coefficient is 1), previously denoted by “ $u(x)$ ”, we introduce a new notation, one that also emphasizes the *number of nodes* that it refers to. So, let

$$\begin{aligned}\psi_n(x) &= (x - x_0) \dots (x - x_{n-1})(x - x_n), \\ \psi_{n-1}(x) &= (x - x_0) \dots (x - x_{n-1}).\end{aligned}$$

Then we have

$$\begin{aligned}\psi_n(x) &= (x - x_n)\psi_{n-1}(x), \\ \psi'_n(x) &= \psi_{n-1}(x) + (x - x_n)\psi'_{n-1}(x).\end{aligned}$$

Hence,

$$\begin{aligned}\psi'_n(x_i) &= (x_i - x_n)\psi'_{n-1}(x_i), \quad i = 0, \dots, n-1, \\ \psi'_n(x_n) &= \psi_{n-1}(x_n).\end{aligned}\tag{1.1}$$

With these new notations, we can write

$$\begin{aligned}L_{n-1}f(x) &= \sum_{i=0}^{n-1} \frac{\psi_{n-1}(x)}{(x - x_i)\psi'_{n-1}(x_i)} f_i, \\ L_n f(x) &= \sum_{i=0}^n \frac{\psi_n(x)}{(x - x_i)\psi'_n(x_i)} f_i.\end{aligned}\tag{1.2}$$

Let us also recall one of the properties of divided differences (Theorem 5.1 a), in Lecture 2):

$$f[x_0, \dots, x_n] = \sum_{i=0}^n \frac{1}{\psi'_n(x_i)} f_i.\tag{1.3}$$

We want to derive a simple recursive formula from  $L_{n-1}f$  to  $L_nf$ , when adding a new node  $x_n$ . Let

$$Q(x) = L_nf(x) - L_{n-1}f(x). \quad (1.4)$$

Obviously,  $Q$  is a polynomial of degree  $n$  and for  $i = 0, \dots, n-1$ ,

$$Q(x_i) = f(x_i) - f(x_i) = 0,$$

so its  $n$  roots are precisely the nodes  $x_0, \dots, x_{n-1}$ . Then  $Q$  is of the form

$$\begin{aligned} Q(x) &= a_n(x - x_0) \dots (x - x_{n-1}) = a_n\psi_{n-1}(x), \\ Q(x_n) &= a_n\psi_{n-1}(x_n), \end{aligned} \quad (1.5)$$

for some constant  $a_n \in \mathbb{R}$  that we want to find.

On the other hand, the polynomial  $L_nf$  also interpolates  $f$  at the node  $x_n$ , so  $L_nf(x_n) = f(x_n) = f_n$  and, thus,

$$Q(x_n) = f_n - L_{n-1}f(x_n). \quad (1.6)$$

By (1.5)–(1.6), it follows that

$$a_n = \frac{f_n - L_{n-1}f(x_n)}{\psi_{n-1}(x_n)}.$$

Now using (1.1)–(1.3), we get:

$$\begin{aligned} a_n &= \frac{f_n}{\psi_{n-1}(x_n)} - \frac{1}{\psi_{n-1}} L_{n-1}f(x_n) \\ &= \frac{f_n}{\psi_{n-1}(x_n)} - \frac{1}{\psi_{n-1}(x_n)} \sum_{i=0}^{n-1} \frac{\psi_{n-1}(x_n)}{(x_n - x_i)\psi'_{n-1}(x_i)} f_i \\ &= \frac{f_n}{\psi_{n-1}(x_n)} + \sum_{i=0}^{n-1} \frac{f_i}{(x_i - x_n)\psi'_{n-1}(x_i)} \\ &= \frac{f_n}{\psi'_n(x_n)} + \sum_{i=0}^{n-1} \frac{f_i}{\psi'_n(x_i)} \\ &= \sum_{i=0}^n \frac{f_i}{\psi'_n(x_i)} = f[x_0, \dots, x_n]. \end{aligned}$$

Thus,

$$Q(x) = f[x_0, \dots, x_n] \psi_{n-1}(x).$$

So, by (1.4)–(1.6), we have the following recurrence relation for the Lagrange polynomial:

$$L_n f(x) = L_{n-1} f(x) + f[x_0, \dots, x_n] \psi_{n-1}(x), \quad n \geq 1. \quad (1.7)$$

Iteratively, we get

$$\begin{aligned} L_0 f(x) &= f(x_0), \\ L_1 f(x) &= f(x_0) + f[x_0, x_1](x - x_0), \\ L_2 f(x) &= f(x_0) + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1), \\ &\dots \\ L_n f(x) &= f(x_0) + f[x_0, x_1](x - x_0) + \dots + f[x_0, \dots, x_n](x - x_0) \dots (x - x_{n-1}), \end{aligned} \quad (1.8)$$

The expression on the right-hand-side of (1.8) is called **Newton's divided difference form** of the interpolation polynomial, or **Newton's interpolation polynomial** and it is denoted by  $N_n f(x)$ . To be clear, by the uniqueness of the interpolation polynomial at  $n + 1$  distinct nodes, the two polynomials *coincide*,  $L_n f(x) = N_n f(x)$ , they are just expressed (written) in different forms.

If we denote by

$$D_i = f[x_0, \dots, x_i], \quad i \geq 0,$$

Newton's polynomial  $N_n f$  can be written in the *nested form*

$$\begin{aligned} N_n f(x) &= D_0 + (x - x_0)D_1 + (x - x_0)(x - x_1)D_2 + \dots + (x - x_0) \dots (x - x_{n-1})D_n \\ &= D_0 + (x - x_0) \left[ D_1 + (x - x_1) \left[ D_2 + \dots \right. \right. \\ &\quad \left. \left. + (x - x_{n-2}) \left[ D_{n-1} + (x - x_{n-1})D_n \right] \dots \right] \right]. \end{aligned} \quad (1.9)$$

Writing it this way, we can see that the evaluation of  $N_n f(x)$  requires only  $n$  multiplications and  $n$  additions (once the divided differences have been computed), so this is a more computationally efficient formula for the interpolation polynomial.

Next, we also want to express the remainder in a new form. Let  $[a, b]$  denote the smallest interval

containing the distinct nodes  $\{x_0, \dots, x_n\}$  and let  $x \in [a, b]$  be fixed. We write recursively:

$$\begin{aligned}
f[x, x_0] &= \frac{f(x) - f(x_0)}{x - x_0} \\
f[x, x_0, x_1] &= \frac{f[x, x_0] - f[x_0, x_1]}{x - x_1} \\
f[x, x_0, x_1, x_2] &= \frac{f[x, x_0, x_1] - f[x_0, x_1, x_2]}{x - x_2} \\
&\dots \dots \\
f[x, x_0, \dots, x_{n-1}] &= \frac{f[x, x_0, \dots, x_{n-2}] - f[x_0, x_1, \dots, x_{n-1}]}{x - x_{n-1}} \\
f[x, x_0, x_1, \dots, x_n] &= \frac{f[x, x_0, \dots, x_{n-1}] - f[x_0, x_1, \dots, x_n]}{x - x_n}.
\end{aligned} \tag{1.10}$$

Multiplying the first equation in (1.10) by  $(x - x_0)$ , the second by  $(x - x_0)(x - x_1)$ , the third by  $(x - x_0)(x - x_1)(x - x_2)$ , ..., the next to last by  $(x - x_0)(x - x_1) \dots (x - x_{n-1})$  and the last one by  $(x - x_0)(x - x_1) \dots (x - x_n)$ , writing the right-hand-side first, we get

$$\begin{aligned}
f(x) - f(x_0) &= (x - x_0)f[x, x_0] \\
(x - x_0)(f[x, x_0] - f[x_0, x_1]) &= (x - x_0)(x - x_1)f[x, x_0, x_1] \\
(x - x_0)(x - x_1)(f[x, x_0, x_1] - f[x_0, x_1, x_2]) &= (x - x_0)(x - x_1)(x - x_2)f[x, x_0, x_1, x_2] \\
&\dots \dots \\
\prod_{i=0}^{n-2} (x - x_i)(f[x, x_0, \dots, x_{n-2}] - f[x_0, x_1, \dots, x_{n-1}]) &= \prod_{i=0}^{n-1} (x - x_i)f[x, x_0, \dots, x_{n-1}] \\
\prod_{i=0}^{n-1} (x - x_i)(f[x, x_0, \dots, x_{n-1}] - f[x_0, x_1, \dots, x_n]) &= \prod_{i=0}^n (x - x_i)f[x, x_0, x_1, \dots, x_n].
\end{aligned}$$

Now, adding all equations above, we obtain

$$\begin{aligned} f(x) &= f(x_0) + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1) + \dots \\ &+ f[x_0, \dots, x_n](x - x_0) \dots (x - x_{n-1}) + f[x, x_0, \dots, x_n]\psi_n(x) \\ &= N_n f(x) + f[x, x_0, \dots, x_n]\psi_n(x), \end{aligned}$$

from which we have

$$R_n f(x) = f[x, x_0, \dots, x_n](x - x_0) \dots (x - x_n). \quad (1.11)$$

By the mean value formula for divided differences (Theorem 5.1 e), in Lecture 2), we find the previous formula for the remainder:

$$R_n f(x) = \frac{(x - x_0) \dots (x - x_n)}{(n + 1)!} f^{(n+1)}(\xi), \quad \xi \in (a, b).$$

**Example 1.1.** Given the data below, find  $N_1(0.15)$  and  $N_2(0.15)$ , the linear and quadratic interpolates evaluated at  $x = 0.15$ . Determine the remainders.

$i$	$x_i$	$f(x_i)$
0	0.1	0.2
1	0.2	0.24
2	0.3	0.3

**Solution.** First, we compute the divided differences:

$$\begin{array}{l|l} x_0 = 0.1 & f[x_0] = 0.2 \\ & \nearrow \\ x_1 = 0.2 & f[x_1] = 0.24 \\ & \nearrow \\ x_2 = 0.3 & f[x_2] = 0.3 \end{array} \quad \begin{array}{l} \longrightarrow f[x_0, x_1] = \frac{0.24 - 0.2}{0.2 - 0.1} = 0.4 \longrightarrow f[x_0, x_1, x_2] = \frac{0.6 - 0.4}{0.3 - 0.1} = 1 \\ \\ \longrightarrow f[x_1, x_2] = \frac{0.3 - 0.24}{0.3 - 0.2} = 0.6 \end{array}$$

The linear interpolate at the nodes  $x_0 = 0.1$  and  $x_1 = 0.2$  is then

$$N_1 f(x) = f(x_0) + f[x_0, x_1](x - x_0) = 0.2 + 0.4(x - 0.1) = 0.4x + 0.16,$$

so we have the approximation

$$f(0.15) \approx N_1(0.15) = 0.22.$$

The error of this approximation is

$$\begin{aligned} R_1 f(0.15) &= \frac{(0.15 - 0.1)(0.15 - 0.2)}{2!} f''(\xi) \\ &= -1.25 \cdot 10^{-3} f''(\xi), \quad \xi \in (a, b). \end{aligned}$$

Using all three nodes, we find the quadratic interpolate

$$\begin{aligned} N_2 f(x) &= f(x_0) + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1) \\ &= 0.2 + 0.4(x - 0.1) + 1 \cdot (x - 0.1)(x - 0.2) \\ &= x^2 + 0.1x + 0.18, \end{aligned}$$

which yields the approximation

$$f(0.15) \approx N_2(0.15) = 0.2175,$$

with an error of

$$\begin{aligned} R_2 f(0.15) &= \frac{(0.15 - 0.1)(0.15 - 0.2)(0.15 - 0.3)}{3!} f'''(\eta) \\ &= 6.25 \cdot 10^{-5} f'''(\eta), \quad \eta \in (a, b). \end{aligned}$$

■

**Example 1.2.** Find the polynomial of minimum degree that interpolates the data  $f(-1)$ ,  $f(0)$  and  $f(2)$ , for some function  $f \in C^3[-1, 2]$ . Determine and discuss the remainder.

**Solution.** We find the divided differences table:

$$\begin{array}{c|ccc} -1 & f(-1) & \longrightarrow & f(0) - f(-1) & \longrightarrow & \frac{1}{6}(2f(-1) - 3f(0) + f(2)) \\ & & \nearrow & & \nearrow & \\ 0 & f(0) & \longrightarrow & \frac{1}{2}(f(2) - f(0)) & & \\ & & \nearrow & & & \\ 2 & f(2) & & & & \end{array}$$



Then the interpolation polynomial is

$$L_2f(x) = f(-1) + \left(f(0) - f(-1)\right)(x+1) + \frac{1}{6}\left(2f(-1) - 3f(0) + f(2)\right)x(x+1).$$

We can now write it as a linear combination of the given function values,

$$L_2f(x) = \frac{1}{3}x(x-2)f(-1) - \frac{1}{2}(x+1)(x-2)f(0) + \frac{1}{6}x(x+1)f(2).$$

When written this way, it is very easy to check that  $L_2f$  satisfies the interpolation conditions, i.e.

$$L_2f(-1) = f(-1), \quad L_2f(0) = f(0), \quad L_2f(2) = f(2).$$

The coefficients of  $f(-1)$ ,  $f(0)$  and  $f(2)$  above are precisely the basis polynomials  $l_0(x)$ ,  $l_1(x)$  and  $l_2(x)$ , as they satisfy  $l_i(x_j) = \delta_{ij}$ , but the computations were much easier to do this way.

Alternatively, we can write the polynomial in the usual form,

$$L_2f(x) = \frac{1}{6}\left(2f(-1) - 3f(0) + f(2)\right)x^2 + \frac{1}{6}\left(-4f(-1) + 3f(0) + f(2)\right)x + f(0).$$

This second form will be more convenient when we also want to differentiate or integrate the polynomial.

Now, the remainder can be written as

$$R_2f(x) = \frac{u(x)}{3!}f'''(\xi) = \frac{x(x+1)(x-2)}{3!}f'''(\xi), \quad \xi \in (-1, 2).$$

If possible, we may try to find a bound for  $|u(x)|$  on  $[-1, 2]$  (but this may require Matlab ...).

In this case, we have

$$\begin{aligned} u(x) &= x(x+1)(x-2) = x^3 - x^2 - 2x, \\ u'(x) &= 3x^2 - 2x - 2. \end{aligned}$$

The zeros of the derivative are  $\frac{1 \pm \sqrt{7}}{3}$ , the smaller being a point of local maximum and the larger, a point of local minimum for  $u(x)$  on  $[-1, 2]$ . For  $|u(x)|$ , we have:

$$\max_{x \in [0, 2]} |u(x)| = \left| u\left(\frac{1 + \sqrt{7}}{3}\right) \right| = \frac{2}{27}(10 + 7\sqrt{7}).$$

Thus, a bound for the error is

$$|R_2 f(x)| \leq \frac{2}{27 \cdot 3!} (10 + 7\sqrt{7}) |f'''(\xi)| \approx 0.3521 |f'''(\xi)|, \xi \in (-1, 2).$$

■

### Newton's forward and backward difference formula

In the case where the interpolating nodes  $x_i$  are not equally spaced, we use Newton's divided difference formula presented above; however, when the nodes are equidistant, we can construct simpler and less expensive algorithms, using *finite* differences. Historically, these algorithms were of great importance in interpolating functions whose values were given in tables, but the availability of more powerful computers diminished their relevance. However, with new processors (fpu's), they have made a comeback.

Assume the values of a function  $f$  are known at the  $h$ -step equidistant nodes

$$x_i = x_0 + ih, i = 0, 1, \dots$$

Recall the forward differences of the function  $f$

$$\begin{aligned} \Delta^1 f(x_i) &= f(x_i + h) - f(x_i) = f_{i+1} - f_i, \\ \Delta^k f(x_i) &= \Delta^{k-1} f(x_i + h) - \Delta^{k-1} f(x_i) = \Delta^{k-1} f_{i+1} - \Delta^{k-1} f_i \end{aligned}$$

and the property (Proposition 5.6 in Lecture 2)

$$f[x_0, x_0 + h, \dots, x_0 + ih] = \frac{1}{i! h^i} \Delta^i f_0.$$

The Newton form of the  $n$ th degree polynomial  $L_n f$  of  $f$  at the nodes  $x_i = x_0 + ih, i = 0, 1, \dots, n$ , can be simplified. Denote by  $s = (x - x_0)/h$ . Then

$$\begin{aligned} (x - x_0) \dots (x - x_{i-1}) f[x_0, \dots, x_0 + ih] &= (sh) \cdot ((s-1)h) \dots ((s-i+1)h) \frac{1}{i! h^i} \Delta^i f_0 \\ &= \frac{s(s-1) \dots (s-i+1)}{i!} \Delta^i f_0. \end{aligned}$$

Using the notation

$$\binom{s}{k} = \frac{s(s-1)\cdots(s-k+1)}{k!}, \quad s \in \mathbb{R}, k \in \mathbb{N}$$

(the generalized binomial coefficient), we find *Newton's forward difference form* of the interpolation polynomial.

$$\begin{aligned} L_n f(x) &= \sum_{i=0}^n \binom{s}{i} \Delta^i f_0 \\ &= f_0 + \binom{s}{1} \Delta f_0 + \binom{s}{2} \Delta^2 f_0 + \cdots + \binom{s}{n} \Delta^n f_0, \end{aligned} \quad (1.12)$$

with  $s = (x - x_0)/h$ .

The error after  $n$  iterations, for  $x = x_0 + sh$ , is given by

$$f(x) - L_n f(x) = h^{n+1} \binom{s}{n+1} f^{(n+1)}(\xi_x), \quad (1.13)$$

where  $\xi_x$  lies in the smallest interval containing  $x_0, \dots, x_n$  and  $x$ .

Similarly, using backward differences  $\nabla$

$$\begin{aligned} \nabla^0 f_i &= f_i, \\ \nabla^1 f_i &= f_i - f_{i-1}, \\ \nabla^k f_i &= \nabla^{k-1} f_i - \nabla^{k-1} f_{i-1} \end{aligned}$$

and the change of variables  $t = (x - x_n)/h$ , we obtain

$$L_n f(x) = f_n + \frac{t}{1!} \nabla f_n + \frac{t(t+1)}{2!} \nabla^2 f_n + \cdots + \frac{t(t+1)\cdots(t+n-1)}{n!} \nabla^n f_n,$$

which can be written as

$$L_n f(x) = f_n + \binom{t}{1} \nabla f_n + \binom{t+1}{2} \nabla^2 f_n + \cdots + \binom{t+n-1}{n} \nabla^n f_n \quad (1.14)$$

This is called *Newton's backward difference formula*.

In this case, the interpolation error is

$$f(x) - L_n f(x) = h^{n+1} \binom{t+n}{n+1} f^{(n+1)}(\eta_x), \quad (1.15)$$

where  $\eta_x$  lies in the smallest interval containing  $x_0, \dots, x_n$  and  $x$ .

**Example 1.3.** Consider again the data in Example 1.1.

$n$	$x_n$	$f_n$
0	0.1	0.2
1	0.2	0.24
2	0.3	0.3

Let us find  $L_2 f(0.15)$  using finite differences.

**Solution.** By (1.12), we have

$$\begin{aligned} L_2 f(x) &= f_0 + \binom{s}{1} \Delta f_0 + \binom{s}{2} \Delta^2 f_0 \\ &= f_0 + \frac{s}{1!} \Delta f_0 + \frac{s(s-1)}{2!} \Delta^2 f_0 \\ &= f_0 + \frac{x-x_0}{h} \Delta f_0 + \frac{(x-x_0)(x-x_0-h)}{2h^2} \Delta^2 f_0, \end{aligned}$$

where  $s = (x - x_0)/h$ ,  $h = 0.1$ .

We compute the forward differences:

$$\begin{array}{l|l} x_0 = 0.1 & \begin{array}{l} \mathbf{f_0} = 0.2 \longrightarrow \Delta \mathbf{f_0} = 0.24 - 0.2 = 0.04 \longrightarrow \Delta^2 \mathbf{f_0} = 0.06 - 0.04 = 0.02 \\ \nearrow \\ \end{array} \\ x_1 = 0.2 & \begin{array}{l} f_1 = 0.24 \longrightarrow \Delta f_1 = 0.3 - 0.24 = 0.06 \\ \nearrow \\ \end{array} \\ x_2 = 0.3 & f_2 = 0.3 \end{array}$$

So,

$$\begin{aligned} L_2 f(x) &= 0.2 + \frac{x-0.1}{0.1} \cdot 0.04 + \frac{(x-0.1)(x-0.2)}{0.02} \cdot 0.02 \\ &= x^2 + 0.1x + 0.18 \end{aligned}$$

and

$$L_2f(0.15) = 0.2175.$$

Using backward differences, by (1.14), we get

$$\begin{aligned} L_2f(x) &= f_2 + \binom{t}{1} \nabla f_2 + \binom{t+1}{2} \nabla^2 f_2 \\ &= f_2 + \frac{t}{1!} \nabla f_2 + \frac{(t+1)t}{2!} \nabla^2 f_2 \\ &= f_2 + \frac{x-x_2}{h} \nabla f_2 + \frac{(x-x_2+h)(x-x_2)}{2h^2} \nabla^2 f_2 \end{aligned}$$

with  $t = (x - x_2)/h$ ,  $h = 0.1$ .

The backward differences are found in the table

$x_0 = 0.1$	$f_0 = 0.2$			
$x_1 = 0.2$	$f_1 = 0.24$	$\searrow$		
		$\longrightarrow$	$\nabla f_1 = 0.24 - 0.2 = 0.04$	
$x_2 = 0.3$	$f_2 = 0.3$	$\searrow$		$\searrow$
		$\longrightarrow$	$\nabla f_2 = 0.3 - 0.24 = 0.06$	$\longrightarrow$
				$\nabla^2 f_2 = 0.06 - 0.04 = 0.02$

Hence,

$$\begin{aligned} L_2f(x) &= 0.3 + \frac{x-0.3}{0.1} \cdot 0.06 + \frac{(x-0.2)(x-0.3)}{0.02} \cdot 0.02 \\ &= x^2 + 0.1x + 0.18 \end{aligned}$$

and

$$L_2f(0.15) = 0.2175.$$

■

**Remark 1.4.** Interpolation algorithms can be classified according to the “step” of the grid (the distance between two consecutive nodes, when sorted in increasing order). There are *variable step* methods (the Lagrange form with fundamental polynomials, the barycentric formulas, Newton’s divided difference formula) and *constant step* algorithms (Newton’s forward and backward formulas). For variable step methods the precision is the same at any intermediate value in the interval covered by the data  $(x_i, f_i)$ . So these methods do not have so-called *preferential precision zones*. In contrast,

Newton's forward formula is particularly useful (i.e. it has higher precision) for interpolating the values of  $f(x)$  near the beginning of the set of values (closer to the first node,  $(x_0, f_0)$ ), whereas the backward formula is preferred when the value of  $f(x)$  is required near the end of the table (in the vicinity of the last node,  $(x_n, f_n)$ ).

### 1.2.3 Aitken-type methods

These are variable step iterative methods and they highlight another important aspect: in many cases, the degree required to attain a certain desired accuracy in polynomial interpolation is *not known*. It can be obtained from the remainder, but that assumes knowledge (or at least knowing a bound) of  $\|f^{(n+1)}\|_\infty$ .

The idea behind these methods is to write an interpolation polynomial of degree  $n$ , iteratively, in terms of two interpolation polynomials of degree  $n - 1$ , that only use a part of the  $n + 1$  nodes. Let us illustrate the idea for a simple case. For two nodes,  $x_0$  and  $x_1$ , the polynomial of degree 1 interpolating these data, can be written successively (using Lagrange basis polynomials) as

$$\begin{aligned} P_{01}(x) &= l_0(x)f_0 + l_1(x)f_1 \\ &= \frac{x - x_1}{x_0 - x_1}f_0 + \frac{x - x_0}{x_1 - x_0}f_1 \\ &= \frac{(x - x_0)f_1 - (x - x_1)f_0}{x_1 - x_0} \\ &= \frac{(x - x_0)P_1(x) - (x - x_1)P_0(x)}{x_1 - x_0}, \end{aligned}$$

where  $P_0$  denotes the polynomial that interpolates  $f$  at the node  $x_0$  (a polynomial of degree 0, hence, a constant,  $f_0$ ),  $P_1$ , the polynomial of degree 0 that interpolates  $f$  at the node  $x_1$  (identically equal to  $f_1$ ), and  $P_{01}$  the polynomial of degree 1 that interpolates  $f$  at the nodes  $x_0, x_1$ . Similarly, if we add another node,  $x_2$ , we can define

$$P_{12}(x) = \frac{(x - x_1)P_2(x) - (x - x_2)P_1(x)}{x_2 - x_1},$$

which is the polynomial of degree 1 that interpolates  $f$  at the nodes  $x_1, x_2$ . We proceed further and define

$$P_{012}(x) = \frac{(x - x_0)P_{12}(x) - (x - x_2)P_{01}(x)}{x_2 - x_0}. \quad (1.16)$$

Let us compute its values at the nodes.

$$\begin{aligned}
P_{012}(x_0) &= \frac{0 - (x_0 - x_2)P_{01}(x_0)}{x_2 - x_0} = P_{01}(x_0) = f_0, \\
P_{012}(x_1) &= \frac{(x_1 - x_0)P_{12}(x_1) - (x_1 - x_2)P_{01}(x_1)}{x_2 - x_0} = \frac{(x_1 - x_0)f_1 - (x_1 - x_2)f_1}{x_2 - x_0} = f_1, \\
P_{012}(x_2) &= \frac{(x_2 - x_0)P_{12}(x_2) - 0}{x_2 - x_0} = P_{12}(x_2) = f_2.
\end{aligned}$$

Since  $P_{012}$  is a polynomial of degree 2 and it interpolates  $f$  at the nodes  $x_0, x_1, x_2$ , it follows by the uniqueness of the Lagrange interpolation polynomial, that  $P_{012} = L_2 f$ .

In a similar fashion, we can construct recursively the polynomials

$$\begin{aligned}
P_{123}(x) &= \frac{(x - x_1)P_{23}(x) - (x - x_3)P_{12}(x)}{x_3 - x_1}, \\
P_{0123}(x) &= \frac{(x - x_0)P_{123}(x) - (x - x_3)P_{012}(x)}{x_3 - x_0} \\
&\dots
\end{aligned}$$

**Proposition 1.5.** *Let  $x_0, \dots, x_k$  be distinct nodes and let  $f_i, i = 0, \dots, k$ , be the values of a function  $f$  at the nodes. Then the Lagrange polynomial interpolating  $f$  at these nodes is given by*

$$\begin{aligned}
P_{01\dots k}(x) &= \frac{1}{x_k - x_0} \begin{vmatrix} x - x_0 & P_{01\dots k-1}(x) \\ x - x_k & P_{12\dots k}(x) \end{vmatrix} \\
&= \frac{(x - x_0)P_{12\dots k}(x) - (x - x_k)P_{01\dots k-1}(x)}{x_k - x_0}.
\end{aligned} \tag{1.17}$$

*Proof.* Obviously, by its construction, the polynomial in (1.17) has degree  $k$ . Its values at the nodes are

$$\begin{aligned}
P_{01\dots k}(x_0) &= \frac{-(x_0 - x_k)P_{01\dots k-1}(x_0)}{x_k - x_0} = P_{01\dots k-1}(x_0) = f_0, \\
P_{01\dots k}(x_j) &= \frac{(x_j - x_0)P_{12\dots k}(x_j) - (x_j - x_k)P_{01\dots k-1}(x_j)}{x_k - x_0} = f_j, \quad j = \overline{1, k-1}, \\
P_{01\dots k}(x_k) &= \frac{(x_k - x_0)P_{12\dots k}(x_k)}{x_k - x_0} = P_{12\dots k}(x_k) = f_k.
\end{aligned}$$

Hence, by the uniqueness of the Lagrange interpolation polynomial, it follows that  $P_{01\dots k} = L_k f$ . □

Thus, we established a recurrence relation between a Lagrange interpolation polynomial of de-

gree  $k$  and two Lagrange interpolation polynomials of degree  $k - 1$ . The computations can be organized in a table, illustrated below for 4 nodes.

$$\begin{array}{cccccc} x_0 & P_0 & & & & \\ x_1 & P_1 & P_{01} & & & \\ x_2 & P_2 & P_{12} & P_{012} & & \\ x_3 & P_3 & P_{23} & P_{123} & P_{0123} & \end{array}$$

Now, if, for instance,  $P_{0123}$  does not provide a desired approximation precision, we can consider a new node and add a new line to the table:

$$\begin{array}{cccccc} x_4 & P_4 & P_{34} & P_{234} & P_{1234} & P_{01234} \end{array}$$

and we can compare neighboring elements on a row, column or diagonal to check if the desired accuracy has been achieved.

The method described above is called *Neville's method*.

The notations can be simplified. We denote now the polynomials above by  $\tilde{P}$  (instead of  $P$ ) and define the new polynomials  $P$  as follows:

$$P_{i,j} = \tilde{P}_{i-j, i-j+1, \dots, i-1, i}, \quad j = i, i-1, \dots, 0,$$

i.e., recursively,

$$\begin{aligned} P_{i,0} &:= f(x_i), \quad i = \overline{0, n}, \\ P_{i,j} &:= \frac{(x - x_{i-j})P_{i,j-1} - (x - x_i)P_{i-1,j-1}}{x_i - x_{i-j}} \\ &= \frac{1}{x_i - x_{i-j}} \begin{vmatrix} x - x_{i-j} & P_{i-1,j-1} \\ x - x_i & P_{i,j-1} \end{vmatrix}, \quad i \geq j > 0. \end{aligned} \tag{1.18}$$

We get a new table

$$\begin{array}{cccccc} x_0 & P_{00} & & & & \\ x_1 & P_{10} & P_{11} & & & \\ x_2 & P_{20} & P_{21} & P_{22} & & \\ x_3 & P_{30} & P_{31} & P_{32} & P_{33} & \end{array}$$

and the Lagrange polynomial will be the one on the diagonal  $L_n f = P_{nn}$ .

If the interpolation converges, then the sequence  $\{P_{ii}\}_{i \geq 0}$  also converges and we can use the stopping criterion

$$|P_{ii} - P_{i-1, i-1}| < \varepsilon.$$



Aitken's method is similar to Neville's method. We construct the table

$$\begin{array}{cccccc} x_0 & P_{00} & & & & \\ x_1 & P_{10} & P_{11} & & & \\ x_2 & P_{20} & P_{21} & P_{22} & & \\ x_3 & P_{30} & P_{31} & P_{32} & P_{33} & \end{array}$$

defining recursively

$$\begin{aligned} P_{i,0} &:= f(x_i), \quad i = \overline{0, n}, \\ P_{i,j+1} &:= \frac{1}{x_i - x_j} \left| \begin{array}{cc} x - x_j & P_{j,j} \\ x - x_i & P_{i,j} \end{array} \right| = \frac{(x - x_j)P_{i,j} - (x - x_i)P_{j,j}}{x_i - x_j}, \quad i > j \geq 0. \end{aligned} \quad (1.19)$$

**Example 1.6.** Approximate  $\sqrt{2}$  interpolating the function  $f(x) = 2^x$  at the nodes  $-1, 0, 1$ , and then at the nodes  $-1, 0, 1, 2$ .

**Solution.**

With Neville's method, we have the table

$$\begin{array}{cccccc} x_0 = -1 & P_{00} = 1/2 & & & & \\ x_1 = 0 & P_{10} = 1 & P_{11} = 5/4 & & & \\ x_2 = 1 & P_{20} = 2 & P_{21} = 3/2 & P_{22} = 23/16, & & \end{array}$$

where, for  $x = 1/2$ ,

$$\begin{aligned} P_{11} &= \frac{(x - x_0)P_{10} - (x - x_1)P_{00}}{x_1 - x_0} = \frac{(1/2 - (-1)) \cdot 1 - (1/2 - 0) \cdot 1/2}{0 - (-1)} = 5/4, \\ P_{21} &= \frac{(x - x_1)P_{20} - (x - x_2)P_{10}}{x_2 - x_1} = \frac{(1/2 - 0) \cdot 2 - (1/2 - 1) \cdot 1}{1 - 0} = 3/2, \\ P_{22} &= \frac{(x - x_0)P_{21} - (x - x_2)P_{11}}{x_2 - x_0} = \frac{(1/2 - (-1)) \cdot 3/2 - (1/2 - 1) \cdot 5/4}{1 - (-1)} = 23/16. \end{aligned}$$

Thus, with linear interpolation, we get the approximation

$$\sqrt{2} \approx 23/16 = 1.4375$$

and

$$|P_{22} - P_{11}| = 3/16 = 0.1875.$$

We add a new node  $x_3 = 2$  and a new line to the table, to get

$$\begin{array}{llllll} x_0 = -1 & P_{00} = 1/2 & & & & \\ x_1 = 0 & P_{10} = 1 & P_{11} = 5/4 & & & \\ x_2 = 1 & P_{20} = 2 & P_{21} = 3/2 & P_{22} = 23/16 & & \\ x_3 = 2 & P_{30} = 4 & P_{31} = 1 & P_{32} = 11/8 & P_{33} = 45/32, & \end{array}$$

with

$$\begin{aligned} P_{31} &= \frac{(x - x_2)P_{30} - (x - x_3)P_{20}}{x_3 - x_2} = \frac{(1/2 - 1) \cdot 4 - (1/2 - 2) \cdot 2}{2 - 1} = 1, \\ P_{32} &= \frac{(x - x_1)P_{31} - (x - x_3)P_{21}}{x_3 - x_1} = \frac{(1/2 - 0) \cdot 1 - (1/2 - 2) \cdot 3/2}{2 - 0} = 11/8 \\ P_{33} &= \frac{(x - x_0)P_{32} - (x - x_3)P_{22}}{x_3 - x_0} = \frac{(1/2 - (-1)) \cdot 11/8 - (1/2 - 2) \cdot 23/16}{2 - (-1)} = 45/32. \end{aligned}$$

The new approximation (using quadratic interpolation) is

$$\sqrt{2} \approx 45/32 = 1.4063,$$

with

$$|P_{33} - P_{22}| = 1/32 = 0.0313.$$

Let us note that the exact value of  $\sqrt{2}$  rounded to 4 correct decimals is 1.4142, so the actual errors of the two approximations are

$$|\sqrt{2} - P_{22}| = 0.0233 \text{ and } |\sqrt{2} - P_{33}| = 0.0079.$$

With Aitken's algorithm, (1.19), we construct the table

$$\begin{array}{llllll} x_0 = -1 & P_{00} = 1/2 & & & & \\ x_1 = 0 & P_{10} = 1 & P_{11} = 5/4 & & & \\ x_2 = 1 & P_{20} = 2 & P_{21} = 13/8 & P_{22} = 23/16 & & \\ x_3 = 2 & P_{30} = 4 & P_{31} = 9/4 & P_{32} = 3/2 & P_{33} = 45/32, & \end{array}$$

where

$$\begin{aligned}
P_{21} &= \frac{(x - x_0)P_{20} - (x - x_2)P_{00}}{x_2 - x_0} = \frac{(1/2 - (-1)) \cdot 2 - (1/2 - 1) \cdot 1/2}{1 - (-1)} = 13/8, \\
P_{22} &= \frac{(x - x_1)P_{21} - (x - x_2)P_{11}}{x_2 - x_1} = \frac{(1/2 - 0) \cdot 13/8 - (1/2 - 1) \cdot 5/4}{1 - 0} = 23/16, \\
\\
P_{31} &= \frac{(x - x_0)P_{30} - (x - x_3)P_{00}}{x_3 - x_0} = \frac{(1/2 - (-1)) \cdot 4 - (1/2 - 2) \cdot 1/2}{2 - (-1)} = 9/4, \\
P_{32} &= \frac{(x - x_1)P_{31} - (x - x_3)P_{11}}{x_3 - x_1} = \frac{(1/2 - 0) \cdot 9/4 - (1/2 - 2) \cdot 5/4}{2 - 0} = 3/2, \\
P_{33} &= \frac{(x - x_2)P_{32} - (x - x_3)P_{22}}{x_3 - x_2} = \frac{(1/2 - 1) \cdot 3/2 - (1/2 - 2) \cdot 23/16}{2 - 1} = 45/32.
\end{aligned}$$

The two algorithms are very similar and they actually yield the same values on the main diagonal of the table (the values  $P_{nn} = L_n f$ ), so the errors of linear/quadratic interpolates are the same as before.

■

### 1.3 Hermite Interpolation

Consider the following situation: For a moving object, we know the distances traveled  $d_0, d_1, \dots, d_m$ , at times  $t_0, t_1, \dots, t_m$ , and we want a polynomial approximation of the distance function  $d = d(t)$  on the entire interval containing the points  $t_0, \dots, t_m$ . Obviously, this is a Lagrange interpolation problem and we already know how to find the interpolation polynomial.

Now, assume that, in addition, we also know the values of the *velocities*  $v_i$  of the object at times  $t_i, i = \overline{0, m}$ . We would expect that this additional information helps us find an *even better* approximation of the function  $d$ . However, from what we know about Lagrange interpolation, there is *no way* to include this data into our approximation. Since the velocity is the derivative with respect to time of the distance traveled, this means that we also have information about the *derivatives* of the function we want to interpolate. This is a **Hermite interpolation** problem. The ideas and computational formulas are similar to the ones we used to determine the Lagrange interpolation polynomial.

#### 1.3.1 Interpolation with double nodes

For a variety of applications, as the one described above, it is convenient to consider polynomials  $P(x)$  that interpolate a function  $f(x)$  and in addition have the derivative polynomial  $P'(x)$  also interpolate the derivative function  $f'(x)$ .

**Hermite interpolation problem with double nodes.** Given  $m + 1$  distinct nodes  $x_i, i = \overline{0, m}$  and the values  $f(x_i), f'(x_i)$  of an unknown function  $f$  and its derivative, find a polynomial  $P(x)$  of minimum degree, satisfying the interpolation conditions

$$\begin{aligned} P(x_i) &= f(x_i), \\ P'(x_i) &= f'(x_i), \quad i = \overline{0, m}. \end{aligned} \tag{1.1}$$

Since for each node there are two values (of the function and of its derivative) given, we call them *double nodes*.

There are  $2m + 2$  conditions in (1.1), so we seek a polynomial of degree (at most)  $n = 2m + 1$ . We determine this polynomial in a similar way to the construction of the Lagrange polynomial.

Recall the notations:

$$\begin{aligned} \psi_m(x) &= (x - x_0) \dots (x - x_{m-1})(x - x_m), \\ l_i(x) &= \frac{(x - x_0) \dots (x - x_{i-1})(x - x_{i+1}) \dots (x - x_m)}{(x_i - x_0) \dots (x_i - x_{i-1})(x_i - x_{i+1}) \dots (x_i - x_m)} = \frac{\psi_m(x)}{(x - x_i)\psi'_m(x_i)}, \end{aligned} \tag{1.2}$$

for  $i = 0, 1, \dots, m$ .

**Theorem 1.1.** *There is a unique polynomial  $H_n f$  of degree at most  $n$ , satisfying the interpolation conditions (1.1). This polynomial can be written as*

$$H_n f(x) = \sum_{i=0}^m \left[ h_{i0}(x) f(x_i) + h_{i1}(x) f'(x_i) \right], \quad (1.3)$$

where

$$\begin{aligned} h_{i0}(x) &= [1 - 2l'_i(x_i)(x - x_i)] [l_i(x)]^2, \\ h_{i1}(x) &= (x - x_i) [l_i(x)]^2, \quad i = 0, \dots, m. \end{aligned} \quad (1.4)$$

$H_n f$  is called the **Hermite interpolation polynomial** of  $f$  at the double nodes  $x_0, x_1, \dots, x_m$ . The functions  $h_{i0}(x), h_{i1}(x)$ ,  $i = \overline{0, m}$  are called **Hermite fundamental (basis) polynomials** associated with these points.

*Proof.* First we will prove that the polynomial in (1.3) does satisfy all interpolation conditions (i.e., existence), and then we will show that it is the only one to do so (i.e., uniqueness).

The degree of polynomials  $l_i$  from (1.2) is  $m$ , so the degree of  $h_{i0}, h_{i1}$  and  $H_n f$  is  $2m + 1 = n$ .

The derivatives of the Hermite fundamental polynomials are

$$\begin{aligned} h'_{i0}(x) &= -2l'_i(x_i)(l_i(x))^2 + 2[1 - 2l'_i(x_i)(x - x_i)]l'_i(x)l_i(x), \\ h'_{i1}(x) &= (l_i(x))^2 + 2(x - x_i)l'_i(x)l_i(x). \end{aligned}$$

Notice that  $l_i(x)$ ,  $i = \overline{0, m}$  are the Lagrange fundamental polynomials, thus,

$$l_i(x_j) = \delta_{ij} = \begin{cases} 0, & i \neq j \\ 1, & i = j \end{cases}.$$

Then,

$$\begin{aligned} h_{i0}(x_j) &= 0, \quad j \neq i, \\ h_{i0}(x_i) &= 1 \cdot (l_i(x_i))^2 = 1, \\ h_{i1}(x_j) &= 0, \quad j \neq i, \\ h_{i1}(x_i) &= 0. \end{aligned}$$

The values of the derivatives at the nodes are

$$\begin{aligned} h'_{i0}(x_j) &= 0, \quad j \neq i, \\ h'_{i0}(x_i) &= -2l'_i(x_i) + 2l'_i(x_i) = 0, \\ h'_{i1}(x_j) &= 0, \quad j \neq i, \\ h'_{i1}(x_i) &= 1 + 0 = 1. \end{aligned}$$

It follows that

$$\begin{aligned} (H_n f)(x_k) &= \sum_{i=0}^m \left[ h_{i0}(x_k) f(x_i) + h_{i1}(x_k) f'(x_i) \right] = f(x_k), \\ (H_n f)'(x_k) &= \sum_{i=0}^m \left[ h'_{i0}(x_k) f(x_i) + h'_{i1}(x_k) f'(x_i) \right] = f'(x_k), \quad k = \overline{0, m}, \end{aligned}$$

hence, the polynomial  $H_n f$  given in (1.3) satisfies the interpolation conditions (1.1).

To prove uniqueness, assume there exists another polynomial  $G_n$  (of degree at most  $n = 2m + 1$ ) satisfying relations (1.1) and consider

$$Q_n = H_n - G_n.$$

Then  $Q_n$  is also a polynomial of degree at most  $n = 2m + 1$ . From the interpolation conditions, it follows that

$$\begin{aligned} Q_n(x_i) &= H_n(x_i) - G_n(x_i) = f(x_i) - f(x_i) = 0, \quad i = 0, \dots, m, \\ Q'_n(x_i) &= H'_n(x_i) - G'_n(x_i) = f'(x_i) - f'(x_i) = 0, \quad i = 0, \dots, m. \end{aligned}$$

So,  $Q_n$ , a polynomial of degree at most  $2m + 1$ , has  $m + 1$  *double* roots. By the Fundamental Theorem of Algebra,  $Q_n$  must be identically zero, thus proving the uniqueness of  $H_n$ . □

**Example 1.2.** One of the most widely used form of Hermite interpolation is the cubic Hermite polynomial, which solves the interpolation problem with two double nodes  $a < b$ ,

$$\begin{aligned} P(a) &= f(a), \quad P(b) = f(b), \\ P'(a) &= f'(a), \quad P'(b) = f'(b). \end{aligned} \tag{1.5}$$

**Solution.** First of all, let us compute the degree. The degree of the polynomial is  $[2 \cdot (\text{number of nodes}) - 1]$ , so, in this case,

$$n = 2 \cdot 2 - 1 = 3.$$

Letting  $x_0 = a$ ,  $x_1 = b$ , with our previous notations and formulas, we have

$$\begin{aligned}\psi_1(x) &= (x - a)(x - b), \\ l_0(x) &= \frac{x - b}{a - b}, \quad l'_0(x) = \frac{1}{a - b}, \\ l_1(x) &= \frac{x - a}{b - a}, \quad l'_1(x) = \frac{1}{b - a}.\end{aligned}$$

The Hermite fundamental polynomials are given by

$$\begin{aligned}h_{00}(x) &= (1 - 2l'_0(a)(x - a))(l_0(x))^2 = \left[1 + 2\frac{x - a}{b - a}\right] \left[\frac{b - x}{b - a}\right]^2, \\ h_{10}(x) &= (1 - 2l'_1(b)(x - b))(l_1(x))^2 = \left[1 + 2\frac{b - x}{b - a}\right] \left[\frac{x - a}{b - a}\right]^2, \\ h_{01}(x) &= (x - a)(l_0(x))^2 = \frac{(x - a)(b - x)^2}{(b - a)^2}, \\ h_{11}(x) &= (x - b)(l_1(x))^2 = -\frac{(x - a)^2(b - x)}{(b - a)^2}.\end{aligned}$$

So the cubic Hermite polynomial is

$$\begin{aligned}H_3f(x) &= \left[1 + 2\frac{x - a}{b - a}\right] \left[\frac{b - x}{b - a}\right]^2 \cdot f(a) + \left[1 + 2\frac{b - x}{b - a}\right] \left[\frac{x - a}{b - a}\right]^2 \cdot f(b) \\ &+ \frac{(x - a)(b - x)^2}{(b - a)^2} \cdot f'(a) - \frac{(x - a)^2(b - x)}{(b - a)^2} \cdot f'(b).\end{aligned}$$

■

### 1.3.2 Newton's divided differences form

Just as in the case of Lagrange interpolation, Newton's divided differences provide a more easily computable form of the Hermite interpolation polynomial.

Consider  $2m + 2$  distinct nodes  $z_0, z_1, \dots, z_{2m}, z_{2m+1}$  and the Newton polynomial interpolating a function  $f$  at these nodes.

$$N_{2m+1}(x) = f(z_0) + f[z_0, z_1](x - z_0) + \dots + f[z_0, \dots, z_{2m+1}](x - z_0) \dots (x - z_{2m}),$$

with the error given by

$$R_{2m+1}(x) = f(x) - N_{2m+1}(x) = f[x, z_0, \dots, z_{2m+1}](x - z_0) \dots (x - z_{2m+1}).$$

We take the limits in the two relations above

$$z_0, z_1 \rightarrow x_0, \quad z_2, z_3 \rightarrow x_1, \quad \dots, \quad z_{2i}, z_{2i+1} \rightarrow x_i, \quad \dots \quad z_{2m}, z_{2m+1} \rightarrow x_m.$$

Denoting by  $n = 2m + 1$ , we get

$$\begin{aligned} N_n(x) &= f(x_0) + f[x_0, x_0](x - x_0) + f[x_0, x_0, x_1](x - x_0)^2 \\ &+ f[x_0, x_0, x_1, x_1](x - x_0)^2(x - x_1) + \dots \\ &+ f[x_0, x_0, \dots, x_m, x_m](x - x_0)^2 \dots (x - x_{m-1})^2(x - x_m) \end{aligned} \quad (1.6)$$

and for the remainder,

$$f(x) - N_n(x) = f[x, x_0, x_0, \dots, x_m, x_m](x - x_0)^2 \dots (x - x_m)^2. \quad (1.7)$$

**Proposition 1.3.** *Let  $[a, b] \subset \mathbb{R}$  be the smallest interval containing the distinct nodes  $x_0, \dots, x_m$  and  $f : [a, b] \rightarrow \mathbb{R}$  be a function of class  $C^{2m+2}[a, b]$ . Then, for the two polynomials in (1.3) and (1.6), we have*

$$H_n f(x) = N_n(x), \forall x \in [a, b], \quad (1.8)$$

with the interpolation error

$$R_n(x) = f(x) - H_n f(x) = [\psi_m(x)]^2 \frac{f^{(n+1)}(\xi_x)}{(n+1)!}, \quad \xi_x \in (a, b). \quad (1.9)$$

*Proof.* By the way it was constructed (in (1.6)), obviously the polynomial  $N_n$  has degree at most  $n$ . Then, by the uniqueness of the Hermite interpolation polynomial, it suffices to show that  $N_n$  satisfies the interpolation conditions (1.1).

From (1.7), it follows that

$$f(x_i) - N_n(x_i) = 0, \quad i = 0, \dots, m.$$



Also, by the same relation, we have for the derivatives,

$$\begin{aligned} f'(x) - N'_n(x) &= (x - x_0)^2 \dots (x - x_m)^2 \frac{\partial}{\partial x} f[x, x_0, x_0, \dots, x_m, x_m] \\ &+ 2f[x, x_0, x_0, \dots, x_m, x_m] \sum_{i=0}^m \left[ (x - x_i) \prod_{\substack{j=0 \\ j \neq i}}^m (x - x_j)^2 \right], \end{aligned}$$

hence,

$$f'(x_i) - N'_n(x_i) = 0, \quad i = 0, \dots, m.$$

Thus,

$$H_n f(x) = N_n(x), \quad \forall x \in [a, b]$$

and the error formula (1.9) follows directly from (1.7) and the mean-value formula for divided differences. □

**Example 1.4.** Let us find the polynomial and the remainder for the Hermite interpolation problem with two double nodes  $a < b$ , from Example 1.2.

**Solution.** We have

$$\begin{aligned} H_3 f(x) &= f(a) + f[a, a](x - a) + f[a, a, b](x - a)^2 \\ &+ f[a, a, b, b](x - a)^2(x - b). \end{aligned}$$

The divided differences table for two double nodes is

$z_0 = a$	$f(a)$	$\longrightarrow$	$f[a, a] = f'(a)$	$\longrightarrow$	$f[a, a, b]$	$\longrightarrow$	$f[a, a, b, b]$
		$\nearrow$		$\nearrow$		$\nearrow$	
$z_1 = a$	$f(a)$	$\longrightarrow$	$f[a, b] = \frac{f(b) - f(a)}{b - a}$	$\longrightarrow$	$f[a, b, b]$		
		$\nearrow$		$\nearrow$			
$z_2 = b$	$f(b)$	$\longrightarrow$	$f[b, b] = f'(b)$				
		$\nearrow$					
$z_3 = b$	$f(b),$						

where

$$\begin{aligned} f[a, a, b] &= \frac{f[a, b] - f'(a)}{b - a}, \\ f[a, b, b] &= \frac{f'(b) - f[a, b]}{b - a}, \\ f[a, a, b, b] &= \frac{f[a, b, b] - f[a, a, b]}{b - a} = \frac{f'(b) - 2f[a, b] + f'(a)}{(b - a)^2}. \end{aligned}$$

The interpolation error is given by

$$\begin{aligned} f(x) - H_3f(x) &= (x - a)^2(x - b)^2f[x, a, a, b, b] \\ &= \frac{(x - a)^2(x - b)^2}{24}f^{(4)}(\xi_x), \end{aligned}$$

with  $\xi_x$  belonging to the smallest interval that contains the points  $a, b$  and  $x$ .

We can find a bound for the error. Considering that on  $[a, b]$ , the maximum of the function  $|(x - a)(x - b)|$  occurs at the midpoint of the interval,  $\frac{a + b}{2}$ , and that the maximum value is  $\frac{(b - a)^2}{4}$ , we have

$$\max_{x \in [a, b]} |f(x) - H_3f(x)| \leq \frac{(b - a)^4}{384} \max_{t \in [a, b]} |f^{(4)}(t)|.$$

■

**Example 1.5** (Continuation of Example 1.1 in Lecture 4). Consider the function  $f : [0.5, 5] \rightarrow \mathbb{R}$ ,  $f(x) = \sqrt{x}$  and the nodes  $a = 1, b = 4$ . Let us compare Lagrange and Hermite approximations.

**Solution.** For the *simple* nodes  $a = 1, b = 4$ , we have the interpolation conditions

$$\begin{aligned} L_1f(a) &= f(a) = 1, \\ L_1f(b) &= f(b) = 2, \end{aligned}$$

satisfied by the Lagrange polynomial of degree 1

$$L_1f(x) = \frac{1}{3}x + \frac{2}{3}.$$

If the nodes are *double*, the interpolation conditions are

$$\begin{aligned} H_3 f(a) &= f(a) = 1, \\ H_3 f(b) &= f(b) = 2, \\ (H_3 f)'(a) &= f'(a) = 1/(2\sqrt{1}) = 1/2, \\ (H_3 f)'(b) &= f'(b) = 1/(2\sqrt{4}) = 1/4. \end{aligned}$$

The divided differences table is

$$\begin{array}{l|l} z_0 = 1 & \begin{array}{ccccccc} f(1) = 1 & \longrightarrow & f'(1) = 1/2 & \longrightarrow & f[1, 1, 4] = -1/18 & \longrightarrow & f[1, 1, 4, 4] = 1/108 \\ & \nearrow & & \nearrow & & \nearrow & \\ z_1 = 1 & f(1) = 1 & \longrightarrow & f[1, 4] = 1/3 & \longrightarrow & f[1, 4, 4] = -1/36 & \\ & \nearrow & & \nearrow & & & \\ z_2 = 4 & f(4) = 2 & \longrightarrow & f'(4) = 1/4 & & & \\ & \nearrow & & & & & \\ z_3 = 4 & f(4) = 2, & & & & & \end{array} \end{array}$$

The corresponding cubic Hermite interpolation polynomial is given by

$$H_3 f(x) = 1 + \frac{1}{2}(x-1) - \frac{1}{18}(x-1)^2 + \frac{1}{108}(x-1)^2(x-4),$$

with derivative

$$(H_3 f)'(x) = \frac{1}{2} - \frac{1}{9}(x-1) + \frac{1}{108}(x-1)[2(x-4) + (x-1)].$$

Check that  $H_3 f$  found above satisfies the interpolation conditions:

$$\begin{aligned} H_3 f(1) &= 1 = f(1), \\ H_3 f(4) &= 1 + \frac{3}{2} - \frac{1}{18} \cdot 9 = 2 = f(4), \\ (H_3 f)'(1) &= \frac{1}{2} = f'(1), \\ (H_3 f)'(4) &= \frac{1}{2} - \frac{1}{3} + \frac{1}{108} \cdot 9 = \frac{1}{4} = f'(4). \end{aligned}$$

The graphs of  $f$  and the two interpolation polynomials,  $L_1$ ,  $H_3$ , on the interval  $[0.5, 5]$ , are shown in Figure 1. The interpolation errors are plotted in Figure 2.

■

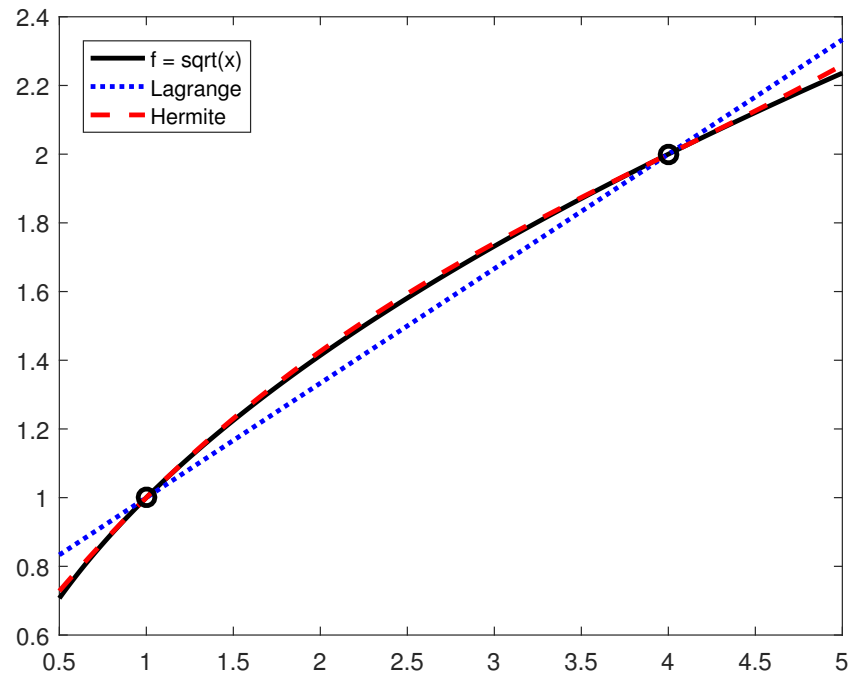


Fig. 1: Lagrange and Hermite interpolation with 2 nodes of the function  $\sqrt{x}$

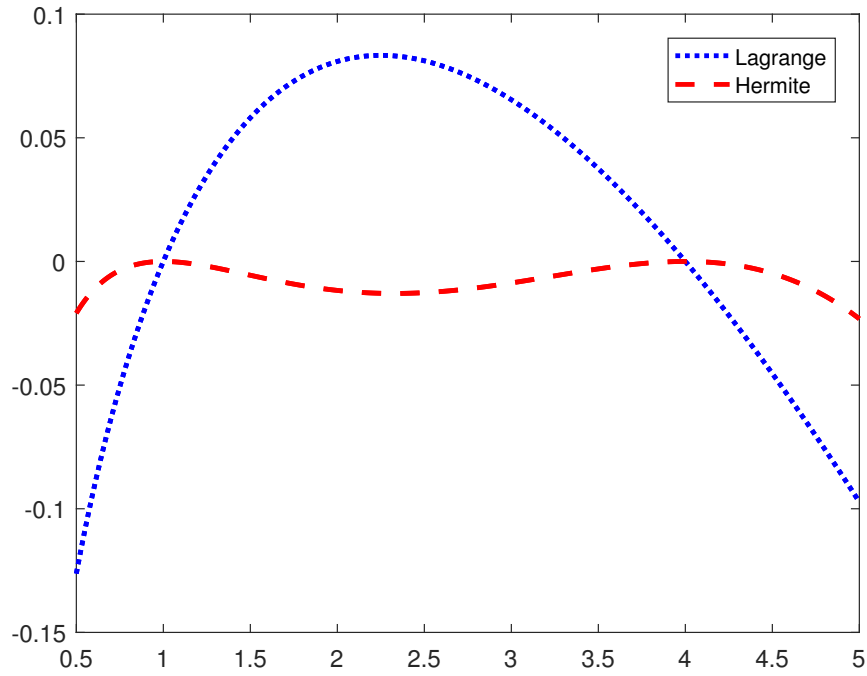


Fig. 2: Error of Lagrange and Hermite interpolation with 2 nodes of the function  $\sqrt{x}$

### 1.3.3 General case

**Hermite interpolation problem.** Given  $m + 1$  distinct nodes  $x_i \in [a, b], i = \overline{0, m}$ ,

$$\begin{aligned} x_0, & \text{ of multiplicity } r_0 + 1, \\ x_1, & \text{ of multiplicity } r_1 + 1, \\ & \dots \\ x_i, & \text{ of multiplicity } r_i + 1, \\ & \dots \\ x_m & \text{ of multiplicity } r_m + 1, \end{aligned}$$

and the values  $f^{(j)}(x_i), i = 0, 1, \dots, m, j = 0, \dots, r_i$ , of an unknown function  $f : [a, b] \rightarrow \mathbb{R}$  whose derivatives of order up to  $r_i$  exist at  $x_i, i = \overline{0, m}$ , find a polynomial  $P(x)$  of minimum degree, satisfying the interpolation conditions

$$P^{(j)}(x_i) = f^{(j)}(x_i), i = \overline{0, m}, j = \overline{0, r_i}. \quad (1.10)$$

Above, there are

$$n + 1 \stackrel{\text{not}}{=} \sum_{i=0}^m (r_i + 1)$$

conditions, so the polynomial satisfying these relations will have degree at most  $n$ .

**Theorem 1.6.** *There is a unique polynomial  $H_n f$  of degree at most  $n$ , satisfying the interpolation conditions (1.10). This polynomial is called the **Hermite interpolation polynomial** of the function  $f$ , relative to the nodes  $x_0, x_1, \dots, x_m$  and the integers  $r_0, r_1, \dots, r_m$ , and it can be written as*

$$H_n f(x) = \sum_{i=0}^m \sum_{j=0}^{r_i} h_{ij}(x) f^{(j)}(x_i). \quad (1.11)$$

**Remark 1.7.**

1. The functions  $h_{ij}(x), i = \overline{0, m}, j = \overline{0, r_i}$ , are called **Hermite fundamental (basis) polynomials** and they satisfy the relations

$$\begin{aligned} h_{ij}^{(k)}(x_l) &= 0, \quad l \neq i, k = \overline{0, r_l}, \\ h_{ij}^{(k)}(x_i) &= \delta_{jk}, \quad k = \overline{0, r_i}. \end{aligned} \quad (1.12)$$

2. If we denote by

$$\begin{aligned} u(x) &= \prod_{i=0}^m (x - x_i)^{r_i+1}, \\ u_i(x) &= \prod_{\substack{j=0 \\ j \neq i}}^m (x - x_j)^{r_j+1} = \frac{u(x)}{(x - x_i)^{r_i+1}}, \end{aligned} \quad (1.13)$$

then the fundamental polynomials  $h_{ij}(x)$  in (1.11) can be written as

$$h_{ij}(x) = \frac{(x - x_i)^j}{j!} \left[ \sum_{k=0}^{r_i-j} \frac{(x - x_i)^k}{k!} \left[ \frac{1}{u_i(x)} \right]_{x=x_i}^{(k)} \right] u_i(x). \quad (1.14)$$

2. A more computable form can be found using Newton divided differences. Re-indexing the nodes according to their multiplicity,

$$\begin{aligned} z_0 &= x_0, \dots, z_{r_0} = x_0, \\ z_{r_0+1} &= x_1, \dots, z_{(r_0+1)+r_1} = x_1, \\ z_{(r_0+1)+(r_1+1)} &= x_2, \dots, z_{(r_0+1)+(r_1+1)+r_2} = x_2, \\ &\dots \\ z_{n-r_m} &= x_m, \dots, z_n = x_m, \end{aligned}$$

the Hermite polynomial can be written in Newton's form as

$$N_n f(x) = f(z_0) + f[z_0, z_1](x - z_0) + \dots + f[z_0, \dots, z_n](x - z_0) \dots (x - z_{n-1}), \quad (1.15)$$

with interpolation error

$$\begin{aligned} R_n(x) &= f(x) - N_n(x) = f[x, z_0, \dots, z_n](x - z_0) \dots (x - z_n) \\ &= \frac{u(x)}{(n+1)!} f^{(n+1)}(\xi_x), \quad \xi_x \in (a, b). \end{aligned} \quad (1.16)$$

**Example 1.8.** Consider the case of a simple node  $x_0$  and a double node  $x_1$ . Find the interpolant for this data and an expression for the remainder.

**Solution.** We have the nodes

$$\begin{aligned} x_0, & \text{ of multiplicity } r_0 + 1 = 1, \\ x_1, & \text{ of multiplicity } r_1 + 1 = 2. \end{aligned}$$

so  $n + 1 = 1 + 2$  and the polynomial has degree  $n = 2$ .

The divided differences table:

$$\begin{array}{c|ccc}
 x_0 & f(x_0) & \longrightarrow & f[x_0, x_1] & \longrightarrow & \frac{f'(x_1) - f[x_0, x_1]}{x_1 - x_0} \\
 & & \nearrow & & \nearrow & \\
 x_1 & f(x_1) & \longrightarrow & f'(x_1) & & \\
 & & \nearrow & & & \\
 x_1 & f(x_1) & & & & 
 \end{array}$$

Then,

$$\begin{aligned}
 H_2 f(x) &= f(x_0) + f[x_0, x_1](x - x_0) + \frac{f'(x_1) - f[x_0, x_1]}{x_1 - x_0}(x - x_0)(x - x_1) \\
 &= f(x_0) + \frac{f(x_1) - f(x_0)}{x_1 - x_0}(x - x_0) + \frac{f'(x_1)}{x_1 - x_0}(x - x_0)(x - x_1) \\
 &\quad - \frac{f(x_1) - f(x_0)}{(x_1 - x_0)^2}(x - x_0)(x - x_1) \\
 &= h_{00}f(x_0) + h_{10}f(x_1) + h_{11}f'(x_1)
 \end{aligned}$$

and the remainder is given by

$$R_2 f(x) = \frac{(x - x_0)(x - x_1)^2}{3!} f'''(\xi),$$

with  $\xi$  belonging to the smallest interval containing  $x_0$  and  $x_1$ .

Now, since  $H_2 f$  has degree 2 (small), we can find it directly: we seek it of the form

$$H_2 f(x) = ax^2 + bx + c$$

and determine coefficients  $a, b$  and  $c$  from the interpolation conditions:

$$\begin{cases}
 H_2 f(x_0) = f(x_0) \\
 H_2 f(x_1) = f(x_1) \\
 (H_2 f)'(x_1) = f'(x_1)
 \end{cases} ,$$

i.e., from the linear system

$$\begin{cases} x_0^2 a + x_0 b + c = f(x_0) \\ x_1^2 a + x_1 b + c = f(x_1) \\ 2x_1 a + b = f'(x_1) \end{cases} . \quad (1.17)$$

The matrix of this system,

$$V = \begin{bmatrix} 1 & x_0 & x_0^2 \\ 1 & x_1 & x_1^2 \\ 0 & 1 & 2x_1 \end{bmatrix},$$

is called a *generalized Vandermonde matrix*. It is invertible and the elements of its inverse are the coefficients of the fundamental polynomials  $h_{00}$ ,  $h_{10}$  and  $h_{11}$ .

If the node  $x_0$  is double and  $x_1$  is simple, the corresponding Hermite polynomial and its error are given by

$$\begin{aligned} H_2 f(x) &= f(x_0) + f'(x_0)(x - x_0) + \frac{f[x_0, x_1] - f'(x_0)}{x_1 - x_0}(x - x_0)^2, \\ R_2 f(x) &= \frac{(x - x_0)^2(x - x_1)}{3!} f'''(\xi). \end{aligned}$$

■

**Example 1.9.** Find a polynomial of minimum degree that interpolates the data  $f(0)$ ,  $f(1)$ ,  $f'(1)$  and  $f''(1)$  (so, a *simple* node and a *triple* one). Evaluate the error.

**Solution.** We have the nodes

$$\begin{aligned} x_0 &= 0, \text{ of multiplicity } r_0 + 1 = 1, \\ x_1 &= 1, \text{ of multiplicity } r_1 + 1 = 3. \end{aligned}$$

Hence, we seek the Hermite polynomial of degree at most

$$n = 1 + 3 - 1 = 3.$$

This will be of the form

$$H_3 f(x) = h_{00}(x)f(0) + h_{10}(x)f(1) + h_{11}(x)f'(1) + h_{12}(x)f''(1).$$



We compute the divided differences

$$\begin{array}{l|l}
 0 & f(0) \longrightarrow f(1) - f(0) \longrightarrow f'(1) - f(1) + f(0) \longrightarrow \frac{f''(1)}{2} - f'(1) + f(1) - f(0) \\
 & \nearrow \qquad \qquad \qquad \nearrow \qquad \qquad \qquad \nearrow \\
 1 & f(1) \longrightarrow f'(1) \longrightarrow \frac{f''(1)}{2} \\
 & \nearrow \qquad \qquad \qquad \nearrow \\
 1 & f(1) \longrightarrow f'(1) \\
 & \nearrow \\
 1 & f(1)
 \end{array}$$

Then the interpolant is

$$\begin{aligned}
 H_3 f(x) &= f(0) + (f(1) - f(0))x + (f'(1) - f(1) + f(0))x(x-1) \\
 &\quad + \left( \frac{f''(1)}{2} - f'(1) + f(1) - f(0) \right) x(x-1)^2 \\
 &= -(x-1)^3 f(0) + x(x^2 - 3x + 3)f(1) - x(x-1)(x-2)f'(1) + \frac{1}{2}x(x-1)^2 f''(1).
 \end{aligned}$$

So the fundamental polynomials are

$$\begin{aligned}
 h_{00}(x) &= -(x-1)^3, \\
 h_{10}(x) &= x(x^2 - 3x + 3), \\
 h_{11}(x) &= -x(x-1)(x-2), \\
 h_{12}(x) &= \frac{1}{2}x(x-1)^2,
 \end{aligned}$$

with derivatives

$$\begin{aligned}
 h'_{00}(x) &= -3(x-1)^2, & h''_{00}(x) &= -6(x-1), \\
 h'_{10}(x) &= 3(x-1)^2, & h''_{10}(x) &= 6(x-1), \\
 h'_{11}(x) &= -(3x^2 - 6x + 2), & h''_{11}(x) &= -6(x-1), \\
 h'_{12}(x) &= \frac{1}{2}(x-1)(3x-2), & h''_{12}(x) &= 3x-2.
 \end{aligned}$$

Now, we can better understand relations (1.12), as we can easily see that

$$\begin{pmatrix} h_{00}(0) \\ h_{00}(1) \\ h'_{00}(1) \\ h''_{00}(1) \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \quad \begin{pmatrix} h_{10}(0) \\ h_{10}(1) \\ h'_{10}(1) \\ h''_{10}(1) \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}, \quad \begin{pmatrix} h_{11}(0) \\ h_{11}(1) \\ h'_{11}(1) \\ h''_{11}(1) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}, \quad \begin{pmatrix} h_{12}(0) \\ h_{12}(1) \\ h'_{12}(1) \\ h''_{12}(1) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}.$$

Also, it is now very easy to check that  $H_3f$  satisfies the interpolation conditions.

Alternatively, we can write the polynomial in the form

$$\begin{aligned} H_3f(x) &= \left(-f(0) + f(1) - f'(1) + \frac{1}{2}f''(1)\right)x^3 + \left(3f(0) - 3f(1) + 3f'(1) - f''(1)\right)x^2 \\ &\quad + \left(-3f(0) + 3f(1) - 2f'(1) + \frac{1}{2}f''(1)\right)x + f(0). \end{aligned}$$

For the remainder, we have

$$R_3f(x) = \frac{u(x)}{4!}f^{(iv)}(\xi) = \frac{x(x-1)^3}{4!}f^{(iv)}(\xi), \quad \xi \in (0, 1).$$

Now,

$$\begin{aligned} u(x) &= x(x-1)^3 = x^4 - 3x^3 + 3x^2 - x, \\ u'(x) &= 4x^3 - 9x^2 + 6x - 1 = (x-1)^2(4x-1), \end{aligned}$$

so  $u(x) \leq 0$  on  $[0, 1]$  and it has a local minimum at  $x = \frac{1}{4}$ . Thus,

$$|u(x)| \leq |u(1/4)| = \left| \frac{1}{4} \left(-\frac{3}{4}\right)^3 \right| = \frac{27}{256}.$$

Then, we find an error bound as

$$|R_3f(x)| \leq \frac{27}{256 \cdot 4!} \max_{t \in [0,1]} |f^{(iv)}(t)| \approx 0.0044 \cdot \|f^{(iv)}\|.$$

■

### Special cases

1. If all  $r_i = 0, i = \overline{0, m}$ , all the nodes are simple and we have the Lagrange interpolation formula.
2. If we consider one single node,  $x_0$ , of multiplicity  $n + 1$ , the Hermite interpolation polynomial is reduced to Taylor's polynomial:

$$\begin{aligned} H_n f(x) &= T_n f(x) = f(x_0) + \frac{x - x_0}{1!} f'(x_0) + \frac{(x - x_0)^2}{2!} f''(x_0) + \dots \\ &+ \frac{(x - x_0)^n}{n!} f^{(n)}(x_0), \end{aligned} \quad (1.18)$$

with remainder

$$R_n(f)(x) = \frac{(x - x_0)^{n+1}}{(n + 1)!} f^{(n+1)}(\xi_x). \quad (1.19)$$

3. Consider two nodes,  $x_0 = a$ , of multiplicity  $m + 1$  and  $x_1 = b$ , of multiplicity  $n + 1$ .

The Hermite polynomial has degree

$$(m + 1) + (n + 1) - 1 = m + n + 1.$$

With the notations from Remark 1.7, we have

$$\begin{aligned} u(x) &= (x - a)^{m+1}(x - b)^{n+1}, \\ u_0(x) &= (x - b)^{n+1}, \\ u_1(x) &= (x - a)^{m+1}. \end{aligned}$$

The Hermite polynomial is of the form

$$H_{m+n+1} f(x) = \sum_{j=0}^m h_{0j}(x) f^{(j)}(a) + \sum_{i=0}^n h_{1i}(x) f^{(i)}(b) \quad (1.20)$$

and the fundamental polynomials are given by

$$\begin{aligned} h_{0j}(x) &= \frac{(x - a)^j}{j!} \left[ \sum_{k=0}^{m-j} \frac{(x - a)^k}{k!} \left[ \frac{1}{(x - b)^{n+1}} \right]_{x=a}^{(k)} \right] (x - b)^{n+1}, \\ h_{1i}(x) &= \frac{(x - b)^i}{i!} \left[ \sum_{k=0}^{n-i} \frac{(x - b)^k}{k!} \left[ \frac{1}{(x - a)^{m+1}} \right]_{x=b}^{(k)} \right] (x - a)^{m+1}. \end{aligned}$$

In Newton's form (1.15),

$$\begin{aligned}
H_{m+n+1}f(x) &= f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \cdots + \frac{f^{(m)}(a)}{m!}(x-a)^m \\
&+ f[\underbrace{a, \dots, a}_{m+1}, b](x-a)^{m+1} + f[\underbrace{a, \dots, a, b}_{m+1}](x-a)^{m+1}(x-b) \\
&+ \cdots + f[\underbrace{a, \dots, a}_{m+1}, \underbrace{b, \dots, b}_{n+1}](x-a)^{m+1}(x-b)^n,
\end{aligned}$$

with remainder

$$\begin{aligned}
R_{m+n+1} &= f[x, \underbrace{a, \dots, a}_{m+1}, \underbrace{b, \dots, b}_{n+1}](x-a)^{m+1}(x-b)^{n+1} \\
&= \frac{f^{(m+n+2)}(\xi_x)}{(m+n+2)!}(x-a)^{m+1}(x-b)^{n+1}, \quad \xi_x \in (a, b).
\end{aligned}$$

## 1.4 Birkhoff Interpolation

Consider the following situation: We have a moving object and the times  $t_0, t_1, \dots, t_m$ . For some of these nodes, we know the *distances* traveled  $d_i = d(t_i), i \in I \subset \{0, 1, \dots, m\}$ , for others, the *velocities*  $v_j = d'(t_j), j \in \tilde{I} \subset \{0, 1, \dots, m\}$  and for others, the *accelerations*  $a_k = d''(t_k), k \in I^* \subset \{0, 1, \dots, m\}$ . Having all these data, can we find a polynomial approximation of the distance function  $d = d(t)$  on the entire interval containing the points  $t_0, \dots, t_m$ ?

Obviously, this is *not* a Lagrange interpolation problem, because we do not have the values of the function at all the nodes. We *cannot* find a Hermite polynomial, either, because at some nodes, only the value of the derivative (or the second derivative) is given (without the values of the function). This is a **Birkhoff interpolation** problem, also known as *lacunary Hermite interpolation* (because not *all* the functional or derivative values for all points are provided) and it is more general than Hermite interpolation.

### 1.4.1 Birkhoff interpolation polynomial

**Birkhoff interpolation problem.** Let  $x_k \in [a, b], k = \overline{0, m}$ , be  $m + 1$  distinct nodes,  $r_k \in \mathbb{N}$  and  $I_k \subseteq \{0, \dots, r_k\}, k = 0, \dots, m$ . Consider the function  $f : [a, b] \rightarrow \mathbb{R}$  whose derivatives  $f^{(j)}(x_k), k = 0, \dots, m, j \in I_k$  exist. Find a polynomial  $P(x)$  of minimum degree, satisfying the interpolation conditions

$$P^{(j)}(x_k) = f^{(j)}(x_k), k = \overline{0, m}, j \in I_k. \quad (1.1)$$

Denote by  $n + 1 = |I_0| + \dots + |I_m|$ , where  $|I_k|$  is the cardinality (number of elements) of  $I_k$ . There are  $n + 1$  interpolation conditions in (1.1), so we seek a polynomial of degree at most  $n$ ,

$$P_n(x) = a_0 + a_1x + \dots + a_nx^n,$$

whose coefficients are found from the linear system generated by the interpolation conditions (1.1). If the determinant of this system is not equal to zero, then the Birkhoff interpolation problem has a unique solution.

**Remark 1.1.** If  $I_k = \{0, 1, \dots, r_k\}$ , for every  $k = 0, \dots, m$ , then the Birkhoff interpolation problem is reduced to Hermite interpolation (which, in turn, is reduced to Lagrange interpolation when  $r_k = 0, k = 0, \dots, m$ ). Hence, Birkhoff interpolation is more general.

Unlike Lagrange and Hermite interpolation, the Birkhoff interpolation problem (1.1) *does not*

always have a solution. When such a polynomial, denoted by  $B_n f$ , exists, it has the form

$$B_n f(x) = \sum_{k=0}^m \sum_{j \in I_k} b_{kj}(x) f^{(j)}(x_k). \quad (1.2)$$

The terms  $b_{kj}(x)$  are called **Birkhoff fundamental polynomials** and they satisfy the relations:

$$\begin{aligned} b_{kj}^{(p)}(x_\nu) &= 0, \quad \nu \neq k, \quad p \in I_\nu, \\ b_{kj}^{(p)}(x_k) &= \delta_{jp}, \quad p \in I_k, \quad \text{for } j \in I_k \text{ and } \nu, k = 0, 1, \dots, m, \end{aligned} \quad (1.3)$$

where

$$\delta_{jp} = \begin{cases} 0, & j \neq p \\ 1, & j = p \end{cases}$$

is Kronecker's symbol.

**Remark 1.2.** Because some of the functional (or derivative) values are missing, finding mathematical expressions for the Birkhoff fundamental polynomials  $b_{kj}$ ,  $k = 0, \dots, m; j \in I_k$ , is, in general, difficult. They can be determined (when possible) directly from the conditions (1.3).

**Example 1.3.** Let  $f \in C^2[0, 1]$  and consider the nodes  $x_0 = 0$ ,  $x_1 = 1$ , for which the values  $f(0) = 1$  and  $f'(1) = 2$  are given. Find the Birkhoff polynomial that interpolates these data.

**Solution.**

We have  $m = 1$ , two nodes, with  $I_0 = \{0\}$ ,  $I_1 = \{1\}$ , so  $n = 1 + 1 - 1 = 1$ . We want a polynomial of degree 1,

$$P(x) = a_0 + a_1 x,$$

satisfying the conditions

$$\begin{aligned} P(0) &= f(0), \\ P'(1) &= f'(1). \end{aligned}$$

From here, we have the linear system

$$\begin{cases} a_0 &= f(0) \\ a_1 &= f'(1) \end{cases}.$$

The determinant of this system is

$$\begin{vmatrix} 1 & 0 \\ 0 & 1 \end{vmatrix} = 1 \neq 0,$$

so this problem has a unique solution

$$\begin{cases} a_0 &= f(0) \\ a_1 &= f'(1) \end{cases},$$

i.e. the polynomial we seek is

$$P(x) = f(0) + f'(1)x = 1 + 2x.$$

On the other hand, by (1.2), the Birkhoff polynomial is of the form

$$B_1 f(x) = b_{00}(x)f(0) + b_{11}(x)f'(1).$$

Let us find the fundamental polynomials  $b_{00}(x)$  și  $b_{11}(x)$ . Both have degree 1, hence,

$$\begin{aligned} b_{00}(x) &= ax + b, \\ b_{11}(x) &= cx + d. \end{aligned}$$

By conditions (1.3), for  $b_{00}$ , we have

$$\begin{cases} b_{00}(x_0) &= 1 \\ b'_{00}(x_1) &= 0 \end{cases} \iff \begin{cases} b_{00}(0) &= 1 \\ b'_{00}(1) &= 0 \end{cases} \iff \begin{cases} b &= 1 \\ a &= 0 \end{cases},$$

thus,

$$b_{00}(x) = 1.$$

Similarly, for  $b_{11}$ , we have

$$\begin{cases} b_{11}(x_0) &= 0 \\ b'_{11}(x_1) &= 1 \end{cases} \iff \begin{cases} b_{11}(0) &= 0 \\ b'_{11}(1) &= 1 \end{cases} \iff \begin{cases} d &= 0 \\ c &= 1 \end{cases},$$

so we get

$$b_{11}(x) = x.$$

Thus,

$$B_1 f(x) = f(0) + x f'(1) = 1 + 2x.$$

■

**Example 1.4.** Find a polynomial of smallest degree (if it exists) satisfying the conditions

$$P(-1) = P(1) = 0, P'(0) = 1.$$

**Solution.**

Here, we have 3 nodes,  $x_0 = -1, x_1 = 0, x_2 = 1$ ,  $m = 2$ , for which  $I_0 = \{0\}, I_1 = \{1\}, I_2 = \{0\}$ . Hence, we seek a polynomial of degree  $n = 1 + 1 + 1 - 1 = 2$ . This is of the form

$$P(x) = a_0 + a_1 x + a_2 x^2$$

and must satisfy the relations

$$P(-1) = 0,$$

$$P'(0) = 1,$$

$$P(1) = 0.$$

We obtain the linear system

$$\begin{cases} a_0 - a_1 + a_2 = 0 \\ a_1 = 1 \\ a_0 + a_1 + a_2 = 0 \end{cases}.$$

Subtracting the first equation from the third, we get  $a_1 = 0$ , which contradicts the second equation. The system is incompatible and, thus, this interpolation problem *does not have* a solution.

■

**Example 1.5. [Abel-Goncharov interpolation]** Let  $f \in C^{n+1}[0, nh]$ , with  $h > 0$ ,  $n \in \mathbb{N}$ . Find a



polynomial of smallest degree satisfying the relations

$$\begin{aligned} P(0) &= f(0), \\ P'(h) &= f'(h), \\ &\dots \\ P^{(n)}(nh) &= f^{(n)}(nh). \end{aligned}$$

This problem has a unique solution for every  $h > 0$ ,  $n \in \mathbb{N}$ . Notice that the problem in Example 1.3 was of this type, with  $n = 1$  and  $h = 1$ .

### 1.4.2 Peano's theorem and the error for Birkhoff interpolation

To find an error formula for Birkhoff interpolation (when the Birkhoff polynomial exists), we need an important result from linear operator theory.

Let us recall some notions and properties:

- Let  $n \in \mathbb{N}^*$ . We define the space

$$H^n[a, b] = \{f : [a, b] \rightarrow \mathbb{R} \mid f \in C^{n-1}[a, b], f^{(n-1)} \text{ absolutely continuous on } [a, b]\}.$$

- A function  $f : [a, b] \rightarrow \mathbb{R}$  is *absolutely continuous* on  $[a, b]$ , if, for instance, it has a derivative  $f'$  almost everywhere, the derivative is Lebesgue integrable, and

$$f(x) = f(a) + \int_a^x f'(t) dt, \quad \forall x \in [a, b].$$

- $H^n[a, b]$  is linear space.
- Any function  $f \in H^n[a, b]$  has a Taylor-type representation, with the remainder in integral form

$$f(x) = \sum_{k=0}^{n-1} \frac{(x-a)^k}{k!} f^{(k)}(a) + \int_a^x \frac{(x-t)^{n-1}}{(n-1)!} f^{(n)}(t) dt.$$

- The function

$$z_+ = \begin{cases} z, & z \geq 0 \\ 0, & z < 0 \end{cases}$$

is called the *positive part* of  $z$ , and  $z_+^n = (z_+)^n$  is called a *truncated power*.

- For  $m \in \mathbb{N}$ ,  $\mathbb{P}_m$  denotes the space of all polynomials of degree at most  $m$ . Obviously,  $\mathbb{P}_m \subset C^\infty[a, b], \forall m \geq 0$ .
- The *kernel* of a linear map  $L : V \rightarrow W$  between two vector spaces  $V$  and  $W$ , is the set of all vectors in  $V$  that are mapped to zero:

$$\ker L = \{v \in V \mid L(v) = \mathbf{0}_W\},$$

where  $\mathbf{0}_W$  is null vector in  $W$ .

The next theorem is **paramount** in Numerical Analysis. It gives a representation of real linear functionals defined on a space  $H^m[a, b]$ . This result provides means for expressing the errors in many approximating procedures.

**Theorem 1.6. [Peano]**

Let  $L : H^{n+1}[a, b] \rightarrow \mathbb{R}$  be a linear functional that commutes with the definite integral operator. If  $\ker L = \mathbb{P}_n$ , then

$$Lf(\mathbf{x}) = \int_a^b K_n(\mathbf{x}, t) f^{(n+1)}(t) dt, \quad (1.4)$$

where

$$K_n(\mathbf{x}, t) = \frac{1}{n!} L\left((\mathbf{x} - t)_+^n\right) \quad (1.5)$$

is called the **Peano kernel**.

So, this is saying that if  $L$  maps *all* polynomials of degree at most  $n$  to the *function identically equal to 0*, i.e.  $Le_k \equiv 0, e_k(x) = x^k, k = 0, 1, \dots, n, Le_{n+1} \not\equiv 0$ , then  $Lf$  can be expressed as in (1.4).

**Corollary 1.7.** If the kernel  $K$  has constant sign on  $[a, b]$  and  $f^{(n+1)}$  is continuous on  $[a, b]$ , then there exists  $\xi \in (a, b)$  such that

$$Lf = \frac{1}{(n+1)!} f^{(n+1)}(\xi) Le_{n+1}, \quad (1.6)$$

where  $e_k(x) = x^k, k \in \mathbb{N}$ .

*Proof.* If the kernel  $K$  has constant sign on  $[a, b]$ , we can apply the mean value theorem in (1.4):

$$Lf = f^{(n+1)}(\xi) \int_a^b K_n(x, t) dt, \quad \xi \in (a, b). \quad (1.7)$$

Notice that the kernel  $K$  *does not* depend on  $f$  and the relation above is true *regardless* of the function  $f$ . Then, taking  $f = e_{n+1}$ , we get

$$\begin{aligned} Le_{n+1} &= e_{n+1}^{(n+1)}(\xi) \int_a^b K_n(x, t) dt \\ &= (n+1)! \int_a^b K_n(x, t) dt, \end{aligned}$$

from which we get

$$\int_a^b K_n(x, t) dt = \frac{1}{(n+1)!} Le_{n+1}.$$

Using this in (1.7), we obtain (1.6). □

**Remark 1.8.** This corollary is the one that is mostly used in applications, to assess the approximation error. We apply Theorem 1.6 or Corollary 1.7 to the *remainder functional*. We derive an approximation formula

$$f(x) = B_n f(x) + R_n f(x).$$

Since  $B_n$  is a polynomial of degree  $n$ , the formula above is *exact* for all polynomials of degree  $n$ , i.e.

$$R_n e_k = (f - B_n) e_k = \mathbf{0}.$$

In this case, we say that the approximation formula has *degree of precision* (or *degree of exactness*)

$d = n$ . Then,

$$\begin{aligned} K_n(x, t) &= \frac{1}{n!} R_n((x - t)_+^n) \\ &= \frac{1}{n!} [(x - t)_+^n - B_n((x - t)_+^n)]. \end{aligned} \quad (1.8)$$

Now, it is easy to check (from the definition) that the function  $F(x) = (x - t)_+^n$  has the *derivative*

$$F'(x) = \frac{\partial [(x - t)_+^n]}{\partial x} = n(x - t)_+^{n-1}$$

and the *integral* (this will only be needed later on, in Chapter 4)

$$\int_a^b F(x) dx = \frac{1}{n+1} (x - t)_+^{n+1} \Big|_{x=a}^{x=b} = \frac{1}{n+1} [(b - t)_+^{n+1} - (a - t)_+^{n+1}]$$

If  $K_n(x, t)$  above has constant sign on  $[a, b]$  (the smallest interval containing the interpolation nodes), then we have an expression for the error of the approximation, as

$$R_n f = \frac{1}{(n+1)!} f^{(n+1)}(\xi) R_n e_{n+1}, \quad \xi \in (a, b). \quad (1.9)$$

**Example 1.9.** Let us find a formula for the rest of the Birkhoff polynomial in Example 1.3.

**Solution.** We found the Birkhoff polynomial

$$B_1 f(x) = f(0) + f'(1)x, \quad x \in [0, 1],$$

so, we have

$$f(x) = B_1 f(x) + R_1 f(x).$$

We apply Peano's theorem to the remainder operator,

$$\begin{aligned} Lf &= R_1 f = f - B_1 f, \\ Lf(x) &= (f - B_1 f)(x) = f(x) - (f(0) + f'(1)x). \end{aligned}$$

We have

$$\begin{aligned} R_1 e_0(x) &= e_0(x) - B_1 e_0(x) = e_0(x) - (e_0(0) + e'_0(1)x) = 1 - (1 + 0) \equiv 0, \\ R_1 e_1(x) &= e_1(x) - B_1 e_1(x) = e_1(x) - (e_1(0) + e'_1(1)x) = x - (0 + 1 \cdot x) \equiv 0, \\ R_1 e_2(x) &= e_2(x) - B_1 e_2(x) = e_2(x) - (e_2(0) + e'_2(1)x) = x^2 - 2x \not\equiv 0, \end{aligned}$$

(the first two were obvious, since  $B_1$  is a polynomial of degree 1, but it is a good computational exercise). Thus,

$$R_1 f(x) = \int_0^1 K_1(x, t) f''(t) dt,$$

with

$$\begin{aligned} K_1(x, t) &= \frac{1}{1!} R_1 ((x - t)_+^1) \\ &= \underbrace{(x - t)_+}_{f(x)} - \left( \underbrace{(0 - t)_+}_{f(0)} + \underbrace{1 \cdot x}_{f'(1)x} \right) \\ &= (x - t)_+ - (-t)_+ - x. \end{aligned}$$

Since  $x, t \in [0, 1]$ , we have  $(-t)_+ = 0$ . Fix an arbitrary  $x \in [0, 1]$ .

If  $0 \leq t \leq x$ , then  $(x - t)_+ = x - t$  and

$$K_1(x, t) = x - t - x = -t \leq 0.$$

If  $x \leq t \leq 1$ , then  $(x - t)_+ = 0$  and

$$K_1(x, t) = 0 - x = -x \leq 0.$$

So, either way,  $K_1$  has constant sign on  $[0, 1]$ . By Corollary 1.7, it follows that

$$\begin{aligned} R_1 f(x) &= \frac{1}{2!} f''(\xi) R_1 e_2(x) \\ &= \frac{x^2 - 2x}{2!} f''(\xi), \quad \xi \in (0, 1). \end{aligned}$$

Now,  $x(x - 2)$  is a parabola with the minimum value at the midpoint value of the interval of the

roots. Thus,

$$|x(x - 2)| \leq 1,$$

for  $x \in [0, 1]$  and we have the following estimate for the interpolation error:

$$|R_1 f(x)| \leq \frac{1}{2} \|f''\|_\infty, \forall x \in [0, 1].$$

■

**Example 1.10.** Let  $f : [0, 2] \rightarrow \mathbb{R}$  be a function in  $C^3[0, 2]$ .

- a) Find a polynomial of minimum degree that interpolates the data  $f'(0)$ ,  $f(1)$  and  $f'(2)$ ;
- b) If  $f'(0) = 1$ ,  $f(1) = 2$  and  $f'(2) = 1$ , use the approximation found in part a) to approximate  $f(1/2)$  and estimate the error.

**Solution.**

a) Since, for instance, for  $x_0 = 0$ , *only* the derivative is given, without the function value, this is Birkhoff interpolation. We have  $m + 1 = 3$  nodes, with  $I_0 = \{1\}$ ,  $I_1 = \{0\}$ ,  $I_2 = \{1\}$ , so  $n = 1 + 1 + 1 - 1 = 2$ . The Birkhoff polynomial of degree (at most) 2,

$$\begin{aligned} B_2(x) &= ax^2 + bx + c, \\ B_2'(x) &= 2ax + b, \end{aligned}$$

must satisfy the relations

$$\begin{aligned} B_2'(0) &= b = f'(0), \\ B_2(1) &= a + b + c = f(1) \\ B_2'(2) &= 4a + b = f'(2). \end{aligned}$$

The determinant of the corresponding linear system is

$$\begin{vmatrix} 0 & 1 & 0 \\ 1 & 1 & 1 \\ 4 & 1 & 0 \end{vmatrix} = - \begin{vmatrix} 1 & 1 \\ 4 & 0 \end{vmatrix} = 4 \neq 0,$$

hence, there exists a unique Birkhoff interpolation polynomial, of the form

$$B_2 f(x) = b_{01}(x) f'(0) + b_{10}(x) f(1) + b_{21}(x) f'(2).$$

The fundamental polynomials (of degree at most 2) must satisfy the conditions

$$\begin{cases} b'_{01}(0) = 1 \\ b_{01}(1) = 0 \\ b'_{01}(2) = 0 \end{cases}, \quad \begin{cases} b'_{10}(0) = 0 \\ b_{10}(1) = 1 \\ b'_{10}(2) = 0 \end{cases} \quad \text{and} \quad \begin{cases} b'_{21}(0) = 0 \\ b_{21}(1) = 0 \\ b'_{21}(2) = 1 \end{cases}.$$

Each of them is of the form  $ax^2 + bx + c$ .

For  $b_{01}$ :

$$\begin{cases} b = 1 \\ a + b + c = 0 \\ 4a + b = 0 \end{cases} \iff \begin{cases} a = -1/4 \\ b = 1 \\ c = -3/4 \end{cases},$$

so

$$b_{01}(x) = -\frac{x^2 - 4x + 3}{4} = -\frac{(x-1)(x-3)}{4}.$$

For  $b_{10}$ , we have:

$$\begin{cases} b = 0 \\ a + b + c = 1 \\ 4a + b = 0 \end{cases} \iff \begin{cases} a = 0 \\ b = 0 \\ c = 1 \end{cases},$$

hence,

$$b_{10}(x) = 1.$$

Finally, for  $b_{21}$ :

$$\begin{cases} b = 0 \\ a + b + c = 0 \\ 4a + b = 1 \end{cases} \iff \begin{cases} a = 1/4 \\ b = 0 \\ c = -1/4 \end{cases}$$

and, thus,

$$b_{21}(x) = \frac{x^2 - 1}{4}.$$

From these, we find the interpolation polynomial

$$B_2f(x) = \frac{1}{4}(x-1)(3-x)f'(0) + f(1) + \frac{1}{4}(x^2-1)f'(2),$$

with derivative

$$(B_2f)'(x) = \frac{1}{2}(2-x)f'(0) + \frac{1}{2}xf'(2).$$

Check that  $B_2f$  satisfies the interpolation conditions.

The remainder is computed as

$$R_2f(x) = f(x) - B_2f(x) = f(x) - \left[ \frac{1}{4}(x-1)(3-x)f'(0) + f(1) + \frac{1}{4}(x^2-1)f'(2) \right].$$

We have

$$\begin{aligned} R_2e_0(x) &= 1 - [0 + 1 + 0] \equiv 0, \\ R_2e_1(x) &= x - \left[ \frac{-x^2 + 4x - 3}{4} \cdot 1 + 1 + \frac{1}{4}(x^2 - 1) \cdot 1 \right] \equiv x - \left[ x - \frac{3}{4} + 1 - \frac{1}{4} \right] = 0, \\ R_2e_2(x) &= x^2 - \left[ 0 + 1 + \frac{1}{4}(x^2 - 1) \cdot 4 \right] = x^2 - 1 - (x^2 - 1) \equiv 0, \\ R_2e_3(x) &= x^3 - \left[ 0 + 1 + \frac{1}{4}(x^2 - 1) \cdot 12 \right] = (x-1)(x^2 - 2x - 2) \not\equiv 0. \end{aligned}$$

(again, the first three relations *needed not* be checked). So the approximation formula

$$f(x) \approx B_2f(x)$$

has degree of precision  $d = 2$ . Then the remainder is given by

$$R_2f(x) = \int_0^2 K_2(x, t) f'''(t) dt,$$

where

$$\begin{aligned} K_2(x, t) &= R_2 \left( \frac{(x-t)_+^2}{2!} \right) \\ &= \frac{1}{2} \left[ (x-t)_+^2 - \left( \frac{1}{4}(x-1)(3-x) \cdot 2(0-t)_+ + (1-t)_+^2 + \frac{1}{4}(x^2-1) \cdot 2(2-t)_+ \right) \right] \\ &= \frac{1}{2} \left[ (x-t)_+^2 - \left( (1-t)_+^2 + \frac{1}{2}(x^2-1)(2-t) \right) \right], \end{aligned}$$

because for  $t \in [0, 2]$ ,  $(-t)_+ = 0$  and  $(2-t)_+ = 2-t$ .



**b)** For the numerical values  $f'(0) = 1$ ,  $f(1) = 2$  and  $f'(2) = 1$ , we have

$$B_2f(x) = \frac{1}{4}(x-1)(3-x) + 2 + \frac{1}{4}(x^2-1) = x+1.$$

Note that, for these particular values, the degree is 1. However, the interpolation formula *still* has degree of exactness  $d = 2$  and everything done in part **a)** still holds.

Then the approximation is

$$f(1/2) \approx B_2f(1/2) = 3/2.$$

We compute the remainder:

$$\begin{aligned} K_2(\textcolor{red}{1}/2, t) &= \frac{1}{2} \left[ (\textcolor{red}{1}/2 - t)_+^2 - \left( (1-t)_+^2 + \frac{1}{2}((\textcolor{red}{1}/2)^2 - 1)(2-t) \right) \right] \\ &= \frac{1}{2} \left( \frac{1}{2} - t \right)_+^2 - \frac{1}{2}(1-t)_+^2 + \frac{3}{16}(2-t). \end{aligned}$$

We have the following cases:

**1.**  $0 \leq t \leq \frac{1}{2}$ , when  $\left(\frac{1}{2} - t\right)_+^2 = \left(\frac{1}{2} - t\right)^2$  and  $(1-t)_+^2 = (1-t)^2$ , so

$$\begin{aligned} K_2(1/2, t) &= \frac{1}{8}(1-2t)^2 - \frac{1}{2}(1-t)^2 + \frac{3}{16}(2-t) \\ &= \frac{1}{16}(2-8t+8t^2 - (8-16t+8t^2) + 3(2-t)) = \frac{5}{16}t \geq 0. \end{aligned}$$

**2.**  $\frac{1}{2} \leq t \leq 1$ , when  $\left(\frac{1}{2} - t\right)_+^2 = 0$ ,  $(1-t)_+^2 = (1-t)^2$  and we have

$$\begin{aligned} K_2(1/2, t) &= -\frac{1}{2}(1-t)^2 + \frac{3}{16}(2-t) \\ &= \frac{1}{16}(- (8-16t+8t^2) + 3(2-t)) = -\frac{1}{16}(8t^2 - 13t + 2) \geq 0, \end{aligned}$$

because the roots of the quadratic polynomial above,  $\frac{13 \pm \sqrt{105}}{16} = \{0.1721, 1.4529\}$ , lie *outside* the interval  $\left[\frac{1}{2}, 1\right]$ , so for  $t \in \left[\frac{1}{2}, 1\right] \subset [0.1721, 1.4529]$ , the quadratic polynomial above has positive sign (opposite to the sign of the leading coefficient).

3.  $1 \leq t \leq 2$ , when  $\left(\frac{1}{2} - t\right)_+^2 = (1 - t)_+^2 = 0$  and

$$K_2(1/2, t) = \frac{3}{16}(2 - t) \geq 0.$$

So, in all three cases,  $K_2$  has constant sign on  $[0, 2]$  and, thus,

$$R_2 f(1/2) = \frac{1}{3!} f'''(\xi) R_2 e_3(1/2) = \frac{1}{6} \cdot \frac{11}{8} f'''(\xi) = \frac{11}{48} f'''(\xi), \quad \xi \in (0, 2).$$

In the end, we have the approximation

$$f(1/2) \approx B_2 f(1/2) = 3/2,$$

with the error

$$|R_2 f(1/2)| \leq \frac{11}{48} \|f'''\|_\infty.$$

■

## 2 Spline Interpolation

Polynomial interpolation has a major setback: the difference between the values of the function  $f$  and the values of the interpolation polynomial *outside* the nodes' interval can be quite large. Choosing more nodes and finding a higher degree polynomial does not solve this problem, but increases the computational cost. So, even though polynomials are smooth and easy to work with functions, they are not always the best choice for approximating functions.

From these considerations came the idea of changing polynomials to *piecewise polynomials* that satisfy some continuity conditions (of the interpolation function and some of its derivatives). Such functions are called *splines*.

Historically, spline functions can be traced all the way back to ancient mathematics. The term “spline” was first used by I. J. Schoenberg in 1946, but a thorough spline function theory started developing in 1964, as their good approximating properties became more evident. They can be used in a large variety of ways in approximation theory, computer graphics, data fitting, numerical integration and differentiation, and the numerical solution of integral, differential, and partial differential

equations.

Over time, there have been several world renowned research groups in spline theory, scattered all over the world. One such group, with remarkable contributions, was a Romanian research group (based especially in Cluj).

The basic idea of approximating a function on an interval  $[a, b]$  with spline functions, is to use different polynomials (of lower degree) on different parts of the interval. The reason for this is the fact that on a sufficiently small interval, functions can be approximated arbitrarily well by polynomials of low degree, even degree 1, or 0.

**Definition 2.1.** Let  $\Delta$  be a grid of the interval  $[a, b]$ ,

$$\Delta : a = x_1 < x_2 < \cdots < x_{n-1} < x_n = b.$$

The set

$$\mathbb{S}_m^k(\Delta) = \{s \mid s \in C^k[a, b], s|_{[x_i, x_{i+1}]} \in \mathbb{P}_m, i = 1, 2, \dots, n-1\} \quad (2.10)$$

is called the **the space of polynomial spline functions of degree  $m$  and class  $k$  on  $\Delta$** .

These are piecewise polynomial functions, of degree  $\leq m$ , continuous at  $x_1, \dots, x_{n-1}$ , together with all their derivatives of order up to  $k$ . In general, we assume  $0 \leq k < m$ . For  $k = m$ ,

$$\mathbb{S}_m^m = \mathbb{P}_m.$$

If  $k = -1$ , we allow discontinuities at the grid points.

## 2.1 Linear Splines

For  $m = 1$  and  $k = 0$ , we have *linear spline functions*. We determine a function  $s_1 \in \mathbb{S}_1^0(\Delta)$  such that

$$s_1(x_i) = f(x_i) = f_i, i = 1, 2, \dots, n.$$

That means that on the interval  $[x_i, x_{i+1}]$ , the function  $s_1$  is the interpolation polynomial of degree 1

$$s_1(f; x) = f_i + f[x_i, x_{i+1}](x - x_i) = f_i + \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i}(x - x_i). \quad (2.11)$$

The graph of this function is shown in Figure 1.

The error is given by

$$|f(x) - s_1(f; x)| \leq \frac{|(x - x_i)(x - x_{i+1})|}{2!} \max_{x \in [x_i, x_{i+1}]} |f''(x)| \leq \frac{h_i^2}{8} \max_{x \in [x_i, x_{i+1}]} |f''(x)|, \quad (2.12)$$

where we denoted by  $h_i = x_{i+1} - x_i$ .

Hence, if  $|\Delta|$  denotes

$$|\Delta| = \max_{i=1, n-1} h_i,$$

we have

$$\|f(\cdot) - s_1(f, \cdot)\|_\infty \leq \frac{|\Delta|^2}{8} \|f''\|_\infty. \quad (2.13)$$

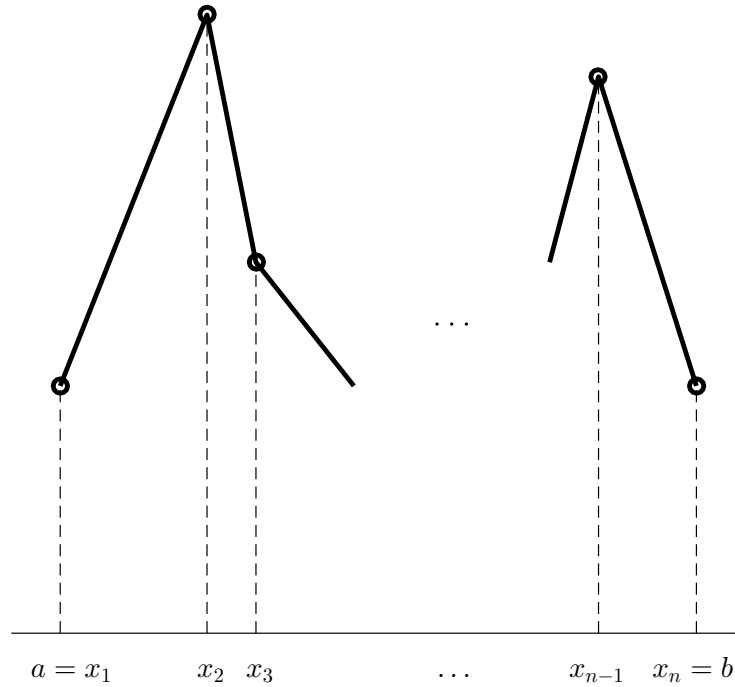


Fig. 1: Linear splines

Obviously,  $\mathbb{S}_1^0(\Delta)$  is a vector space. To find its dimension, we count the number of degrees of freedom and the number of constraints. There are  $n - 1$  subintervals and 2 coefficients to be determined (i.e. 2 degrees of freedom) on each, for a total of  $2(n - 1)$ . We have continuity conditions

at each interior node, so  $n - 2$  constraints. Thus, in the end we have

$$\dim \mathbb{S}_1^0(\Delta) = 2(n - 1) - (n - 2) = n.$$

A basis for this space is given by the so-called *B-spline functions*. Taking  $x_0 = x_1 = a$ ,  $x_{n+1} = x_n = b$ , for  $i = \overline{1, n}$ , we define

$$B_i(x) = \begin{cases} \frac{x - x_{i-1}}{x_i - x_{i-1}}, & \text{If } x_{i-1} \leq x \leq x_i \\ \frac{x_{i+1} - x}{x_{i+1} - x_i}, & \text{If } x_i \leq x \leq x_{i+1} \\ 0, & \text{in rest.} \end{cases} \quad (2.14)$$

Note that the first equation for  $i = 1$ , and the second equation for  $i = n$ , are to be ignored. The functions  $B_i$  are sometimes referred to as “hat functions” (Chinese hats), but note that the first and the last hat are cut in half. Their graphs are depicted in Figure 2. They are linearly independent and have the property

$$B_i(x_j) = \delta_{ij}.$$

Any function  $s \in \mathbb{S}_1^0(\Delta)$  can be written uniquely as

$$s(x) = \sum_{i=1}^n c_i B_i(x).$$

*B-spline functions* play the same role as fundamental Lagrange polynomials  $l_i$ .

A linear spline agrees with the data, but it has the disadvantage of not having a smooth graph. Most data will represent a smooth curved graph, one without the corners of a linear spline. Consequently, we usually want to construct a smooth curve that interpolates the given data points, but one that follows the shape of the linear spline.

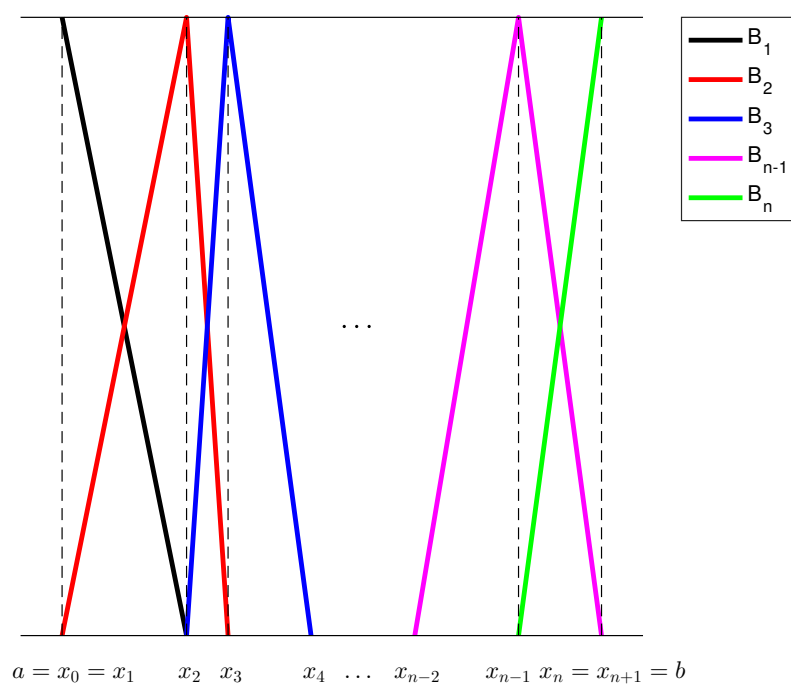


Fig. 2: Linear B-splines

## 2.2 Cubic Splines

Recall the **the space of polynomial spline functions** of degree  $m$  and class  $k$  on  $\Delta$

$$\begin{aligned} \mathbb{S}_m^k(\Delta) &= \{s \mid s \in C^k[a, b], s|_{[x_i, x_{i+1}]} \in \mathbb{P}_m, i = 1, 2, \dots, n-1\}, \\ \Delta &: a = x_1 < x_2 < \dots < x_{n-1} < x_n = b. \end{aligned} \quad (2.1)$$

Now we focus on the case  $m = 3$ .

*Cubic splines* are the most widely used. In general, cubic splines are fairly smooth functions that are convenient to work with. Cubic spline interpolation is a useful tool in mathematical modeling of curves and surfaces of complex geometric shapes in aircraft construction, shipbuilding, production of hydro turbines and many more areas of science and technology. They have also come to be widely used in the past several decades in computer graphics. There are several types of cubic spline functions, depending on the smoothness conditions they satisfy.

### Interpolation with cubic splines $s \in \mathbb{S}_3^1(\Delta)$

We impose the continuity of the first order derivative of  $s_3(f; \cdot)$  by prescribing the values of the first derivative at each node  $x_i, i = 1, 2, \dots, n$ . Given  $n$  arbitrary numbers  $m_1, m_2, \dots, m_n$ , we seek a function  $s_3(f; \cdot)$  that satisfies the conditions

$$\begin{aligned} s_3|_{[x_i, x_{i+1}]} &= p_i(x) \in \mathbb{P}_3, i = 1, 2, \dots, n-1, \\ s_3(f; x_i) &= f_i, i = 1, 2, \dots, n, \\ s'_3(f; x_i) &= m_i, i = 1, 2, \dots, n. \end{aligned} \quad (2.2)$$

This means that on each subinterval  $[x_i, x_{i+1}]$ ,  $s_3(f; \cdot)$  is the unique solution of the Hermite interpolation problem

$$\begin{aligned} p_i(x_i) &= f_i, \quad p_i(x_{i+1}) = f_{i+1}, \\ p'_i(x_i) &= m_i, \quad p'_i(x_{i+1}) = m_{i+1}, i = \overline{1, n-1}. \end{aligned} \quad (2.3)$$

The divided differences are computed from the table

$x_i$	$f_i$	$\longrightarrow$	$m_i$	$\longrightarrow$	$\frac{f[x_i, x_{i+1}] - m_i}{h_i}$	$\longrightarrow$	$\frac{m_{i+1} - 2f[x_i, x_{i+1}] + m_i}{h_i^2}$
		$\nearrow$		$\nearrow$		$\nearrow$	
$x_i$	$f_i$	$\longrightarrow$	$f[x_i, x_{i+1}]$	$\longrightarrow$	$\frac{m_{i+1} - f[x_i, x_{i+1}]}{h_i}$		
		$\nearrow$		$\nearrow$			
$x_{i+1}$	$f_{i+1}$	$\longrightarrow$	$m_{i+1}$				
		$\nearrow$					
$x_{i+1}$	$f_{i+1},$						

Using the Newton form of the Hermite polynomial, we have

$$\begin{aligned}
p_i(x) &= f_i + m_i(x - x_i) + \frac{f[x_i, x_{i+1}] - m_i}{h_i}(x - x_i)^2 \\
&\quad + \frac{m_{i+1} - 2f[x_i, x_{i+1}] + m_i}{h_i^2}(x - x_i)^2(x - x_{i+1}).
\end{aligned}$$

Alternatively, we can write it in Taylor's form around  $x_i$ . Considering that  $x - x_{i+1} = x - x_i - h_i$ , for  $x \in [x_i, x_{i+1}]$ , we get

$$p_i(x) = c_{i,0} + c_{i,1}(x - x_i) + c_{i,2}(x - x_i)^2 + c_{i,3}(x - x_i)^3, \quad (2.4)$$

with

$$\begin{aligned}
c_{i,0} &= f_i, \\
c_{i,1} &= m_i, \\
c_{i,2} &= \frac{f[x_i, x_{i+1}] - m_i}{h_i} - c_{i,3}h_i = \frac{3f[x_i, x_{i+1}] - 2m_i - m_{i+1}}{h_i}, \\
c_{i,3} &= \frac{m_{i+1} - 2f[x_i, x_{i+1}] + m_i}{h_i^2}.
\end{aligned} \quad (2.5)$$

Hence, to compute  $s_3(f; x)$  at a point  $x \in [a, b]$  that is not a node, we first identify the interval  $[x_i, x_{i+1}]$  that contains  $x$ , then compute the coefficients in (2.5) and evaluate the spline using (2.4).

Next, we discuss some possible choices for the parameters  $m_1, m_2, \dots, m_n$ .



### Piecewise cubic Hermite interpolation

Assuming that the derivatives  $f'(x_i)$ ,  $i = 1, \dots, n$ , are known, we choose  $m_i = f'(x_i)$ . This way, we obtain a strictly local scheme, where the polynomial on each subinterval  $[x_i, x_{i+1}]$  is completely determined by the interpolation data at node points inside, independently of the other pieces. The error in this case (see Example 1.4, in Lecture 6) is

$$\|f(\cdot) - s_3(f, \cdot)\|_\infty \leq \frac{1}{384} |\Delta|^4 \|f^{(4)}\|_\infty. \quad (2.6)$$

For equally spaced nodes, we have

$$|\Delta| = (b - a)/(n - 1)$$

and, therefore,

$$\|f(\cdot) - s_3(f, \cdot)\|_\infty = O(n^{-4}), \quad n \rightarrow \infty. \quad (2.7)$$

### Interpolation with cubic splines $s \in \mathbb{S}_3^2(\Delta)$

To have  $s_3(f; \cdot) \in \mathbb{S}_3^2(\Delta)$ , we require continuity of the second derivatives at the nodes, i.e.

$$p''_{i-1}(x_i) = p''_i(x_i), \quad i = 2, \dots, n - 1,$$

which, for the Taylor coefficients in (2.4), means

$$2c_{i-1,2} + 6c_{i-1,3}h_{i-1} = 2c_{i,2}, \quad i = 2, \dots, n - 1. \quad (2.8)$$

Substituting in (2.5), we obtain the linear system

$$h_i m_{i-1} + 2(h_{i-1} + h_i)m_i + h_{i-1}m_{i+1} = b_i, \quad i = 2, \dots, n - 1, \quad (2.9)$$

where

$$b_i = 3 \left( h_i f[x_{i-1}, x_i] + h_{i-1} f[x_i, x_{i+1}] \right). \quad (2.10)$$

Thus, we have a system of  $n - 2$  linear equations with  $n$  unknowns,  $m_1, m_2, \dots, m_n$ . Once  $m_1$  and  $m_n$  are chosen, the system is *tridiagonal* and can be solved efficiently by several methods.

Next, we discuss possible choices for  $m_1$  and  $m_n$ .

**1. Complete (clamped) splines.** We take

$$m_1 = f'(a), \quad m_n = f'(b).$$

For this type of spline, it can be shown that, if  $f \in C^4[a, b]$ , then

$$\|f^{(r)}(\cdot) - s_3^{(r)}(f, \cdot)\|_\infty \leq C_r |\Delta|^{4-r} \|f^{(4)}\|_\infty, \quad r = 0, 1, 2, 3, \quad (2.11)$$

where

$$C_0 = \frac{5}{384}, \quad C_1 = \frac{1}{24}, \quad C_2 = \frac{3}{8},$$

and  $C_3$  depends on the ratio  $|\Delta|/\min_i h_i$ .

**2. Endpoint second derivative splines.** We require

$$s_3''(f, a) = f''(a), \quad s_3''(f, b) = f''(b).$$

These lead to two more equations,

$$\begin{aligned} 2m_1 + m_2 &= 3f[x_1, x_2] - \frac{1}{2}f''(a)h_1, \\ m_{n-1} + 2m_n &= 3f[x_{n-1}, x_n] - \frac{1}{2}f''(b)h_{n-1}. \end{aligned} \quad (2.12)$$

We place the first equation at the beginning of the system (2.9) and the second at the end of it, thus preserving the tridiagonal structure of the system.

**3. Natural cubic splines.** Imposing

$$s_3''(f; a) = s_3''(f; b) = 0,$$

we get the same two equations as above, with  $f''(a) = f''(b) = 0$ :

$$\begin{aligned} 2m_1 + m_2 &= 3f[x_1, x_2], \\ m_{n-1} + 2m_n &= 3f[x_{n-1}, x_n]. \end{aligned} \quad (2.13)$$

Motivation for these boundary conditions can be given by looking at the physics of *bending thin beams of flexible materials* to pass thru the given data. To the left of  $x_1$  and to the right of  $x_n$ , the beam is straight and therefore the second derivatives are zero at the transition points  $x_1$  and  $x_n$ .

The advantage of this type of spline is that it requires only the function values of  $f$  – no derivatives – but the price paid is a decrease in the accuracy to  $O(|\Delta|^2)$  near the endpoints (unless indeed  $f''(a) = f''(b) = 0$ ).

**4. “Not-a-knot” (deBoor) splines.** Here we impose the conditions that the first two pieces and the last two, coincide, i.e.

$$p_1(x) \equiv p_2(x), \quad p_{n-2}(x) \equiv p_{n-1}(x).$$

This means that the first and last interior nodes,  $x_2$  and  $x_{n-1}$ , are both inactive (hence, the name). We get two more equations expressing the continuity of  $s_3'''(f; x)$  at  $x = x_2$  and  $x = x_{n-1}$ . This comes down to the equality of the leading coefficients  $c_{1,3} = c_{2,3}$  and  $c_{n-2,3} = c_{n-1,3}$ . Thus, we get

$$\begin{aligned} h_2^2 m_1 + (h_2^2 - h_1^2)m_2 - h_1^2 m_3 &= \beta_1, \\ h_{n-1}^2 m_{n-2} + (h_{n-1}^2 - h_{n-2}^2)m_{n-1} - h_{n-2}^2 m_n &= \beta_2, \end{aligned} \quad (2.14)$$

where

$$\begin{aligned} \beta_1 &= 2(h_2^2 f[x_1, x_2] - h_1^2 f[x_2, x_3]), \\ \beta_2 &= 2(h_{n-1}^2 f[x_{n-2}, x_{n-1}] - h_{n-2}^2 f[x_{n-1}, x_n]). \end{aligned}$$

Again, we place the first equation at the beginning of the system (2.9) and the second at the end of it. Even so, the resulting system is *no longer* tridiagonal, but it can be transformed into a tridiagonal one, by combining equations 1 and 2, and  $n-1$  and  $n$ , respectively. Consequently, the first and the last equations become

$$\begin{aligned} h_2 m_1 + (h_2 + h_1)m_2 &= \gamma_1, \\ (h_{n-1} - h_{n-2})m_{n-1} + h_{n-2} m_n &= \gamma_2, \end{aligned} \quad (2.15)$$

where

$$\begin{aligned} \gamma_1 &= \frac{1}{h_2 + h_1} \left[ f[x_1, x_2] h_2 (h_1 + 2(h_1 + h_2)) + h_1^2 f[x_2, x_3] \right], \\ \gamma_2 &= \frac{1}{h_{n-1} + h_{n-2}} \left[ h_{n-1}^2 f[x_{n-2}, x_{n-1}] + (2(h_{n-1} + h_{n-2}) + h_{n-1}) h_{n-2} f[x_{n-1}, x_n] \right]. \end{aligned}$$

### Finding cubic splines using the second derivatives

Computational formulas for finding cubic splines  $s \in \mathbb{S}_3^2(\Delta)$  can be derived (in a similar way) when the arbitrary numbers  $M_1, M_2, \dots, M_n$  are given and forced to satisfy the conditions

$$\begin{aligned} s_3|_{[x_i, x_{i+1}]} &= p_i(x) \in \mathbb{P}_3, \quad i = 1, 2, \dots, n-1, \\ s_3(f; x_i) &= f_i, \quad i = 1, 2, \dots, n, \\ s_3''(f; x_i) &= M_i, \quad i = 1, 2, \dots, n. \end{aligned} \tag{2.16}$$

Since  $s_3$  is a cubic polynomial, its second derivative is linear. Hence, on  $[x_i, x_{i+1}]$ , we have

$$s_3''(f; x) = ax + b,$$

satisfying the conditions

$$s_3''(f; x_i) = M_i, \quad s_3''(f; x_{i+1}) = M_{i+1}, \quad i = 1, 2, \dots, n-1.$$

The values  $a$  and  $b$  are determined from the system

$$\begin{cases} ax_i + b = M_i \\ ax_{i+1} + b = M_{i+1} \end{cases}.$$

Integrating successively, then imposing (2.16) and the continuity conditions at the nodes,

$s_3'(f; x_i) = s_3'(f; x_{i+1})$ ,  $i = \overline{1, n-1}$ , we get the linear system

$$h_{i-1} M_{i-1} + 2(h_{i-1} + h_i) M_i + h_i M_{i+1} = 6(f[x_i, x_{i+1}] - f[x_{i-1}, x_i]), \tag{2.17}$$

for  $i = \overline{2, n-1}$ .

The two extra conditions needed for a closed system can be imposed, e.g., on  $M_1$  and  $M_n$ . If  $M_1 = M_n = 0$ , we get the natural cubic spline.

Other conditions can be enforced, such as the continuity of  $s_3'''(f; x)$  at  $x = x_2$  and  $x = x_{n-1}$ , which lead to deBoor cubic splines.

If the first and last equations are

$$\begin{aligned} 2M_1 + M_2 &= 6(f[x_1, x_2] - f'_1), \\ M_{n-1} + 2M_n &= 6(f'_n - f[x_{n-1}, x_n]), \end{aligned} \tag{2.18}$$

where  $f'_1 = f'(a)$ ,  $f'_n = f'(b)$ , then the resulting function is the complete cubic spline.

**Example 2.1.** Find the natural cubic spline that interpolates the data

$x_i$	1	2	4	5
$f_i$	3	5	9	10

**Solution.**

We have  $n = 4$  nodes and  $h_1 = 1, h_2 = 2, h_3 = 1$ .

From (2.9)–(2.13), the linear system for the unknowns  $m_i$ , also called *slopes*, is

$$\begin{cases} 2m_1 + m_2 = 6 \\ 2m_1 + 6m_2 + m_3 = 18 \\ 2m_2 + 6m_3 + 2m_4 = 12 \\ m_3 + 2m_4 = 3 \end{cases}$$

with solution

$$m_1 = \frac{87}{46}, m_2 = \frac{51}{23}, m_3 = \frac{21}{23}, m_4 = \frac{24}{23}.$$

The system (from (2.17) together with the conditions  $M_1 = M_4 = 0$ ) for the *moments*  $M_i$  becomes

$$\begin{cases} M_1 = 0 \\ M_1 + 6M_2 + 2M_3 = 0 \\ 2M_2 + 6M_3 + M_4 = -6 \\ M_4 = 0 \end{cases}$$

whose solution is

$$M_1 = 0, M_2 = \frac{3}{8}, M_3 = -\frac{9}{8}, M_4 = 0.$$

Hence, both ways, we get the natural cubic spline function

$$s_3(x) = \begin{cases} \frac{x^3}{16} - \frac{3x^2}{16} + \frac{17x}{8} + 1, & x \in [1, 2] \\ -\frac{x^3}{8} + \frac{15x^2}{16} - \frac{x}{8} + \frac{5}{2}, & x \in [2, 4] \\ \frac{3x^3}{16} - \frac{45x^2}{16} + \frac{119x}{8} - \frac{35}{2}, & x \in [4, 5] \end{cases}.$$

■

### Minimality properties of cubic spline interpolants

Natural and complete splines have interesting optimality properties. Henceforth, we denote them by  $s_{nat}(f; \cdot)$  and  $s_{compl}(f; \cdot)$ , respectively.

**Theorem 2.2.** *Let  $g \in C^2[a, b]$  be any function that interpolates  $f$  on  $\Delta$ . Then*

$$\int_a^b |s''_{nat}(f; x)|^2 dx \leq \int_a^b |g''(x)|^2 dx, \quad (2.19)$$

*with equality if and only if  $g(\cdot) = s_{nat}(f; \cdot)$ .*

For the next minimality result, we slightly change the subdivision  $\Delta$ . Consider the grid

$$\Delta' : a = x_0 = x_1 < x_2 < \cdots < x_{n-1} < x_n = x_{n+1} = b, \quad (2.20)$$

where the endpoints are *double* nodes. That means that when we use  $\Delta'$ , we interpolate the function values at all interior points, and, both the functional and the derivative values, at the endpoints.

**Theorem 2.3.** *Let  $g \in C^2[a, b]$  be any function that interpolates  $f$  on  $\Delta'$ . Then*

$$\int_a^b |s''_{compl}(f; x)|^2 dx \leq \int_a^b |g''(x)|^2 dx, \quad (2.21)$$

*with equality if and only if  $g(\cdot) = s_{compl}(f; \cdot)$ .*

**Remark 2.4.** Taking  $g(\cdot) = s_{compl}(f; \cdot)$  in Theorem 2.2, we get

$$\int_a^b |s''_{nat}(f; x)|^2 dx \leq \int_a^b |s''_{compl}(f; x)|^2 dx. \quad (2.22)$$

So, in a sense, the natural cubic spline is the “smoothest” interpolant.

**Remark 2.5.** These minimality properties are at the origin of the name “spline”. A *spline* is a flexible strip of wood used in drawing curves (or a musical instrument in that shape).

**Example 2.6.** Consider the function  $f(x) = \arctan x$ ,  $x \in [-2, 2]$  and the nodes  $\{-2, -1, 0, 1, 2\}$ . Figure 1 shows the graphs of the function  $f$ , the nodes and the complete, natural, deBoor and piecewise Hermite cubic splines interpolating  $f$ . In Figure 2 we have the interpolation errors.

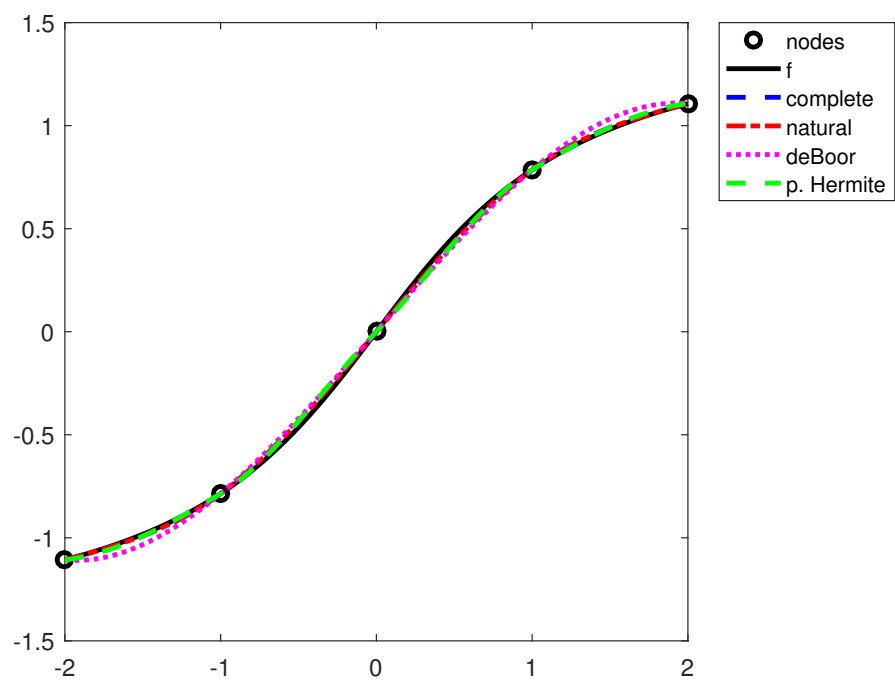


Fig. 1: Interpolation with cubic splines,  $f(x) = \arctan x$

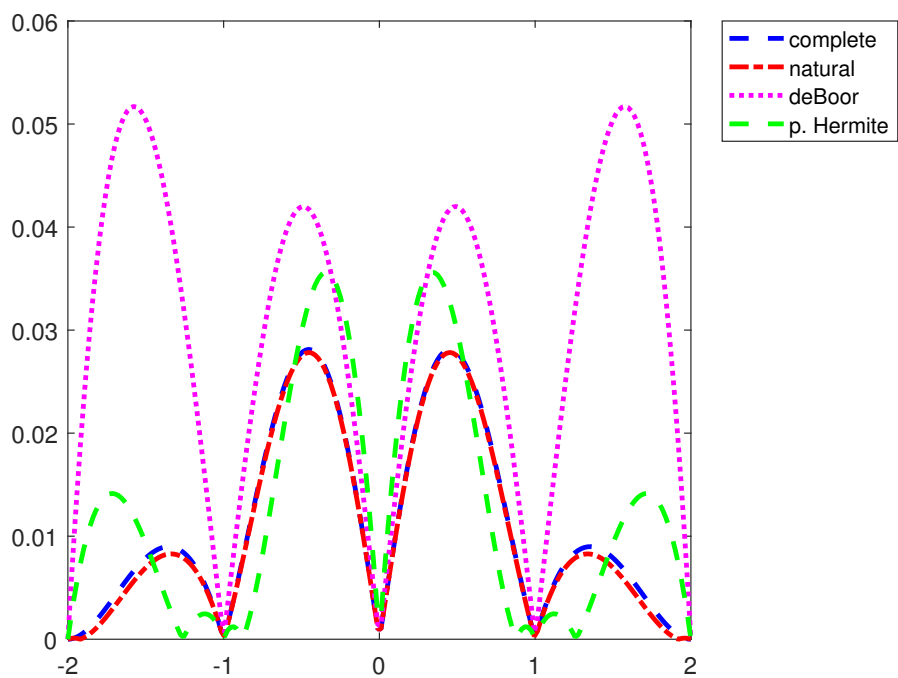


Fig. 2: Errors in cubic spline interpolation,  $f(x) = \arctan x$

## 3 Least Squares Approximation

### 3.1 Best approximation problem

In general, an approximation problem can be described as follows: Let  $f \in X$  be a function,  $\Phi$ , a family of approximants and  $\|\cdot\|$  a norm on  $X$ . We seek an approximation  $\hat{\varphi} \in \Phi$  of  $f$  that approximates the given function “as well as possible”.

$$\|f - \hat{\varphi}\| \leq \|f - \varphi\|, \forall \varphi \in \Phi. \quad (3.1)$$

This is called a *best approximation problem* of  $f$  with elements of  $\Phi$ . The function  $\hat{\varphi}$  is called a *best approximation* of  $f$  relative to the norm  $\|\cdot\|$ . Given a basis  $\{\pi_j\}_{j=1}^m$  of  $\Phi$ , we can write

$$\Phi = \Phi_m = \left\{ \varphi \mid \varphi(t) = \sum_{j=1}^m c_j \pi_j(t), c_j \in \mathbb{R} \right\}. \quad (3.2)$$

$\Phi$  is a finite dimensional linear space or a subset of one.

In the preceding sections we gave a polynomial (i.e., we used  $\Phi = \mathbb{P}_m$ ) or piecewise polynomial (with  $\Phi = S_m^k(\Delta)$ ) approximation based on using interpolation at suitably chosen node points. Another approach is to seek an approximation with a small “average error” over the interval of approximation. That “average error” can be best expressed in terms of inner products and norms.

### 3.2 Scalar Products and Norms

The functions we want to approximate can be defined *continuously*, i.e., on an interval  $[a, b]$ , or *discreetly*, on a set of points  $\{t_1, \dots, t_N\}$ . A *measure* will be defined accordingly, as an *integral* in the continuous case and as a *sum* for discrete functions.

Many measures also involve *weight functions*. An intuitive, physical justification for a *weighted measure* would be that some observations are more important than others, or they are more common, so they “weigh more”.



**Definition 3.1.** A **weight function**  $w$  is defined as follows:

– **continuous case**, a function  $w : [a, b] \rightarrow \mathbb{R}_+$ , satisfying the conditions

- (i)  $\int_a^b |x|^n w(x) dx$  exists and is finite,  $\forall n \geq 0$ ,
- (ii) if  $\int_a^b w(x)g(x) dx = 0$ ,  $g(x) \geq 0$ , then  $g \equiv 0$ ;

– **discrete case**,  $w_i \geq 0$ , satisfying the conditions

- (i)  $\sum_{i=1}^N |t_i|^n w(t_i)$  exists and is finite,  $\forall n \geq 0$ ,
- (ii) if  $\sum_{i=1}^N w_i g_i = 0$ ,  $g_i \geq 0$ , then  $g_i = 0, \forall i = \overline{1, N}$ .

A few commonly used continuous weights:

$$\begin{aligned}
 w(x) &\equiv 1 && \text{on } [-1, 1], \\
 w(x) &= \frac{1}{\sqrt{1-x^2}} && \text{on } [-1, 1], \\
 w(x) &= \sqrt{1-x^2} && \text{on } [-1, 1], \\
 w(x) &= e^{-x} && \text{on } [0, \infty), \\
 w(x) &= e^{-x^2} && \text{on } (-\infty, \infty).
 \end{aligned}$$

**Definition 3.2.** Let  $w$  be a weight function. The **scalar (inner) product** of two functions  $u$  and  $v$  is defined as

$$\begin{aligned}
 \langle u, v \rangle &= \int_a^b w(x)u(x)v(x) dx \text{ for continuous functions and} \\
 \langle u, v \rangle &= \sum_{i=1}^N w_i u_i v_i \text{ in the discrete case.}
 \end{aligned}$$

The **norm** of a function  $u$  is

$$||u|| = (\langle u, u \rangle)^{\frac{1}{2}}.$$

The most frequently used (discrete and continuous) norms are given in Table 1.

Discrete norm	Continuous norm
$\ u\ _p = \left( \sum_{i=1}^N w_i  u(t_i) ^p \right)^{1/p}, \quad p \geq 1$	$\ u\ _p = \left( \int_a^b w(x)  u(x) ^p dx \right)^{1/p}, \quad p \geq 1$
$\ u\ _\infty = \max_{i=1, \dots, m}  u(t_i) $	$\ u\ _\infty = \max_{x \in [a, b]}  u(x) $

Table 1: Commonly used discrete and continuous norms

Let us recall the main properties of scalar products:

- 1. Symmetry:**  $\langle u, v \rangle = \langle v, u \rangle$ ;
- 2. Homogeneity:**  $\langle \alpha u, v \rangle = \langle u, \alpha v \rangle = \alpha \langle u, v \rangle$ ,  $\alpha \in \mathbb{R}$ ;
- 3. Additivity:**  $\langle u + v, z \rangle = \langle u, z \rangle + \langle v, z \rangle$ ;
- 4. Positive definiteness:**  $\langle u, u \rangle \geq 0$  and  $\langle u, u \rangle = 0 \iff u = 0$ ;
- 5. Cauchy–Bunyakovsky–Schwarz inequality:**  $|\langle u, v \rangle| \leq \|u\| \cdot \|v\|$ .

**Definition 3.3.** We say that two functions  $u$  and  $v$  are **orthogonal** if

$$\langle u, v \rangle = 0.$$

More generally, we say that a family of functions  $\{u_k\}_{k=1, \dots, n}$  is an **orthogonal system** if

$$\langle u_i, u_j \rangle = 0, \quad i \neq j.$$

They are called **orthonormal** if

$$\langle u_i, u_j \rangle = \delta_{ij}.$$

### 3.3 Least Squares Approximation

We will consider a particular case for problem (3.1) by choosing the (discrete or continuous) 2-norm. In what follows,  $\|\cdot\|$  will mean  $\|\cdot\|_2$ . So, we want to minimize the *square of the error*

$$\begin{aligned} E^2(\varphi) &= \|\varphi - f\|^2 = \langle \varphi - f, \varphi - f \rangle \\ &= \langle \varphi, \varphi \rangle - 2\langle \varphi, f \rangle + \langle f, f \rangle \\ &= \|\varphi\|^2 - 2\langle \varphi, f \rangle + \|f\|^2. \end{aligned} \tag{3.3}$$

This is then called a *least squares approximation problem* or *mean square approximation problem*. Its solution was given by Gauss and Legendre at the beginning of the 19th century.

#### 3.3.1 Normal equations

Recall that we seek a function  $\varphi \in \Phi_m$  that minimizes  $E^2$  in (3.3). Then

$$\varphi(t) = \sum_{j=1}^m c_j \pi_j(t).$$

Substitute  $\varphi$  into (3.3) to get

$$\begin{aligned} E^2(\varphi) &= \left\langle \sum_{i=1}^m c_i \pi_i, \sum_{j=1}^m c_j \pi_j \right\rangle - 2 \left\langle \sum_{j=1}^m c_j \pi_j, f \right\rangle + \|f\|^2 \\ &= \sum_{i=1}^m \sum_{j=1}^m c_i c_j \langle \pi_i, \pi_j \rangle - 2 \sum_{j=1}^m c_j \langle \pi_j, f \rangle + \|f\|^2. \end{aligned}$$

We want to find the values  $c_j \in \mathbb{R}$  that minimize the error above. We solve this problem (finding the minimum of a function) by taking partial derivatives with respect to each unknown  $c_i$  and setting them equal to 0. Thus, we get

$$\frac{\partial E^2}{\partial c_i} = 2 \sum_{j=1}^m c_j \langle \pi_i, \pi_j \rangle - 2 \langle \pi_i, f \rangle = 0, \quad i = 1, \dots, m,$$

or

$$\sum_{j=1}^m \langle \pi_i, \pi_j \rangle c_j = \langle \pi_i, f \rangle, \quad i = 1, \dots, m. \tag{3.4}$$

The equations in (3.4) are called **normal equations** and they form a linear system

$$Ac = b, \quad (3.5)$$

with

$$a_{ij} = \langle \pi_i, \pi_j \rangle \quad \text{and} \quad b_i = \langle \pi_i, f \rangle. \quad (3.6)$$

Now, since the scalar product is symmetric, so is matrix  $A$ . Also, for every  $x \in \Phi, x \neq 0$ ,

$$\begin{aligned} x^T Ax &= \sum_{i=1}^m \sum_{j=1}^m a_{ij} x_i x_j = \sum_{i=1}^m \sum_{j=1}^m x_i x_j \langle \pi_i, \pi_j \rangle \\ &= \sum_{i=1}^m \sum_{j=1}^m \langle x_i \pi_i, x_j \pi_j \rangle = \left\| \sum_{i=1}^m x_i \pi_i \right\|^2 > 0. \end{aligned}$$

So  $A$  is symmetric and positive definite, therefore, nonsingular. Thus, the system (3.5) has a unique solution  $c_j^*$ ,  $j = 1, \dots, m$ , and hence, so does the least squares approximation problem,

$$\varphi^*(t) = \sum_{j=1}^m c_j^* \pi_j(t). \quad (3.7)$$

**Example 3.4.** Find the linear least squares approximation of the function  $f(x) = \cos \pi t$  on  $[0, 1]$ , using the canonical basis  $\pi_j(t) = t^{j-1}$ ,  $j \in \mathbb{N}^*$ .

**Solution.** The function is given on an interval, so we want the *continuous* least squares approximation. Since we want a *linear* approximation, i.e., a polynomial of degree 1, we have the basis

$$\pi_1(t) = 1, \quad \pi_2(t) = t$$

and we seek an approximation polynomial

$$\varphi(t) = c_1 \pi_1(t) + c_2 \pi_2(t) = a + bt,$$

with simplified notation  $c_1 = a, c_2 = b$ . The normal equations (3.4) are

$$\begin{aligned} \langle \pi_1, \pi_1 \rangle a + \langle \pi_1, \pi_2 \rangle b &= \langle \pi_1, f \rangle, \\ \langle \pi_2, \pi_1 \rangle a + \langle \pi_2, \pi_2 \rangle b &= \langle \pi_2, f \rangle, \end{aligned}$$

with

$$\begin{aligned}
\langle \pi_1, \pi_1 \rangle &= \int_0^1 dt = 1, \quad \langle \pi_1, \pi_2 \rangle = \int_0^1 t dt = 1/2, \quad \langle \pi_2, \pi_2 \rangle = \int_0^1 t^2 dt = 1/3, \\
\langle \pi_1, f \rangle &= \int_0^1 \cos \pi t dt = \frac{1}{\pi} \sin \pi t \Big|_0^1 = 0, \\
\langle \pi_2, f \rangle &= \int_0^1 t \cos \pi t dt = \frac{1}{\pi} t \sin \pi t \Big|_0^1 - \frac{1}{\pi} \int_0^1 \sin \pi t dt = \frac{1}{\pi^2} \cos \pi t \Big|_0^1 = -\frac{2}{\pi^2},
\end{aligned}$$

the last one being integrated by parts ( $u = t$ ,  $dv = \cos \pi t dt$ ). Then we solve the system

$$\begin{cases} a + \frac{1}{2}b = 0, \\ \frac{1}{2}a + \frac{1}{3}b = -\frac{2}{\pi^2} \end{cases},$$

with solution  $a = \frac{12}{\pi^2}$ ,  $b = -\frac{24}{\pi^2}$ .

So, we found the linear least squares approximation

$$\varphi^*(t) = \frac{12}{\pi^2}(1 - 2t).$$

Let us look at the error at some points:

$t$	$f(t)$	$\varphi^*(t)$	$ f(t) - \varphi^*(t) $
0	1	$12/\pi^2$	$12/\pi^2 - 1 \approx 0.22$
1/6	$\sqrt{3}/2$	$8/\pi^2$	$\sqrt{3}/2 - 8/\pi^2 \approx 0.06$
1/4	$\sqrt{2}/2$	$6/\pi^2$	$\sqrt{2}/2 - 6/\pi^2 \approx 0.01$
1/3	1/2	$4/\pi^2$	$1/2 - 4/\pi^2 \approx 0.09$
1/2	0	0	0
1	-1	$-12/\pi^2$	$12/\pi^2 - 1 \approx 0.22$

■

When we seek the *discrete* least squares approximation of some scattered data, the problem is known as *data fitting*, and it arises in many applications.

**Example 3.5.** Find the least squares polynomial approximation that best fits the following data:

$x_i$	-5	-3	-1	1	3	5
$y_i$	4.8	3.0	2.0	2.8	3.2	10

**Solution.** The scatterplot is shown in Figure 3. We see from the graph that the best fit is given by a quadratic function,

$$\varphi(x) = a + bx + cx^2,$$

i.e., the basis is  $1, x, x^2$  (the canonical basis again).

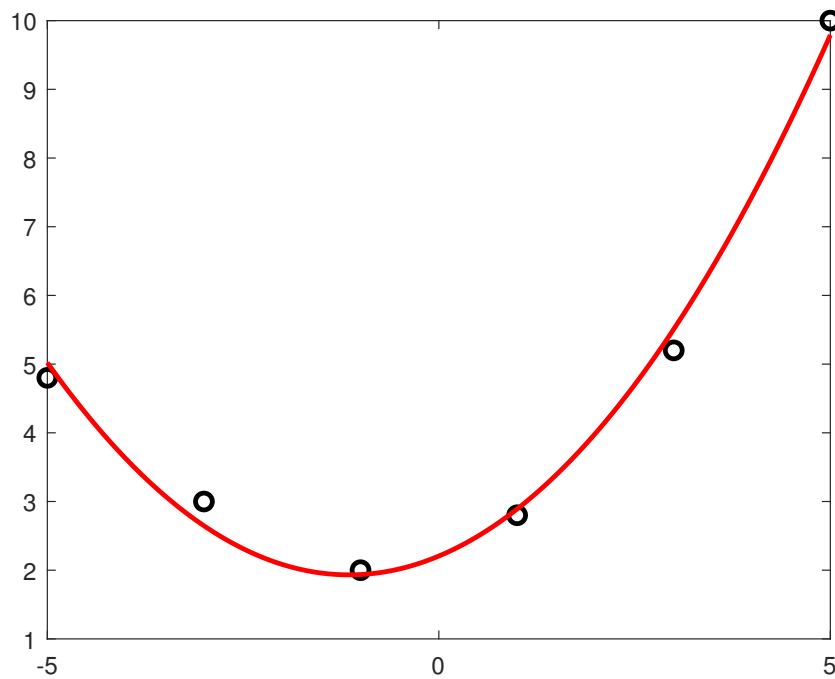


Fig. 3: Data fitting, Example 3.5

The normal equations (3.4) are

$$\begin{aligned} \langle \pi_1, \pi_1 \rangle a + \langle \pi_1, \pi_2 \rangle b + \langle \pi_1, \pi_3 \rangle c &= \langle \pi_1, y \rangle, \\ \langle \pi_2, \pi_1 \rangle a + \langle \pi_2, \pi_2 \rangle b + \langle \pi_2, \pi_3 \rangle c &= \langle \pi_2, y \rangle, \\ \langle \pi_3, \pi_1 \rangle a + \langle \pi_3, \pi_2 \rangle b + \langle \pi_3, \pi_3 \rangle c &= \langle \pi_3, y \rangle, \end{aligned}$$

with

$$\begin{aligned} \langle \pi_i, \pi_j \rangle &= \sum_{k=1}^6 x_k^{i-1} x_k^{j-1} = \sum_{k=1}^6 x_k^{i+j-2} \text{ and} \\ \langle \pi_i, y \rangle &= \sum_{k=1}^6 x_k^{i-1} y_k, \quad i, j = \overline{1, 3}, \quad k = 1, \dots, 6. \end{aligned}$$

So, the normal equations are

$$\begin{aligned} a \sum_{k=1}^6 1 + b \sum_{k=1}^6 x_k + c \sum_{k=1}^6 x_k^2 &= \sum_{k=1}^6 y_k, \\ a \sum_{k=1}^6 x_k + b \sum_{k=1}^6 x_k^2 + c \sum_{k=1}^6 x_k^3 &= \sum_{k=1}^6 x_k y_k, \\ a \sum_{k=1}^6 x_k^2 + b \sum_{k=1}^6 x_k^3 + c \sum_{k=1}^6 x_k^4 &= \sum_{k=1}^6 x_k^2 y_k. \end{aligned}$$

The sums are computed in the following table (on the last row)

	$x_k$	$y_k$	$x_k^2$	$x_k^3$	$x_k^4$	$x_k y_k$	$x_k^2 y_k$
	-5	4.8	25	-125	625	-24.0	120
	-3	3.0	9	-27	81	-9.0	27
	-1	2.0	1	-1	1	2.0	2
	1	2.8	1	1	1	2.8	2.8
	3	5.2	9	27	81	15.6	46.8
	5	10.0	25	125	625	50	250
$\sum_{k=1}^6$	0	27.8	70	0	1414	33.4	448.6

The resulting linear system is

$$\begin{cases} 6a & + & 70c & = & 27.8 \\ & 70b & & = & 33.4 \\ 70a & + & 1414c & = & 448.6 \end{cases},$$

with solution  $a = 2.206$ ,  $b = 0.477$ ,  $c = 0.208$ . The best fit approximation of this data is

$$\varphi^*(x) = 0.208x^2 + 0.477x + 2.206.$$

■

### 3.3.2 Orthogonal polynomials

Least square approximations can use *other functions*, not just polynomials, and bases *other than canonical* can be used, as well. The ideas and procedures described in the previous section still apply.

As we have seen in this chapter, polynomials or piecewise polynomials work quite well and it is rarely the case when other, more complicated approximating functions have to be used.

As far as the basis is concerned, in the continuous case, some choices are better than others and things *can be improved* there. Let us point out a few troublesome aspects:

- For continuous least squares approximations, the linear system  $Ac = b$  in (3.5)-(3.6) can be ill-conditioned. If the canonical basis,  $\pi_j(t) = t^{j-1}$ ,  $j = \overline{1, m}$ , is used on the interval  $[0, 1]$  (as in Example 3.4), then

$$a_{ij} = \langle \pi_i, \pi_j \rangle = \int_0^1 t^{i+j-2} dt = \frac{1}{i+j-1}, \quad i, j = \overline{1, m},$$

i.e.,  $A = H_m$ , the Hilbert matrix, which is known to be very ill-conditioned. The basis functions become almost linearly dependent, as the exponent grows. The solution of the linear system (3.5) is extremely sensitive to small changes in the coefficients or right-hand constants and as a consequence, when  $m \geq 4$ , the solutions will be completely unsatisfactory.

- Another disadvantage is that all the coefficients  $c_j$  found this way depend on  $m$ ,  $c_j = c_j^{(m)}$ . Increasing  $m$  will produce an enlarged system of normal equations with a *completely new* solution vector. There is no relation between  $c_j^{(m)}$  and  $c_j^{(m+n)}$ , so no way of using previous computations.

Both these problems can be overcome if the basis  $\{\pi_j\}_{j=1}^m$  is chosen to be **orthogonal**. If  $\langle \pi_i, \pi_j \rangle = 0$ ,  $i \neq j$ , then the coefficients  $a_{ij} = 0$ ,  $i \neq j$ , which means the system  $Ac = b$  is *diagonal* with solution

$$c_j^* = \frac{\langle \pi_j, f \rangle}{\langle \pi_j, \pi_j \rangle}, \quad j = 1, \dots, m \quad (3.8)$$

and the least squares approximation is

$$\varphi^*(t) = \sum_{j=1}^m \frac{\langle \pi_j, f \rangle}{\langle \pi_j, \pi_j \rangle} \pi_j(t). \quad (3.9)$$



Now, instead of solving a system of normal equations, we can use formula (3.8) directly. Obviously, the coefficients  $c_j^*$  are independent of  $m$  and once computed, they remain the same for any larger  $m$ . We now have what is called *permanence of the coefficients*.

Another aspect: recall from linear algebra that any linearly independent system can be orthogonalized using the *Gram-Schmidt procedure*. So using an orthogonal basis is *not* restrictive at all.

In fact, applying that procedure to the canonical basis  $1, t, t^2, \dots$  on an interval  $[a, b]$ , with respect to an appropriate weight function  $w$ , several well-known families of orthogonal polynomials can be obtained, including the Chebyshev polynomials (of the first and second kind) that were already discussed. Also, that procedure provides a linear recurrence relation between 3 consecutive such orthogonal polynomials:

$$\pi_{k+1}(t) = (t - \alpha_k)\pi_k(t) - \beta_k\pi_{k-1}(t), \quad k = 0, 1, \dots, \quad \pi_{-1}(t) = 0, \quad \pi_0(t) = 1, \quad (3.10)$$

where

$$\alpha_k = \frac{\langle t\pi_k, \pi_k \rangle}{\|\pi_k\|^2}, \quad k = 0, 1, \dots, \quad \beta_k = \frac{\|\pi_k\|^2}{\|\pi_{k-1}\|^2}, \quad k = 1, 2, \dots, \quad \beta_0 = \mu_0. \quad (3.11)$$

Such examples are given in Table 2.

Name	Notation	Polynomial	Weight fn.	Interval	$\alpha_k$	$\beta_k$
Legendre	$l_m$	$[(x^2 - 1)^m]^{(m)}$	1	$[-1, 1]$	0	$\beta_0 = 2,$ $\beta_k = (4 - k^{-2})^{-1}, k \geq 1$
Chebyshev 1 <sup>st</sup>	$T_m$	$\cos(m \arccos x)$	$(1 - x^2)^{-\frac{1}{2}}$	$[-1, 1]$	0	$\beta_0 = \pi,$ $\beta_1 = \frac{1}{2},$ $\beta_k = \frac{1}{4}, k \geq 2$
Chebyshev 2 <sup>nd</sup>	$Q_m$	$\frac{\sin[(m+1) \arccos x]}{\sqrt{1-x^2}}$	$(1 - x^2)^{\frac{1}{2}}$	$[-1, 1]$	0	$\beta_0 = \frac{\pi}{2},$ $\beta_k = \frac{1}{4}, k \geq 1$
Laguerre	$L_m^a$	$x^{-a}e^x (x^{m+a}e^{-x})^{(m)}$	$x^a e^{-x}, a > -1$	$[0, \infty)$	$2k + a + 1$	$\beta_0 = \Gamma(1+a),$ $\beta_k = k(k+a), k \geq 1$
Hermite	$H_m$	$(-1)^m e^{x^2} (e^{-x^2})^{(m)}$	$e^{-x^2}$	$\mathbb{R}$	0	$\beta_0 = \sqrt{\pi},$ $\beta_k = \frac{k}{2}, k \geq 1$

Table 2: Orthogonal polynomials and recurrence coefficients

**Example 3.6.** Let  $f : [-1, 1] \rightarrow [0, \pi], f(t) = \arccos t$ . Find the least squares polynomial approximation  $\varphi^*$  of  $f$  relative to the weight function  $w(t) = \frac{1}{\sqrt{1-t^2}}$ .

**Solution.** For this weight function, we use Chebyshev polynomials of the first kind,  $\pi_j(t) = T_j(t)$ . Recall that

$$\langle \pi_i, \pi_j \rangle = \int_{-1}^1 \frac{T_i(t)T_j(t)}{\sqrt{1-t^2}} dt = \begin{cases} 0, & i \neq j \\ \pi, & i = j = 0 \\ \frac{\pi}{2}, & i = j \neq 0 \end{cases},$$

so

$$\langle \pi_0, \pi_0 \rangle = \pi \text{ and } \langle \pi_j, \pi_j \rangle = \frac{\pi}{2}, j \neq 0.$$

Now, compute the numerators of the coefficients  $c_j^*$ .

$$\langle \pi_j, f \rangle = \int_{-1}^1 \frac{1}{\sqrt{1-t^2}} \cos(j \arccos t) \arccos t dt.$$

With the change of variables  $t = \cos x$ , we have  $\arccos t = x$ ,  $\sqrt{1-t^2} = \sin x$ ,  $dt = -\sin x dx$  and the interval  $[-1, 1]$  is mapped into  $[\pi, 0]$ . So,

$$\langle \pi_j, f \rangle = \int_0^\pi \frac{1}{\sin x} x \cos(jx) \sin x dx = \int_0^\pi x \cos(jx) dx.$$

Then,

$$\langle \pi_0, f \rangle = \int_0^\pi x dx = \frac{1}{2}x^2 \Big|_0^\pi = \frac{\pi^2}{2}.$$

For  $j \neq 0$ , we use integration by parts with  $u = x$ ,  $dv = \cos(jx) dx$  (so  $du = dx$ ,  $v = \frac{1}{j} \sin(jx)$ ):

$$\begin{aligned} \langle \pi_j, f \rangle &= \int_0^\pi x \cos(jx) dx = \frac{1}{j} x \sin(jx) \Big|_0^\pi - \frac{1}{j} \int_0^\pi \sin(jx) dx \\ &= 0 + \frac{1}{j^2} \cos(jx) \Big|_0^\pi = \frac{1}{j^2} (\cos(j\pi) - \cos 0) = \frac{1}{j^2} ((-1)^j - 1) \\ &= \begin{cases} 0, & j \text{ even} \\ -\frac{2}{j^2}, & j \text{ odd} \end{cases}. \end{aligned}$$

So,

$$c_0^* = \frac{\pi}{2}, \quad c_j^* = \begin{cases} 0, & j \neq 0 \text{ even} \\ -\frac{4}{j^2\pi}, & j \text{ odd} \end{cases}.$$

Then the least squares approximating polynomial is given by

$$\varphi_{2n+1}^*(t) = \frac{\pi}{2} - \sum_{i=0}^n \frac{4}{(2i+1)^2\pi} T_{2i+1}(t).$$

This is the solution of the least squares problem:

$$\int_{-1}^1 \frac{1}{\sqrt{1-t^2}} (f(t) - \varphi_n^*(t))^2 dt = \min \left\{ \int_{-1}^1 \frac{1}{\sqrt{1-t^2}} (f(t) - \varphi(t))^2 dt : \varphi \in \mathbb{P}_n \right\}.$$

The approximation is quite good, even for small degrees, as seen in Figure 4.

■

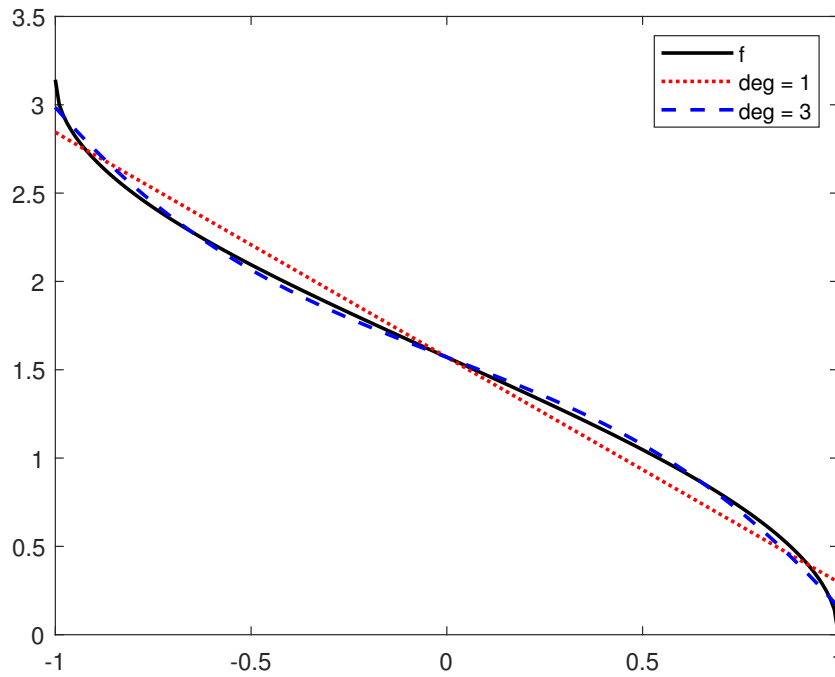


Fig. 4: Example 3.6

# Chapter 4. Numerical Differentiation and Integration

## 1 Approximation of Linear Functionals, Basic Notions

Let  $X$  be a linear space and  $L, L_1, \dots, L_m : X \rightarrow \mathbb{R}$  be real, linear functionals, that are linearly independent.

**Definition 1.1.** An approximation formula of  $L$  using  $L_1, \dots, L_m$ , is a formula of the type

$$L(f) = \sum_{i=1}^m A_i L_i(f) + R(f), \quad f \in X. \quad (1.1)$$

The real parameters  $A_i$  are called **coefficients**, and  $R(f)$  is the **remainder** of the formula.

For an approximation formula of the form (1.1), given  $L_i$ , we want to determine the coefficients  $A_i$  and study the corresponding remainder (error).

The functionals  $L_i$  express the available information on  $f$  and they also depend on the particular type of approximation we seek, i.e. on  $L$ .

**Example 1.2.** Let  $X = \{f \mid f : [a, b] \rightarrow \mathbb{R}\}$ ,  $L_i(f) = f(x_i)$ , for some distinct nodes  $x_i \in [a, b]$ ,  $i = \overline{0, m}$  and  $L(f) = f(\alpha)$ , for an arbitrary  $\alpha \in [a, b]$ . Formula (1.1) becomes

$$f(\alpha) = \sum_{i=0}^m l_i(\alpha) f(x_i) + (Rf)(\alpha),$$

i.e. the *Lagrange interpolation formula*. We have

$$A_i = l_i(\alpha),$$

where  $l_i$  are the Lagrange fundamental polynomials. One of the expressions for the remainder is

$$(Rf)(\alpha) = \frac{u(\alpha)}{(m+1)!} f^{(m+1)}(\xi), \quad \xi \in [a, b], \quad u(x) = (x - x_0) \dots (x - x_m),$$

if  $f^{(m+1)}$  exists on  $[a, b]$ .

**Example 1.3.** Let  $X$  and  $L_i$  be defined as in the previous example. Assuming that  $f^{(k)}(\alpha)$ ,  $k \in \mathbb{N}^*$  exists, define  $L(f) = f^{(k)}(\alpha)$ . We get an approximation formula for the derivative of order  $k$  of  $f$

at  $\alpha$ ,

$$f^{(k)}(\alpha) = \sum_{i=0}^m A_i f(x_i) + R(f),$$

called a *numerical differentiation formula*.

**Example 1.4.** Let  $x_k \in [a, b]$ ,  $k = \overline{0, m}$  be distinct nodes and  $I_k$  some sets of indices. Consider  $X = \{f \mid f : [a, b] \rightarrow \mathbb{R}, f \text{ integrable on } [a, b], \text{ for which } f^{(j)}(x_k), k = \overline{0, m}, j \in I_k \text{ exist}\}$ ,  $L_{kj}(f) = f^{(j)}(x_k)$  and  $L(f) = \int_a^b f(x)dx$ . Formula (1.1) becomes

$$\int_a^b f(x)dx = \sum_{k=0}^m \sum_{j \in I_k} A_{kj} f^{(j)}(x_k) + R(f),$$

called a *numerical integration (quadrature) formula*.

In general, there are two approaches for solving the approximation problem (1.1):

- the **interpolation method**: apply the functional  $L$  to a suitable interpolation polynomial of  $f$ , instead of  $f$  itself;
- the **method of undetermined coefficients**: find the coefficients in (1.1), by making the remainder  $R(f)$  be 0 for polynomials of degree as high as possible, i.e., imposing the conditions  $R(e_k) = 0$ ,  $e_k(x) = x^k$ ,  $k = 0, 1, \dots, d$ , for as large a  $d$  as possible.

## 2 Numerical Differentiation

Numerical approximation of derivatives is used when the values of a function  $f$  are given in tables, as empirical data, or the expression of  $f$  is complicated.

We can derive simple, immediate numerical differentiation rules using divided and finite differences. Let  $f : [a, b] \rightarrow \mathbb{R}$  be differentiable on  $[a, b]$ ,  $x \in [a, b]$ , arbitrary and  $h > 0$ . We have

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = \lim_{h \rightarrow 0} f[x, x+h].$$

From here, we immediately get the approximation

$$f'(x) \approx \frac{f(x+h) - f(x)}{h} \equiv D_h f(x), \quad (2.1)$$

called the *forward difference numerical derivative*.

Expanding  $f(x + h)$  in a Taylor's series around  $x$ , we get

$$\begin{aligned} f(x + h) &= f(x) + hf'(x) + \frac{h^2}{2}f''(\xi), \\ \frac{f(x + h) - f(x)}{h} &= f'(x) + \frac{h}{2}f''(\xi), \quad \xi \in (x, x + h), \end{aligned}$$

from which we have the error formula

$$(RD_h f)(x) = f'(x) - D_h f(x) = -\frac{h}{2}f''(\xi), \quad \xi \in (x, x + h). \quad (2.2)$$

The error is proportional to  $h$ , so formula (2.2) can be used for small steps  $h$ .

Similarly, we obtain the *backward difference numerical derivative*,

$$f'(x) \approx \frac{f(x) - f(x - h)}{h} \equiv \tilde{D}_h f(x), \quad (2.3)$$

with approximation error

$$(R\tilde{D}_h f)(x) = f'(x) - \frac{f(x) - f(x - h)}{h} = \frac{h}{2}f''(\xi), \quad \xi \in (x - h, x). \quad (2.4)$$

Interpolating  $f$  at the nodes  $x - h$ ,  $x + h$  and then taking the derivative, we obtain

$$f'(x) \approx \frac{f(x + h) - f(x - h)}{2h} \equiv \hat{D}_h f(x), \quad (2.5)$$

known as the *central difference numerical derivative formula*, with remainder given by

$$(R\hat{D}_h f)(x) = -\frac{h^2}{6}f'''(\xi), \quad \xi \in (x - h, x + h). \quad (2.6)$$

This says that for small values of  $h$ , the formula (2.5) should be more accurate than the earlier approximations, because the error term of (2.6) decreases more rapidly with  $h$ .

**Example 2.1.** Use  $D_h f$  and  $\hat{D}_h f$  to approximate the derivative of  $f(x) = \cos x$  at  $x = \pi/6$ . Study the error of each approximation.

**Solution.** The exact value is  $f'\left(\frac{\pi}{6}\right) = -\sin \frac{\pi}{6} = -\frac{1}{2}$ .

By (2.2), when using  $D_h f$ , the error is

$$(RD_h f)\left(\frac{\pi}{6}\right) = f'\left(\frac{\pi}{6}\right) - D_h\left(\frac{\pi}{6}\right) = \frac{h}{2} \cos \xi,$$

thus,

$$\left|(RD_h f)\left(\frac{\pi}{6}\right)\right| \leq \frac{h}{2}.$$

Similarly, for  $\widehat{D}_h f$ , the error is bounded by

$$\left|(R\widehat{D}_h f)\left(\frac{\pi}{6}\right)\right| \leq \frac{h^2}{6}.$$

Table 1 contains the approximation results yielded by the two methods, for various values of  $h$ . Indeed, both the value of  $D_h f$  and that of  $\widehat{D}_h f$  are approaching  $-0.5$ . Moreover, looking at the errors, we see that when  $h$  is halved, the error is almost halved (see the first ratio column) for the first approximation. This confirms the fact that the error is proportional to  $h$  (relation (2.2)). For the approximation  $\widehat{D}_h f$ , we see that the errors decrease more rapidly and the last column of ratios confirms the (superior) rate of convergence of  $O(h^2)$  given in (2.6).

$h$	$D_h f$	Error	Ratio	$\widehat{D}_h f$	Error	Ratio
0.1	-0.54243	$4.243e-2$		-0.49917	$-8.329e-4$	
0.05	-0.52144	$2.144e-2$	1.98	-0.49979	$-2.083e-4$	4.00
0.025	-0.51077	$1.077e-2$	1.99	-0.49995	$-5.208e-5$	4.00
0.0125	-0.50540	$5.403e-3$	1.99	-0.49998	$-1.302e-5$	4.00
0.00625	-0.50270	$2.701e-3$	2.00	-0.49999	$-3.255e-6$	4.00

Table 1: Example 2.1,  $f(x) = \cos x$

■

**Remark 2.2.** One must be very cautious in using numerical differentiation, because of the sensitivity to errors in the function values. This is especially true if the function values are obtained empirically with relatively large experimental errors, as is common in practice. Numerical differentiation is an *unstable* operation, meaning that even if the approximation of a function is good, that *does not* guarantee that its derivative will be a good approximation for the derivative of the function.

Here is such an example: Let

$$f(x) = g(x) + \frac{x^{n^2}}{n}, \quad n \geq 1, \quad x \in [0, 1], \quad f, g \in C[0, 1].$$

Notice that

$$\begin{aligned} \|f - g\|_\infty &= \max_{x \in [0, 1]} \frac{x^{n^2}}{n} = \frac{1}{n} \rightarrow 0, \quad n \rightarrow \infty, \\ \|f' - g'\|_\infty &= \max_{x \in [0, 1]} nx^{n^2-1} = n \rightarrow \infty. \end{aligned}$$

Numerical derivatives can be used to find numerical methods for (ordinary or partial) differential equations. This is done in order to reduce the differential equation to a form that can be solved more easily than the original equation.

### 3 Numerical Integration

Let  $f : [a, b] \rightarrow \mathbb{R}$  be integrable on  $[a, b]$ ,  $F_k(f)$ ,  $k = \overline{0, m}$  give information on  $f$  (usually, linear functionals, such as values or derivatives) and let  $w : [a, b] \rightarrow \mathbb{R}_+$  be a weight function which is integrable on  $[a, b]$ .

**Definition 3.1.** A formula of the type

$$\int_a^b w(x)f(x)dx = \sum_{j=0}^m A_j F_j(f) + R(f), \quad (3.1)$$

is called a **numerical integration formula** for the function  $f$  or a **quadrature formula**. The parameters  $A_j$ ,  $j = \overline{0, m}$  are called the **coefficients** of the formula, and  $R(f)$  the **remainder**.

The weight function can be very useful in, among other things, “absorbing” any singularities the integrand has on  $[a, b]$  (since on the right-hand-side *only* values related to  $f$  are used).

**Definition 3.2.** The natural number  $d$  satisfying the property that  $\forall f \in \mathbb{P}_d, R(f) = 0$  and  $\exists g \in \mathbb{P}_{d+1}$  such that  $R(g) \neq 0$  is called **degree of precision** (or **degree of exactness**) of the quadrature formula (3.1).



**Remark 3.3.** Since  $R$  is a linear functional, it follows that a quadrature formula has degree of precision  $d$  if and only if

$$R(e_j) = 0, \quad j = 0, 1, \dots, d, \quad R(e_{d+1}) \neq 0. \quad (3.2)$$

If the degree of precision of a quadrature formula is known, then the remainder can be determined using Peano's Theorem.

### 3.1 Interpolatory Quadratures, Newton-Cotes Formulas

For now, we will restrict our discussion to the case  $w(x) \equiv 1$ . Many numerical integration formulas are based on the idea of replacing  $f$  by an approximating function whose integral can be evaluated. Most of the times, that approximating function is an interpolation polynomial. Then, we obtain a quadrature formula of the form

$$\int_a^b f(x)dx = \sum_{k=0}^m A_k f(x_k) + R(f), \quad (3.3)$$

called an **interpolatory quadrature**. If, in addition, the nodes used are equally spaced, it is called a **Newton-Cotes quadrature**. If the nodes include the endpoints of the interval,  $a$  and  $b$ , then we have a *closed* Newton-Cotes formula, otherwise, an *open* one.

There are  $2m + 2$  unknowns ( $m + 1$  nodes and  $m + 1$  coefficients) in formula (3.3). Imposing conditions (3.2), it follows that the maximum possible degree of precision can be obtained for a polynomial with  $2m + 2$  coefficients, i.e. of degree  $2m + 1$ , hence,  $e_{2m+1}$ . Thus, the maximum degree of precision of a quadrature formula (3.3) with  $m + 1$  nodes is

$$d_{\max} = 2m + 1 = 2 * (\text{nr. of nodes}) - 1.$$

Any interpolatory numerical integration scheme (3.3) has degree of precision at least  $m$  (since the interpolation formula has that degree of exactness).

We start with three of the most widely used (but also, simplest) quadratures, obtained from low degree polynomial interpolation.

### Rectangle (Midpoint) Rule

We interpolate  $f$  at a single *double* node,  $x_0 = \frac{a+b}{2}$ , the midpoint of the interval (hence, the name of the method). So we use the Taylor polynomial of degree 1. Assuming that  $f$  has second order continuous derivatives on  $(a, b)$ , we have

$$\begin{aligned} f(x) &= T_1 f(x) + R_1 f(x) \\ &= f\left(\frac{a+b}{2}\right) + \left(x - \frac{a+b}{2}\right) f'\left(\frac{a+b}{2}\right) + \frac{1}{2!} \left(x - \frac{a+b}{2}\right)^2 f''(\xi), \quad \xi \in (a, b). \end{aligned}$$

Integrating, we get

$$\begin{aligned} \int_a^b f(x) dx &= (b-a) f\left(\frac{a+b}{2}\right) + f'\left(\frac{a+b}{2}\right) \int_a^b \left(x - \frac{a+b}{2}\right) dx + R(f) \\ &= (b-a) f\left(\frac{a+b}{2}\right) + f'\left(\frac{a+b}{2}\right) \frac{1}{2} \left(x - \frac{a+b}{2}\right)^2 \Big|_a^b + R(f) \\ &= (b-a) f\left(\frac{a+b}{2}\right) + R(f), \end{aligned}$$

because the second integral is  $\frac{1}{2} \left[ \left(\frac{b-a}{2}\right)^2 - \left(\frac{b-a}{2}\right)^2 \right] = 0$ .

Check the conditions (3.2). We have

$$\begin{aligned} R(e_0) &= \int_a^b e_0(x) dx - (b-a) e_0\left(\frac{a+b}{2}\right) = b-a - (b-a) = 0, \\ R(e_1) &= \int_a^b x dx - (b-a) \frac{a+b}{2} = \frac{b^2-a^2}{2} - \frac{b^2-a^2}{2} = 0. \end{aligned}$$

So, we found the formula

$$\int_a^b f(x) dx = (b-a) f\left(\frac{a+b}{2}\right) + R(f), \quad (3.4)$$

called the **rectangle rule**, an *open* Newton-Cotes formula, having degree of precision  $d = 1$ , which is the *maximum* possible for a formula with a single node ( $m = 0$ ).

We compute the remainder by

$$R(f) = \frac{f''(\xi)}{2!} \int_a^b \left(x - \frac{a+b}{2}\right)^2 dx = \frac{(b-a)^3}{24} f''(\xi), \quad \xi \in (a, b). \quad (3.5)$$

Let us see a geometrical interpretation of this formula. Recall that, if  $f(x) \geq 0$  for  $x \in [a, b]$ , the definite integral in (3.4) represents the area of the region that lies below the graph of  $f(x)$ , above the  $Ox$  axis and between the lines  $x = a$  and  $x = b$ . This area is approximated by the area of the *rectangle* with base  $b - a$  and height  $f\left(\frac{a+b}{2}\right)$  (see Figure 1). Hence, the other name of the method.

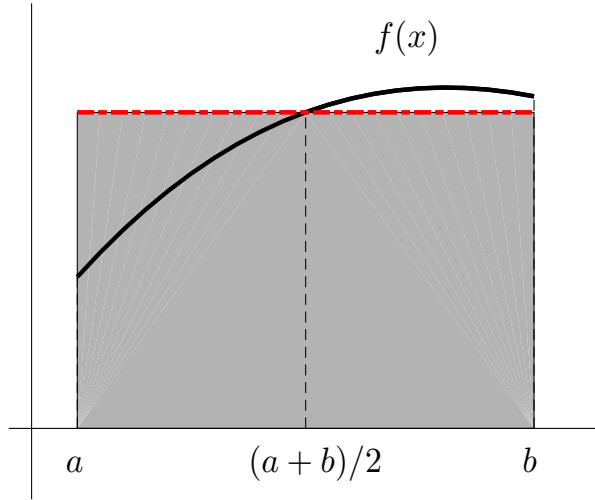


Fig. 1: Geometrical illustration of the rectangle rule

**Remark 3.4.** The rectangle rule (3.4) can also be obtained using the *method of undetermined coefficients*. We seek a quadrature formula with one node, i.e., of the form

$$\int_a^b f(x) dx = A_0 f(x_0) + R(f).$$

Then impose conditions (3.2) and go as far as possible.

**Solution.** First, we want  $R(e_0) = 0$ , which means

$$\int_a^b e_0(x) dx = A_0 e_0(x_0), \text{ i.e., } \int_a^b dx = A_0.$$

We get the first equation,  $A_0 = b - a$ .

Then, from  $R(e_1) = 0$ , we obtain

$$\int_a^b e_1(x) dx = A_0 e_1(x_0), \text{ i.e., } \int_a^b x dx = A_0 x_0,$$

so, the second equation is  $A_0 x_0 = \frac{b^2 - a^2}{2}$ . The two equations have the solution

$$\begin{aligned} A_0 &= b - a, \\ x_0 &= \frac{a + b}{2}. \end{aligned}$$

Can we go further? Let's check.

$$\begin{aligned} R(e_2) &= \int_a^b e_2(x) dx - (b - a) e_2\left(\frac{a + b}{2}\right) = \int_a^b x^2 dx - (b - a) \left(\frac{a + b}{2}\right)^2 \\ &= \frac{b^3 - a^3}{3} - (b - a) \frac{(a + b)^2}{4} = \frac{(b - a)^3}{12} \neq 0, \end{aligned}$$

so the degree of precision is  $d = 1$ . From here, we can obtain the expression of the remainder (3.5) using **Peano's theorem**. Let us recall this important result and see how it is used for quadratures formulas. For a numerical integration formula

$$\int_a^b f(x) dx = Q(f) + R(f),$$

with **degree of precision**  $d = n$ , assuming  $f \in C^{n+1}[a, b]$ , the remainder has the form

$$R(f) = \int_a^b K_n(t) f^{(n+1)}(t) dt,$$

with

$$\begin{aligned}
K_n(t) &= \frac{1}{n!} R_n((x-t)_+^n) = \frac{1}{n!} \left[ \int_a^b (x-t)_+^n dx - Q((x-t)_+^n) \right] \\
&= \frac{1}{n!} \left[ \frac{1}{n+1} (x-t)_+^{n+1} \Big|_{x=a}^{x=b} - Q((x-t)_+^n) \right] \\
&= \frac{1}{n!} \left[ \frac{(b-t)_+^{n+1} - (a-t)_+^{n+1}}{n+1} - Q((x-t)_+^n) \right].
\end{aligned}$$

If  $K_n$  has constant sign on  $[a, b]$ , then

$$R(f) = \frac{1}{(n+1)!} f^{(n+1)}(\xi) R(e_{n+1}), \quad \xi \in (a, b).$$

So, for the midpoint formula

$$\int_a^b f(x) dx = (b-a) f\left(\frac{a+b}{2}\right) + R(f),$$

with degree of precision  $d = 1$ , we have

$$R(f) = \int_a^b K_1(t) f''(t) dt,$$

with

$$\begin{aligned}
K_1(t) &= \frac{1}{1!} R((x-t)_+) = \int_a^b (x-t)_+ dx - (b-a) \left( \frac{a+b}{2} - t \right)_+ \\
&= \frac{1}{2} [(b-t)_+^2 - (a-t)_+^2] - (b-a) \left( \frac{a+b}{2} - t \right)_+.
\end{aligned}$$

Now, since  $t \in [a, b]$ , it follows that  $(b-t)_+ = b-t$  and  $(a-t)_+ = 0$ . The third term depends on the sign of  $\frac{a+b}{2} - t$ . So, we have two cases:

1.  $a \leq t \leq \frac{a+b}{2}$ , when  $\left(\frac{a+b}{2} - t\right)_+ = \frac{a+b}{2} - t$  and, hence,

$$\begin{aligned} K_1(t) &= \frac{1}{2}(b-t)^2 - (b-a) \left(\frac{a+b}{2} - t\right) = \frac{1}{2}b^2 - bt + \frac{1}{2}t^2 - \frac{1}{2}(b^2 - a^2) + bt - at \\ &= \frac{1}{2}t^2 - at + \frac{1}{2}a^2 = \frac{1}{2}(t-a)^2 \geq 0; \end{aligned}$$

2.  $\frac{a+b}{2} < t \leq b$ , when  $\left(\frac{a+b}{2} - t\right)_+ = 0$  and we have

$$K_1(t) = \frac{1}{2}(b-t)^2 \geq 0.$$

So, in both cases,  $K_1$  has a constant sign over  $[a, b]$ , and thus, the remainder can be expressed as

$$R(f) = \frac{1}{2!}f''(\xi)R(e_2) = \frac{(b-a)^3}{24}f''(\xi), \quad \xi \in (a, b),$$

as in (3.5). ■

To improve on the approximation of the integral, break the interval  $[a, b]$  into  $n$  smaller subintervals determined by the equidistant nodes  $x_i = a + ih, i = \overline{0, n}, h = (b-a)/n$ , and apply the rectangle rule (3.4) on each subinterval, i.e.,

$$\int_{x_i}^{x_{i+1}} f(x)dx = hf\left(\frac{x_i + x_{i+1}}{2}\right) + \frac{h^3}{24}f''(\xi_i), \quad \xi_i \in [x_i, x_{i+1}].$$

We have

$$\int_a^b f(x)dx = \sum_{i=0}^{n-1} \int_{x_i}^{x_{i+1}} f(x)dx = h \sum_{i=0}^{n-1} f\left(\frac{x_i + x_{i+1}}{2}\right) + \frac{h^3}{24} \sum_{i=0}^{n-1} f''(\xi_i), \quad \xi_i \in [x_i, x_{i+1}].$$

Using a mean value formula for the continuous function  $f''$ ,

$$f''(\xi) = \frac{f''(\xi_0) + \cdots + f''(\xi_{n-1})}{n}, \quad \xi \in (a, b),$$

we get

$$\int_a^b f(x)dx = h \sum_{i=0}^{n-1} f\left(a + \left(i + \frac{1}{2}\right)h\right) + \frac{h^2(b-a)}{24} f''(\xi), \quad \xi \in (a, b), \quad (3.6)$$

called the **composite (repeated) rectangle (midpoint) formula**.

### Trapezoidal Rule

We proceed similarly, approximating the integrand by the Lagrange interpolation polynomial with 2 nodes,  $x_0 = a, x_1 = b$ , the endpoints of the interval. If  $f$  is twice continuously differentiable on  $(a, b)$ , we have

$$f(x) = \frac{x-b}{a-b}f(a) + \frac{x-a}{b-a}f(b) + \frac{f''(\xi)}{2!}(x-a)(x-b), \quad \xi \in (a, b).$$

Integrating, after doing all the computations, we get

$$\int_a^b f(x)dx = \frac{b-a}{2} \left( f(a) + f(b) \right) - \frac{(b-a)^3}{12} f''(\xi), \quad \xi \in (a, b), \quad (3.7)$$

called the **trapezoidal (or trapezium) rule**, a *closed* Newton-Cotes formula. Again, the name comes from the geometrical interpretation (see Figure 2), where the area of the region that lies between the graph of  $f$ , the  $x$ -axis and the lines  $x = a$  and  $x = b$ , is approximated by the area of the trapezoid with bases  $f(a), f(b)$  and height  $b - a$ .

Since this rule is derived from Lagrange interpolation with two nodes (the degree of the interpolation polynomial being 1), we know that its degree of precision is *at least*  $d = 1$  (without checking  $R(e_0) = R(e_1) = 0$ ). Let us check if  $d > 1$ .

$$R(e_2) = \int_a^b x^2 dx - \frac{b-a}{2}(a^2 + b^2) = \frac{1}{3}(b^3 - a^3) - \frac{b-a}{2}(a^2 + b^2) = -\frac{(b-a)^3}{6} \neq 0.$$

Thus, the degree of precision is  $d = 1$ .

Now, just as we did with the rectangle rule, we divide the interval  $[a, b]$  into  $n$  subintervals

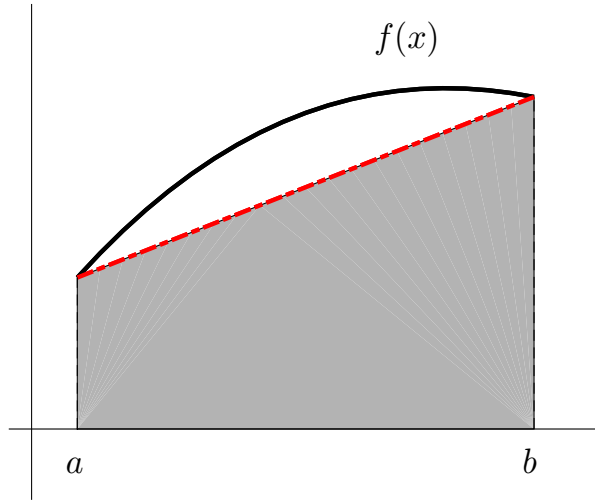


Fig. 2: Geometrical illustration of the trapezoidal rule

$[x_i, x_{i+1}]$ ,  $x_i = a + ih$ ,  $i = \overline{0, n}$ , of length  $h = \frac{b-a}{n}$ . We have

$$\int_a^b f(x)dx = \sum_{i=0}^{n-1} \int_{x_i}^{x_{i+1}} f(x)dx = \frac{h}{2} \sum_{i=0}^{n-1} (f(x_i) + f(x_{i+1})) - \frac{h^3}{12} \sum_{i=0}^{n-1} f''(\xi_i), \quad \xi_i \in [x_i, x_{i+1}]$$

Using again the mean value theorem and denoting by  $f_i = f(x_i)$ , we get the **composite (repeated) trapezoidal (trapezium) rule**,

$$\int_a^b f(x)dx = \frac{h}{2} [f(a) + 2(f_1 + \cdots + f_{n-1}) + f(b)] - \frac{h^2(b-a)}{12} f''(\xi), \quad \xi \in (a, b). \quad (3.8)$$

**Remark 3.5.** Obviously, for larger  $n$ , we get increasingly accurate approximations of the definite integral. But which sequence of values of  $n$  should be used? If  $n$  is doubled repeatedly,  $n \rightarrow 2n$ , then the function values used in each approximation (3.8) will include all of the earlier function values used in the preceding approximation. Thus, the *doubling* of  $n$  will ensure that all previously computed information is used in the new calculation, making the trapezoidal rule less expensive than it would be otherwise.



## Simpson's Rule

For this formula, we consider Hermite interpolation at the nodes  $x_0 = a, x_1 = \frac{a+b}{2}$ , *double* and  $x_2 = b$ . Then the corresponding Hermite interpolation polynomial has degree 3 and is of the form

$$\begin{aligned} H_3(x) &= f(x_0) + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_1](x - x_0)(x - x_1) + \\ &+ f[x_0, x_1, x_1, x_2](x - x_0)(x - x_1)^2. \end{aligned}$$

If  $f$  has continuous derivatives of order 4 on  $[a, b]$ , the error of the approximation can be written as

$$R_3(x) = \frac{(x - x_0)(x - x_1)^2(x - x_2)}{4!} f^{(4)}(\xi), \quad \xi \in (a, b).$$

Integrating on  $[a, b]$  the relation  $f(x) = H_3(x) + R_3(x)$ , we get a new closed Newton-Cotes formula,

$$\int_a^b f(x)dx = \frac{b-a}{6} \left[ f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right] - \frac{(b-a)^5}{2880} f^{(4)}(\xi), \quad \xi \in (a, b), \quad (3.9)$$

called the **(Cavalieri-) Simpson rule**. Its degree of precision is  $d = 3$ .

Dividing the interval  $[a, b]$  into an *even* number  $n = 2m$  of subintervals of length  $h = \frac{b-a}{2m}$ , and denoting by  $x_i = a + ih, f_i = f(x_i), i = \overline{0, 2m}$ , we have

$$\begin{aligned} \int_a^b f(x)dx &= \sum_{i=1}^m \int_{x_{2i-2}}^{x_{2i}} f(x)dx \\ &= \sum_{i=1}^m \left[ \frac{h}{3} (f_{2i-2} + 4f_{2i-1} + f_{2i}) - \frac{h^5}{90} f^{(4)}(\xi_i) \right], \quad \xi_i \in [x_{2i-2}, x_{2i}]. \end{aligned}$$

By the mean value theorem, we get the **composite (repeated) Simpson's rule**

$$\begin{aligned} \int_a^b f(x)dx &= \frac{h}{3} \left[ f(a) + 4 \sum_{i=1}^m f_{2i-1} + 2 \sum_{i=1}^{m-1} f_{2i} + f(b) \right] \\ &- \frac{h^4(b-a)}{180} f^{(4)}(\xi), \quad \xi \in (a, b). \end{aligned} \quad (3.10)$$

**Remark 3.6.**

1. The trapezoidal and Simpson's rules can also be derived using the method of undetermined coefficients for a two-point and three-point, respectively, quadrature formula.
2. Simpson's formula can be derived by considering interpolation with 3 *simple* nodes, so a polynomial of degree 2. We get the same coefficients, but the integral of the remainder will be zero. This is why Hermite interpolation was used instead.
3. These are three of the simplest quadrature formulas. The rectangle and trapezoidal rules are comparable precision-wise ( $O(h^2)$ ) and also from the computational cost point of view (number of flops per iteration). The trapezoidal rule is usually preferred when the number of nodes is doubled at each iteration (see Remark 3.5). Simpson's rule is superior in precision ( $O(h^4)$ ), but it also incurs a higher computational load.

**Example 3.7.** Approximate the integral

$$\int_0^1 \frac{1}{1+x} dx$$

using the three methods above.

**Solution.** The exact value of the integral is

$$\int_0^1 \frac{1}{1+x} dx = \ln(1+x) \Big|_0^1 = \ln 2 = 0.693147180559945.$$

By the rectangle rule, we have the approximation

$$\int_0^1 \frac{1}{1+x} dx \approx 1 \cdot f\left(\frac{1}{2}\right) = \frac{2}{3} = 0.6667,$$

with error  $E_1 = 0.0265$ . Using the trapezoidal rule, we obtain

$$\int_0^1 \frac{1}{1+x} dx \approx \frac{1}{2}(f(0) + f(1)) = \frac{3}{4} = 0.75,$$

with error  $E_2 = -0.0569$ . Finally, with Simpson's rule, we get

$$\int_0^1 \frac{1}{1+x} dx \approx \frac{1}{6} \left[ f(0) + 4f\left(\frac{1}{2}\right) + f(1) \right] = \frac{25}{36} = 0.6944,$$

with approximation error  $E_3 = -0.0013$ . ■

**Example 3.8.** Let us approximate

$$\int_0^1 e^{-x^2} dx = 0.746824132812427,$$

with the composite trapezoidal and Simpson's rules.

**Solution.** The approximation errors (as well as the ratios of successive approximations) for the two methods are given in Table 2, for various values of  $n$ . These confirm the higher rate of convergence,  $O(h^4)$ , of Simpson's repeated method over the composite trapezoidal rule.

$n$	Composite Trapezoidal		Repeated Simpson	
	Error	Ratio	Error	Ratio
2	$1.55e-2$		$3.56e-4$	
4	$3.84e-3$	4.02	$3.12e-5$	11.4
8	$9.59e-4$	4.01	$1.99e-6$	15.7
16	$2.40e-4$	4.00	$1.25e-7$	15.9
32	$5.99e-5$	4.00	$7.79e-9$	16.0
64	$1.50e-5$	4.00	$4.87e-10$	16.0
128	$3.74e-6$	4.00	$3.04e-11$	16.0

Table 2: Example 3.8 ■

Let us see another example of obtaining a quadrature formula two ways.

**Example 3.9.** Consider a quadrature formula of the type

$$\int_{-1}^1 f(x) dx = Af'(-1) + Bf(1) + R(f).$$

- a) Find  $A$  and  $B$  such that the formula has the maximum degree of exactness possible.
- b) Let  $B_1 f$  be the Birkhoff polynomial interpolating  $f$ , given  $f'(-1)$  and  $f(1)$ . Compute  $\int_{-1}^1 (B_1 f)(x) dx$  and compare it to the formula found in a).
- c) Express the remainder  $R(f)$  in the form

$$R(f) = \text{const} \cdot f''(\xi), \xi \in (-1, 1).$$

**Solution.**

- a) We set  $R(e_k) = 0$ ,  $e_k(x) = x^k$  and go as far as possible.

$$\begin{aligned} R(e_0) &= 2 - [A \cdot 0 + B \cdot 1] = 2 - B = 0, \\ R(e_1) &= 0 - [A \cdot 1 + B \cdot 1] = -A - B = 0. \end{aligned}$$

From these two equations, we get  $A = -2$  and  $B = 2$ . Check further:

$$R(e_2) = \frac{2}{3} - [-2 \cdot (-2) + 2 \cdot 1] = \frac{2}{3} - 6 = -\frac{16}{3} \neq 0,$$

so the maximum degree of precision possible is  $d = 1$  and the quadrature formula is

$$\int_{-1}^1 f(x) dx = 2(-f'(1) + f(1)) + R(f).$$

- b) For Birkhoff interpolation, we have the nodes  $x_0 = -1, x_1 = 1$  and  $I_0 = \{1\}, I_1 = \{0\}$ . Then the degree of the polynomial is  $n = 1 + 1 - 1 = 1$ . The polynomial

$$B_1 f(x) = ax + b$$

must satisfy the interpolation conditions

$$\begin{cases} (B_1 f)'(-1) = f'(-1) \\ (B_1 f)(1) = f(1) \end{cases} \iff \begin{cases} a = f'(-1) \\ a + b = f(1) \end{cases} \iff \begin{cases} a = f'(-1) \\ b = f(1) - f'(-1) \end{cases}.$$

So, we have

$$\begin{aligned}
f(x) &= (B_1 f)(x) + (R_1 f)(x) \\
&= x f'(-1) + (f(1) - f'(-1)) + (R_1 f)(x) \\
&= (x - 1) f'(-1) + f(1) + (R_1 f)(x) \\
&= b_{01}(x) f'(-1) + b_{10}(x) f(1) + (R_1 f)(x).
\end{aligned}$$

Integrating, we get

$$\begin{aligned}
\int_{-1}^1 f(x) dx &= \int_{-1}^1 (B_1 f)(x) dx + R(f) \\
&= f'(-1) \int_{-1}^1 (x - 1) dx + f(1) \int_{-1}^1 dx + R(f) \\
&= -2 f'(-1) + 2 f(1) + R(f),
\end{aligned}$$

the same quadrature formula as before.

c) The degree of precision is  $d = 1$ . Then,

$$R(f) = \int_{-1}^1 K_1(t) f''(t) dt,$$

with

$$\begin{aligned}
K_1(t) &= R((x - t)_+) = \int_{-1}^1 (x - t)_+ d\mathbf{x} - 2 \left[ - \frac{\partial (x - t)_+}{\partial \mathbf{x}} \Big|_{\mathbf{x}=-1} + (1 - t)_+ \right] \\
&= \frac{1}{2} (\mathbf{x} - t)_+^2 \Big|_{\mathbf{x}=-1}^{\mathbf{x}=1} - 2 \left[ -1 + (1 - t)_+ \right] \\
&= \frac{1}{2} \left( (1 - t)_+^2 - (-1 - t)_+^2 \right) + 2 - 2(1 - t)_+.
\end{aligned}$$

Since  $t \in [-1, 1]$ ,  $(1 - t)_+ = 1 - t$ ,  $(-1 - t)_+ = 0$  and further we have

$$K_1(t) = \frac{1}{2}(1 - t)^2 + 2 - 2(1 - t) = \frac{1}{2}(1 - t)^2 + 2t = \frac{1}{2}(1 + t)^2 \geq 0.$$

So,

$$R(f) = \frac{1}{2}f''(\xi)R(e_2) = \frac{1}{2}\left(-\frac{16}{3}\right)f''(\xi) = -\frac{8}{3}f''(\xi), \quad \xi \in (-1, 1).$$

Alternatively, we can find the remainder in the Birkhoff interpolation formula, using Peano's Theorem:

$$(R_1f)(x) = \frac{(x - 1)(x + 3)}{2}f''(\xi), \quad \xi \in (-1, 1).$$

(try it, it's a good exercise). Then, integrating, we get

$$R(f) = \int_{-1}^1 (R_1f)(x) dx = -\frac{8}{3}f''(\xi), \quad \xi \in (-1, 1),$$

same as before. ■

### 3.2 Adaptive Quadratures

As seen so far, the errors in numerical integration methods depend not only on the size of the interval, but also on values of certain higher order derivatives of the function to be integrated. Newton-Cotes methods (including the three simple ones, that use low degree polynomial interpolation) work well for smooth integrands (even with a small number of nodes), but perform poorly for functions having large values of higher order derivatives – especially for functions having large oscillations on some subintervals or on the whole interval. As a simple example, consider

$$\int_0^1 \sqrt{x} dx = \frac{2}{3}.$$

This integrand has infinite derivative at  $x = 0$ , but is smooth at points close to  $x = 1$ .

Generally, numerical integration schemes use evenly spaced nodes. When the function to be integrated has a singularity at some point  $\alpha \in [a, b]$ , this requires many nodes in the vicinity of that

point, to reduce the errors caused by the chaotic behaviour of the function in that neighborhood. But this implies that many more nodes (more than necessary) are used throughout the *entire* interval of integration, increasing (unnecessarily) the computational cost of the method. Ideally, we want to use small subintervals where the derivatives are large, and larger subintervals where the derivatives are small and well-behaved.

A method that does this systematically is called **adaptive quadrature**. The general approach in an adaptive quadrature is to use two different methods on each subinterval, compare the results, and divide the interval when the differences are large. The structure of such an algorithm would be “Divide and conquer”.

In Algorithm 3.1 we present an example of a general structure for a recursive adaptive quadrature. The parameter “met” is a function that implements a composite quadrature rule, such as the trapezoidal or Simpson’s rule, and  $m$  is the number of subintervals.

Unlike other methods, that decide what amount of work is needed to achieve a desired precision, an adaptive quadrature computes only as much as is necessary.

**Algorithm 3.1.** [Adaptive quadrature]

```
function  $I = \text{adquad}(f, a, b, \varepsilon, \text{met}, m)$ 
     $I1 = \text{met}(f, a, b, m);$ 
     $I2 = \text{met}(f, a, b, 2m);$ 
    if  $|I1 - I2| < \varepsilon$  % success
         $I = I2;$ 
        return
    else % recursive subdivision
         $I = \text{adquad}(f, a, \frac{a+b}{2}, \varepsilon, \text{met}, m) + \text{adquad}(f, \frac{a+b}{2}, b, \varepsilon, \text{met}, m);$ 
    end
end
```

### 3.3 Iterated Quadratures; Romberg's Method

#### 3.3.1 Richardson Extrapolation

*Extrapolation* is a method for generating high-accuracy numerical schemes using low-order formulas. The most widely used is *Richardson extrapolation*.

Consider the integral

$$I := \int_a^b f(x) dx$$

and a numerical integration scheme

$$I \approx I_n,$$

for which we have an asymptotic error formula of the form

$$I - I_n \approx \frac{c}{n^p}, \quad (3.1)$$

where  $c$  depends on  $a, b$  and the derivatives of a certain order of the function  $f$  on  $[a, b]$ . The difficulty in using this estimate is not knowing the value of the constant  $c$ . We can obtain a computable estimate of the error without needing to know  $c$  explicitly. We write (3.1) for a larger  $n$ :

$$I - I_{2n} \approx \frac{c}{(2n)^p} = \frac{c}{2^p n^p} \quad (3.2)$$

and eliminate the unknown  $c$  from relations (3.1)-(3.2). We obtain

$$I - I_n \approx 2^p (I - I_{2n}),$$

and then the approximation

$$I \approx \frac{2^p I_{2n} - I_n}{2^p - 1} = I_{2n} + \frac{I_{2n} - I_n}{2^p - 1} \stackrel{\text{not}}{=} R_{2n}, \quad (3.3)$$

called **Richardson's extrapolation formula**. From this, we can get another error estimate for  $I_{2n}$ ,

$$I - I_{2n} \approx \frac{I_{2n} - I_n}{2^p - 1}, \quad (3.4)$$

called **Richardson's error estimate**. The term  $R_{2n}$  is an improved estimate of  $I$ , based on using



$I_n, I_{2n}, p$  and the assumption (3.1). It is a more accurate approximation to  $I$  than is  $I_{2n}$ . How much more accurate it is depends on the validity of (3.1)–(3.2).

**Example 3.1.** Let us consider a few simple Newton-Cotes formulas.

**Solution.** For the composite trapezoidal rule, if  $f$  has continuous second order derivatives on  $[a, b]$ , we have

$$\int_a^b f(x)dx = \frac{b-a}{n} \left[ f(a) + f(b) + \sum_{i=1}^{n-1} f\left(a + \frac{b-a}{n}i\right) \right] - \frac{(b-a)^3}{12n^2} f''(\xi) = T_n + \frac{c}{n^2},$$

hence,  $p = 2$ . Using Richardson extrapolation, we get

$$I \approx \frac{4T_{2n} - T_n}{3} = T_{2n} + \frac{1}{3}(T_{2n} - T_n) \quad (3.5)$$

and the error formula

$$I - T_{2n} \approx \frac{1}{3}(T_{2n} - T_n). \quad (3.6)$$

With Simpson's repeated rule, if  $f$  has continuous fourth order derivatives on  $[a, b]$  and  $f_j = f\left(a + \frac{b-a}{2n}j\right)$ ,  $j = \overline{0, 2n}$ , we have

$$\int_a^b f(x)dx = \frac{b-a}{6n} \left[ f(a) + f(b) + 4 \sum_{i=1}^n f_{2i-1} + 2 \sum_{i=1}^{n-1} f_{2i} \right] - \frac{(b-a)^5}{2880n^4} f^{(4)}(\xi) = S_n + \frac{c}{n^4},$$

so in this case,  $p = 4$ . With Richardson extrapolation, we obtain

$$I \approx \frac{16S_{2n} - S_n}{15} = S_{2n} + \frac{1}{15}(S_{2n} - S_n) \quad (3.7)$$

and the error

$$I - S_{2n} \approx \frac{1}{15}(S_{2n} - S_n). \quad (3.8)$$

■

**Example 3.2.** Consider again the problem in Example 3.8 in Lecture 9: the approximation of the

integral

$$I = \int_0^1 e^{-x^2} dx = 0.746824132812427,$$

using the composite trapezoidal rule.

**Solution.** Last time we obtained the values

$n$	Approx. value $T_n$	Error	Ratio
2	0.7313702518	$1.55e-2$	
4	0.7429840978	$3.84e-3$	4.02
8	0.7458656148	$9.59e-4$	4.01
16	0.7465845968	$2.40e-4$	4.00
32	0.7467642547	$5.99e-5$	4.00

Table 1: Approximations and errors for repeated trapezoidal rule, Example 3.2

We have

$$T_2 = 0.7313702518,$$

$$T_4 = 0.7429840978.$$

By (3.5), we get the approximation

$$I \approx R_4 = \frac{1}{3}(4T_4 - T_2) = 0.7468553798,$$

with absolute error

$$0.0000312 = 3.12e-5.$$

Notice in Table 1 that the error of  $R_4$  is smaller than that of  $T_{32}$ , so  $R_4$  (after only 2 steps) gives a better approximation of  $I$  than  $T_{32}$ , obtained after 5 steps!

Now, let us estimate the error in  $T_4$  using Richardson extrapolation. By (3.6), we have

$$I - T_4 \approx \frac{1}{3}(T_4 - T_2) = 0.00387.$$

The actual error in  $T_4$  is 0.00384, so we obtained a very accurate error estimate.

■

**Remark 3.3.** Richardson's extrapolation and error estimation are not always as accurate as this example might suggest, but it is usually a fairly accurate procedure. The main assumption that must be satisfied is (3.1), with a known  $p$ . And the extrapolation itself provides a way of testing whether this assumption is valid for the actual values of  $I_n$  being used: Continue the ideas in (3.1)–(3.3) and write successively

$$\begin{aligned} I - I_n &\approx 2^p(I - I_{2n}), \\ I - I_{2n} &\approx 2^p(I - I_{4n}). \end{aligned}$$

We get

$$\begin{aligned} I_{2n} - I_n &= (I - I_n) - (I - I_{2n}) \\ &\approx 2^p(I - I_{2n}) - (I - I_{2n}) = (2^p - 1)(I - I_{2n}). \end{aligned}$$

Similarly, we have

$$\begin{aligned} I_{4n} - I_{2n} &= (I - I_{2n}) - (I - I_{4n}) \\ &\approx (I - I_{2n}) - 2^{-p}(I - I_{2n}) = (1 - 2^{-p})(I - I_{2n}). \end{aligned}$$

Then,

$$\frac{I_{2n} - I_n}{I_{4n} - I_{2n}} \approx \frac{2^p - 1}{1 - \frac{1}{2^p}} = 2^p.$$

We obtained the (computable) estimate

$$2^p \approx \frac{I_{2n} - I_n}{I_{4n} - I_{2n}},$$

or

$$p \approx \log_2 \left( \frac{I_{2n} - I_n}{I_{4n} - I_{2n}} \right) = \frac{1}{\ln 2} \ln \left( \frac{I_{2n} - I_n}{I_{4n} - I_{2n}} \right). \quad (3.9)$$

This gives a practical means of checking/finding the value of  $p$  in (3.1), using three successive values  $I_n, I_{2n}, I_{4n}$ .

**Example 3.4.** Let us use the approximations in Table 1 to estimate the value of  $p$  for which

$$I - T_n \approx \frac{c}{n^p}.$$

**Solution.** We use (3.9).

For  $n = 1$ , we get the estimate

$$p_1 = \log_2 \left( \frac{T_4 - T_2}{T_8 - I_4} \right) = 2.0109.$$

If we use  $n = 2$ , we have

$$p_2 = \log_2 \left( \frac{T_8 - T_4}{T_{16} - I_8} \right) = 2.0028.$$

And, finally, if  $n = 4$ , we obtain

$$p_3 = \log_2 \left( \frac{T_{16} - T_8}{T_{32} - I_{16}} \right) = 2.0007.$$

We see that the estimates converge to  $p = 2$ , which is consistent with the theoretical value determined for the composite trapezoidal rule. ■

### 3.3.2 Iterated Quadratures; Romberg's Method

Just like in the case of Lagrange interpolation, we want algorithms for which it is easy to go from one step (iteration) to the next, by using previously computed values of the function.

One drawback of adaptive quadratures is that they compute repeatedly the function values at the nodes and when such an algorithm is executed, there is an extra computational cost due to recursion. *Iterated quadratures* overcome this shortcoming. They apply at the first step a composite quadrature rule and then divide the interval into equal parts using at each step the previously computed approximations. **Romberg's method** is such an iterative algorithm, starting with the composite trapezoidal (or midpoint) rule and then improving the convergence by using Richardson extrapolation.

The initial approximations are obtained by applying either the trapezoid or midpoint rule with  $n_k = 2^{k-1}$ ,  $k \in \mathbb{N}$ . Then the value of the step  $h_k$  is

$$h_k = \frac{b - a}{n_k} = \frac{b - a}{2^{k-1}}.$$

With these notations, we have (for the trapezium rule)

$$\int_a^b f(x)dx = \frac{h_k}{2} \left[ f(a) + f(b) + 2 \sum_{i=1}^{2^{k-1}-1} f(a + ih_k) \right] - \frac{b-a}{12} h_k^2 f''(\xi_k), \quad \xi_k \in [a, b].$$

Denote by  $R_{k,1}$  the approximation above, i.e.,

$$\begin{aligned} R_{1,1} &= \frac{h_1}{2} [f(a) + f(b)] = \frac{b-a}{2} [f(a) + f(b)], \\ R_{2,1} &= \frac{h_2}{2} [f(a) + f(b) + 2f(a + h_2)] \\ &= \frac{b-a}{4} [f(a) + f(b) + 2f(a + \frac{1}{2}h_1)] \\ &= \frac{1}{2} \left[ \frac{b-a}{2} (f(a) + f(b)) + (b-a)f(a + \frac{1}{2}h_1) \right] \\ &= \frac{1}{2} \left[ R_{1,1} + h_1 f(a + \frac{1}{2}h_1) \right] \\ &\dots \\ R_{k,1} &= \frac{1}{2} \left[ R_{k-1,1} + h_{k-1} \sum_{i=1}^{2^{k-2}} f(a + (i - \frac{1}{2})h_{k-1}) \right], \quad k = \overline{2, n}. \end{aligned} \tag{3.10}$$

Since  $h_k = \frac{1}{2}h_{k-1}$ , each successive level of improvement increases the order of the error term from  $O(h^{2k-2})$  to  $O(h^{2k})$ , so by

$$O(h^2) = O\left(\frac{1}{n^2}\right).$$

Then we can use Richardson extrapolation with  $p = 2$ , by eliminating the term in  $h_k^2$  from the approximation of  $I$  by  $R_{k-1,1}$  and  $R_{k,1}$ , respectively. We obtain

$$I = \frac{4R_{k,1} - R_{k-1,1}}{3} + O(h_k^4)$$

and define

$$R_{k,2} = \frac{4R_{k,1} - R_{k-1,1}}{3}. \tag{3.11}$$

We apply Richardson extrapolation to these values, too. In general, if  $f \in C^{2n+2}[a, b]$ , then, for  $k = \overline{1, n}$ , we can write

$$\int_a^b f(x)dx = \frac{h_k}{2} \left[ f(a) + f(b) + 2 \sum_{i=1}^{2^{k-1}-1} f(a + ih_k) \right] + \sum_{i=1}^k K_i h_k^{2i} + O(h_k^{2k+2}),$$

where  $K_i$  does not depend on  $h_k$ .

Successively eliminating the powers of  $h$  from the relation above, we get

$$R_{k,j} = \frac{4^{j-1} R_{k,j-1} - R_{k-1,j-1}}{4^{j-1} - 1}, \quad k = \overline{2, n}, \quad j = \overline{2, k}. \quad (3.12)$$

The computations can be arranged in a table (from (3.10) and (3.12)):

$$\begin{array}{ccccccc} R_{1,1} & & & & & & \\ R_{2,1} & R_{2,2} & & & & & \\ R_{3,1} & R_{3,2} & R_{3,3} & & & & \\ \vdots & \vdots & \vdots & \ddots & & & \\ R_{n,1} & R_{n,2} & R_{n,3} & \dots & R_{n,n} & & \end{array}$$

If the sequence  $\{R_{n,1}\}_n$  (which is just the repeated trapezium rule) converges, so does  $\{R_{n,n}\}_n$ , but at a *faster* rate. We can use the stopping criterion

$$|R_{n-1,n-1} - R_{n,n}| < \varepsilon.$$

**Remark 3.5.** The second column in Romberg's method corresponds to Simpson's composite rule. We introduce the notation

$$S_{k,1} = R_{k,2}.$$

Then, the values in the third column are

$$R_{k,3} = \frac{4^2 R_{k,2} - R_{k-1,2}}{4^2 - 1} = \frac{16 S_{k,1} - S_{k-1,1}}{15},$$

which is Richardson's extrapolation for Simpson's rule. The relation

$$S_{k,1} = \frac{16 S_{k,1} - S_{k-1,1}}{15} \quad (3.13)$$

is at the core of a well-known (and oftenly used) adaptive quadrature algorithm (due to Gander and Gautschi).

**Example 3.6.** Approximate the integral

$$I = \int_0^{\pi} \sin x \, dx$$

with precision  $\varepsilon = 10^{-1}$ , using Romberg's method, .

**Solution.** The exact value of the integral is

$$I = -\cos x \Big|_0^{\pi} = 2.$$

Using the repeated trapezoidal rule with  $n_1 = 2^0$ ,  $h_1 = \pi$  (i.e. nodes  $x_0 = 0, x_1 = \pi$ ) and  $n_2 = 2^1$ ,  $h_2 = \pi/2$  (so nodes  $x_0 = 0, x_1 = \pi/2, x_2 = \pi$ ), we get

$$\begin{aligned} R_{1,1} &= \frac{\pi}{2}(\sin 0 + \sin \pi) = 0, \\ R_{2,1} &= \frac{1}{2} \left[ R_{1,1} + h_1 f \left( a + \frac{1}{2} h_1 \right) \right] = \frac{1}{2} \left( 0 + \pi \sin \frac{\pi}{2} \right) = \frac{\pi}{2} = 1.5708. \end{aligned}$$

Richardson extrapolation is next:

$$R_{2,2} = \frac{4R_{2,1} - R_{1,1}}{3} = \frac{2\pi}{3} = 2.0944.$$

We have

$$|R_{2,2} - R_{1,1}| = 2.0944 > 0.1,$$

so we continue. We compute

$$\begin{aligned} R_{3,1} &= \frac{1}{2} \left[ R_{2,1} + h_2 \left[ f \left( a + \frac{1}{2} h_2 \right) + f \left( a + \frac{3}{2} h_2 \right) \right] \right] \\ &= \frac{1}{2} \left[ R_{2,1} + \frac{\pi}{2} \left( \sin \frac{\pi}{4} + \sin \frac{3\pi}{4} \right) \right] = 1.8961, \\ R_{3,2} &= \frac{4R_{3,1} - R_{2,1}}{3} = 2.0046, \\ R_{3,3} &= \frac{16R_{3,2} - R_{2,2}}{15} = 1.9986 \end{aligned}$$

and

$$|R_{3,3} - R_{2,2}| = 0.0958 < 0.1.$$

Hence, we obtained the approximation

$$I \approx R_{3,3} = 1.9986,$$

(with an error of  $1.4e - 3$ ), which is obviously better than the trapezoidal rule with  $n = 4$ ,  $R_{3,1}$  (with the error of 0.1039). Also, it is more accurate than Simpson's approximation with 4 nodes,  $I \approx 2.005$  (with error  $5e - 3$ ).

In fact, for this example, the algorithm converges very fast, as seen below:

$$\begin{array}{cccc} 0 & & & \\ 1.5708 & 2.0944 & & \\ 1.8961 & 2.0046 & 1.9986 & \\ 1.9742 & 2.0003 & 2.0000 & 2.0000 \end{array}$$

■

### 3.4 Weighted Gaussian Quadratures

The numerical methods studied so far were based on integrating linear and quadratic interpolating polynomials, and the resulting formulas were applied on subdivisions of ever smaller subintervals. In this section, we consider a numerical method that is based on the *exact integration* of polynomials of increasing degree; no subdivision of the integration interval is used. The motivation of this approach is the following: if we have a numerical integration formula to integrate low- to moderate-degree polynomials *exactly*, then the hope is that the same formula will integrate other functions  $f(x)$  *almost exactly*, if  $f(x)$  is well approximable by such polynomials.

#### 3.4.1 General Framework

**Definition 3.7.** *An interpolatory formula of the form*

$$\int_a^b w(x) f(x) dx = \sum_{k=1}^m A_k f(x_k) + R_m(f) \quad (3.14)$$



is called (**weighted**) **Gaussian quadrature** if it has maximum degree of precision,  $d = 2m - 1$ .

The function  $w : (a, b) \rightarrow \mathbb{R}_+$  is a *weight function*, a function for which the *moments*

$$\mu_j = \int_a^b w(x) x^j dx \quad (3.15)$$

exist and are finite for each  $j \in \mathbb{N}$ . The purpose of a weight function is to “absorb” some singularities of the integrand.

We want to determine the coefficients  $A_k$  and the nodes  $x_k$  such that

$$R_m(e_0) = R_m(e_1) = \dots = R_m(e_{2m-1}) = 0. \quad (3.16)$$

Let us start with a simple example. Consider the integral

$$\int_{-1}^1 f(x) dx. \quad (3.17)$$

So, in this case,  $w(x) \equiv 1$ .

**Case  $m = 1$ .**

We seek a numerical integration formula

$$\int_{-1}^1 f(x) dx \approx A_1 f(x_1).$$

From the first two relations in (3.16), we get

$$\begin{aligned} A_1 &= 2, \\ A_1 x_1 &= 0, \end{aligned}$$

and, thus, the formula

$$\int_{-1}^1 f(x) dx \approx 2f(0),$$

which is the midpoint (rectangle) rule. Recall that the rectangle rule had indeed the maximum

degree of precision possible with just one node,  $d = 1$ .

**Case  $m = 2$ .**

Now, we want a quadrature of the form

$$\int_{-1}^1 f(x) dx \approx A_1 f(x_1) + A_2 f(x_2),$$

with 4 unknowns, which are determined from the first 4 relations (3.16). This leads to the system

$$\begin{aligned} A_1 + A_2 &= 2 \\ A_1 x_1 + A_2 x_2 &= 0 \\ A_1 x_1^2 + A_2 x_2^2 &= \frac{2}{3} \\ A_1 x_1^3 + A_2 x_2^3 &= 0 \end{aligned} \tag{3.18}$$

with solution

$$A_1 = A_2 = 1, x_1 = -\frac{\sqrt{3}}{3}, x_2 = \frac{\sqrt{3}}{3}. \tag{3.19}$$

Hence, we found the quadrature formula

$$\int_{-1}^1 f(x) dx \approx f\left(-\frac{\sqrt{3}}{3}\right) + f\left(\frac{\sqrt{3}}{3}\right). \tag{3.20}$$

Being exact for the monomials  $1, x, x^2$  and  $x^3$ , this formula will be exact for *all* polynomials of degree  $\leq 3$ . Hence, its degree of exactness is  $d = 3$ . Compare this with Simpson's rule, which uses *three* nodes to attain the same degree of precision.

**General case  $m > 2$ .**

We seek now the formula

$$\int_{-1}^1 f(x) dx \approx \sum_{k=1}^m A_k f(x_k),$$

which has  $2m$  unspecified parameters, the nodes  $x_1, \dots, x_m$ , and the coefficients  $A_1, \dots, A_m$ . They are found by forcing the integration formula to be exact for the  $2m$  monomials  $1, x, x^2, \dots, x^{2m-1}$ . In turn, this forces the quadrature formula to be exact for *all* polynomials of degree  $2m - 1$ . This leads to the following system of  $2m$  nonlinear equations in  $2m$  unknowns:

$$\begin{aligned}
A_1 + A_2 + \cdots + A_{2m-1} &= 2 \\
A_1x_1 + A_2x_2 + \cdots + A_{2m-1}x_{2m-1} &= 0 \\
A_1x_1^2 + A_2x_2^2 + \cdots + A_{2m-1}x_{2m-1}^2 &= \frac{2}{3} \\
A_1x_1^3 + A_2x_2^3 + \cdots + A_{2m-1}x_{2m-1}^3 &= 0 \\
&\vdots \quad \vdots \\
A_1x_1^{2m-2} + A_2x_2^{2m-2} + \cdots + A_{2m-1}x_{2m-1}^{2m-2} &= \frac{2}{2m-1} \\
A_1x_1^{2m-1} + A_2x_2^{2m-1} + \cdots + A_{2m-1}x_{2m-1}^{2m-1} &= 0.
\end{aligned} \tag{3.21}$$

**Example 3.8.** Consider again the integral in Example [3.2](#).

$$I = \int_0^1 e^{-x^2} dx = 0.746824132812427.$$

**Solution.** The linear change of variables

$$x = \frac{b+a+t(b-a)}{2} \tag{3.22}$$

maps the interval  $[-1, 1]$  to  $[a, b]$ . So, with the substitution

$$x = \frac{1+t}{2}, \quad t = 2x-1,$$

we get

$$I = \frac{1}{2} \int_{-1}^1 e^{-\frac{1}{4}(1+t)^2} dt.$$

We apply Gaussian quadratures to the above integral. The errors are given in Table [2](#).

Comparing these with the ones given by the composite trapezoid or Simpson's rules (even with extrapolation), we see that these approximations are much more accurate, with fewer nodes.

■

$m$	Error
2	$2.29e - 4$
3	$9.55e - 6$
4	$3.35e - 7$
5	$6.05e - 9$
6	$7.77e - 11$
7	$7.89e - 13$

Table 2: Gaussian quadratures errors, Example 3.8

Gaussian quadrature formulas for a general weight function  $w$  can be found completely similarly. From relations (3.16), we obtain the system

$$\begin{aligned}
A_1 + A_2 + \cdots + A_{2m-1} &= \mu_0 \\
A_1 x_1 + A_2 x_2 + \cdots + A_{2m-1} x_{2m-1} &= \mu_1 \\
A_1 x_1^2 + A_2 x_2^2 + \cdots + A_{2m-1} x_{2m-1}^2 &= \mu_2 \\
&\vdots \\
A_1 x_1^{2m-1} + A_2 x_2^{2m-1} + \cdots + A_{2m-1} x_{2m-1}^{2m-1} &= \mu_{2m-1}.
\end{aligned} \tag{3.23}$$

**Example 3.9.** Find a Gaussian quadrature formula with 1 node, on the interval  $[0, 1]$ , with respect to the weight function  $w(x) = \frac{1}{\sqrt{x}}$ .

**Solution.** First off, let us notice that the function  $w(x) = \frac{1}{\sqrt{x}}$  is indeed a weight function on  $[0, 1]$ , since the moments

$$\mu_j = \int_0^1 \frac{1}{\sqrt{x}} x^j dx = \int_0^1 x^{j-\frac{1}{2}} dx = \frac{1}{j+\frac{1}{2}} x^{j+\frac{1}{2}} \Big|_0^1 = \frac{2}{2j+1}$$

exist and are finite for every  $j \in \mathbb{N}$ .

We want a formula of the form

$$\int_0^1 \frac{f(x)}{\sqrt{x}} dx \approx A_1 f(x_1),$$

having the maximum degree of precision possible, i.e.,  $d = 1$ .

Forcing equality for  $e_0(x) = 1$  and  $e_1(x) = x$  leads to the system

$$\begin{aligned} A_1 &= 2, \\ A_1 x_1 &= \frac{2}{3}, \end{aligned}$$

with solution  $A_1 = 2, x_1 = \frac{1}{3}$ . We obtain the formula

$$\int_0^1 \frac{f(x)}{\sqrt{x}} dx \approx 2f\left(\frac{1}{3}\right),$$

with degree of precision  $d = 1$ . ■

**Example 3.10.** Determine a Gaussian quadrature formula with 2 nodes, with respect to the weight function  $w(x) = e^{-x}$  on the interval  $[0, \infty)$ .

**Solution.** To compute the moments, recall *Euler's Gamma function*  $\Gamma : (0, \infty) \rightarrow (0, \infty)$ ,

$$\Gamma(a) = \int_0^\infty x^{a-1} e^{-x} dx, \quad \Gamma(n+1) = n!, \quad n \in \mathbb{N}.$$

Then the moments are

$$\mu_j = \int_0^\infty e^{-x} x^j dx = \int_0^\infty x^{(j+1)-1} e^{-x} dx = \Gamma(j+1) = j!, \quad j \in \mathbb{N}.$$

A quadrature formula with 2 nodes is of the form

$$\int_0^\infty e^{-x} f(x) dx \approx A_1 f(x_1) + A_2 f(x_2).$$

The nodes and coefficients are determined by having the formula above be exact for  $1, x, x^2$  and  $x^3$ , i.e. from the equations

$$\begin{aligned} A_1 + A_2 &= 1, \\ A_1 x_1 + A_2 x_2 &= 1, \end{aligned}$$

$$\begin{aligned} A_1 x_1^2 + A_2 x_2^2 &= 2, \\ A_1 x_1^3 + A_2 x_2^3 &= 6. \end{aligned}$$

The solution of the system above is

$$A_1 = \frac{2 + \sqrt{2}}{4}, A_2 = \frac{2 - \sqrt{2}}{4}, x_1 = 2 - \sqrt{2}, x_2 = 2 + \sqrt{2}$$

and yields the numerical integration formula

$$\int_0^\infty e^{-x} f(x) dx \approx \frac{2 + \sqrt{2}}{4} f(2 - \sqrt{2}) + \frac{2 - \sqrt{2}}{4} f(2 + \sqrt{2}),$$

with degree of exactness  $d = 3$ . ■

We see from these examples that solving the system (3.23) is *not* an easy task, even for a small number of nodes. This system is not linear in the nodes. It is linear in the coefficients, but with a Vandermonde system matrix, which is known to have *conditioning* (stability) problems. Even when a solution can be found (numerically), it is possible that some of the nodes are complex, or have values outside the interval  $[a, b]$ . Which is why we use another approach, one that involves *orthogonal polynomials*.

### 3.4.2 Orthogonal Polynomials

The use of orthogonal polynomials is justified by the following result.

**Theorem 3.11.** *Let  $u(x) = (x - x_1)(x - x_2) \dots (x - x_m)$ . Then, the quadrature formula (3.14) is exact for all polynomials  $p \in \mathbb{P}_{2m-1}$  if and only if  $u$  is orthogonal to the set  $\mathbb{P}_{m-1}$ ,  $u \perp \mathbb{P}_{m-1}$ , with respect to the inner product*

$$\langle f, g \rangle_w = \int_a^b w(x) f(x) g(x) dx. \quad (3.24)$$

*Proof.* “ $\Rightarrow$ ” Let  $p \in \mathbb{P}_{m-1}$ . Since  $u$  has degree  $m$ , it follows that  $up \in \mathbb{P}_{2m-1}$ , so formula (3.14) is

exact for  $up$ , i.e.,

$$\int_a^b w(x)u(x)p(x) dx = \sum_{k=1}^m A_k u(x_k)p(x_k) = 0,$$

because  $u(x_k) = 0, \forall k = \overline{1, m}$ . Hence,  $u \perp p$  and, further,  $u \perp \mathbb{P}_{m-1}$ .

“ $\Leftarrow$ ” Let  $f \in \mathbb{P}_{2m-1}$ , arbitrary. By the division algorithm, there exist  $q, r \in \mathbb{P}_{m-1}$  such that  $f = uq + r$ . Thus, we have

$$\int_a^b w(x)f(x) dx = \int_a^b w(x)u(x)q(x) dx + \int_a^b w(x)r(x) dx = 0 + \int_a^b w(x)r(x) dx,$$

since  $u \perp q$ .

Now, formula (3.14) is an interpolatory one, and as such, has degree of exactness at least  $d = m - 1$ . Since  $r \in \mathbb{P}_{m-1}$ , we have

$$\int_a^b w(x)r(x) dx = \sum_{k=1}^m A_k r(x_k).$$

But for any  $k = \overline{1, m}$ ,  $f(x_k) = u(x_k)q(x_k) + r(x_k) = r(x_k)$  and, thus,

$$\int_a^b w(x)f(x) dx = \sum_{k=1}^m A_k f(x_k),$$

i.e. formula (3.14) is exact for every  $f \in \mathbb{P}_{2m-1}$ . □

**Remark 3.12.** So we now know that the nodes of a Gaussian quadrature are the roots of a polynomial orthogonal to  $\mathbb{P}_{m-1}$  with respect to the weight  $w$ . Such families of orthogonal polynomials have been studied extensively. Table 3 contains such examples. A few immediate conclusions:

1. A first consequence is the fact that all the nodes in (3.14) are real, distinct and interior to the interval  $(a, b)$ .
2. Another consequence: the nodes can be obtained from the equation

$$u(x) = \begin{vmatrix} \mu_0 & \mu_1 & \dots & \mu_m \\ \mu_1 & \mu_2 & \dots & \mu_{m+1} \\ \vdots & & & \\ \mu_{m+1} & \mu_{m+2} & \dots & \mu_{2m-1} \\ 1 & x & \dots & x^m \end{vmatrix} = 0. \quad (3.25)$$

**3.** Recall that there exists a linear recurrence relation between 3 consecutive monic orthogonal polynomials on the interval  $[a, b]$  with respect to the weight  $w$ :

$$\pi_{k+1}(t) = (t - \alpha_k)\pi_k(t) - \beta_k\pi_{k-1}(t), \quad k = 0, 1, \dots, \quad \pi_{-1}(t) = 0, \quad \pi_0(t) = 1, \quad (3.26)$$

where

$$\alpha_k = \frac{\langle t\pi_k, \pi_k \rangle}{\|\pi_k\|^2}, \quad k = 0, 1, \dots, \quad \beta_k = \frac{\|\pi_k\|^2}{\|\pi_{k-1}\|^2}, \quad k = 1, 2, \dots, \quad \beta_0 = \mu_0. \quad (3.27)$$

Name	Notation	Polynomial	Weight fn.	Interval	$\alpha_k$	$\beta_k$
Legendre	$l_m$	$[(x^2 - 1)^m]^{(m)}$	1	$[-1, 1]$	0	$\beta_0 = 2,$ $\beta_k = (4 - k^2)^{-1}, k \geq 1$
Chebyshev 1 <sup>st</sup>	$T_m$	$\cos(m \arccos x)$	$(1 - x^2)^{-\frac{1}{2}}$	$[-1, 1]$	0	$\beta_0 = \pi,$ $\beta_1 = \frac{1}{2},$ $\beta_k = \frac{1}{4}, k \geq 2$
Chebyshev 2 <sup>nd</sup>	$Q_m$	$\frac{\sin[(m+1) \arccos x]}{\sqrt{1-x^2}}$	$(1 - x^2)^{\frac{1}{2}}$	$[-1, 1]$	0	$\beta_0 = \frac{\pi}{2},$ $\beta_k = \frac{1}{4}, k \geq 1$
Laguerre	$L_m^a$	$x^{-a}e^x (x^{m+a}e^{-x})^{(m)}$	$x^a e^{-x}, a > -1$	$[0, \infty)$	$2k + a + 1$	$\beta_0 = \Gamma(1 + a),$ $\beta_k = k(k + a), k \geq 1$
Hermite	$H_m$	$(-1)^m e^{x^2} (e^{-x^2})^{(m)}$	$e^{-x^2}$	$\mathbb{R}$	0	$\beta_0 = \sqrt{\pi},$ $\beta_k = \frac{k}{2}, k \geq 1$

Table 3: Orthogonal polynomials and recurrence coefficients

**Example 3.13.** Let us revisit some previous examples.

**Solution.** For the weight function  $w \equiv 1$ , we can now solve system (3.18) (a 2-point formula) much easier. We know that the nodes are the roots of the Legendre polynomial

$$l_2(x) = [(x^2 - 1)^2]'' = [x^4 - 2x^2 + 1]'' = [4x^3 - 4x]' = 4(3x^2 - 1),$$



i.e.  $\pm \frac{\sqrt{3}}{3}$ . Then the coefficients  $A_0 = A_1 = 1$  are immediately found.

In Example 3.10, the nodes are the roots of the Laguerre polynomial (with  $a = 0$ )

$$L_2^0(x) = e^x [x^2 e^{-x}]'' = e^x [(2x - x^2)e^{-x}]' = x^2 - 4x + 2,$$

so,  $2 \pm \sqrt{2}$ . Alternatively, by (3.25),

$$u(x) = \begin{vmatrix} \mu_0 & \mu_1 & \mu_2 \\ \mu_1 & \mu_2 & \mu_3 \\ 1 & x & x^2 \end{vmatrix} = \begin{vmatrix} 1 & 1 & 2 \\ 1 & 2 & 6 \\ 1 & x & x^2 \end{vmatrix} = (2x^2 - 6x) - (x^2 - 6) + 2(x - 2) = x^2 - 4x + 2.$$

Once the nodes are known, the coefficients can be easily found to be  $A_{1,2} = \frac{2 \mp \sqrt{2}}{4}$ . ■

Other properties of Gaussian quadratures:

**Proposition 3.14.** *The coefficients  $A_k, k = \overline{1, m}$  in (3.14) are all positive.*

Regarding the convergence of Gaussian quadratures, we have:

**Theorem 3.15.** *If  $[a, b]$  is bounded and  $f \in C[a, b]$ , then the Gaussian formula (3.14) converges,  $R_m(f) \rightarrow 0, m \rightarrow \infty$ .*

The proof is based on Weierstrass' theorem.

For the remainder of the quadrature formula, the following holds:

**Proposition 3.16.** *If  $f \in C^{2m}[a, b]$ , then there exists  $\xi \in (a, b)$  such that*

$$R_m(f) = \frac{f^{(2m)}(\xi)}{(2m)!} \int_a^b w(x) u^2(x) dx. \quad (3.28)$$

The proof is based on writing the Hermite interpolation polynomial at the double nodes  $x_1, \dots, x_m$ , multiplying it by the weight function and integrating.

**Example 3.17.** Find the error in the Gauss-Legendre and Gauss-Laguerre quadrature formulas with 2 nodes.

**Solution.** We have the 2-point numerical integration formula (3.20) (Gauss-Legendre):

$$\int_{-1}^1 f(x) dx = f\left(-\frac{\sqrt{3}}{3}\right) + f\left(\frac{\sqrt{3}}{3}\right) + R_2(f),$$

with  $u(x) = \left(x + \frac{\sqrt{3}}{3}\right)\left(x - \frac{\sqrt{3}}{3}\right) = x^2 - \frac{1}{3}$ . If  $f \in C^4[-1, 1]$ , for some  $\xi \in (-1, 1)$ , we have

$$\begin{aligned} R_2(f) &= \frac{f^{(iv)}(\xi)}{4!} \int_{-1}^1 u^2(x) dx = \frac{f^{(iv)}(\xi)}{4!} \int_{-1}^1 \left(x^4 - \frac{2}{3}x^2 + \frac{1}{9}\right) dx \\ &= 2 \frac{f^{(iv)}(\xi)}{24} \int_0^1 \left(x^4 - \frac{2}{3}x^2 + \frac{1}{9}\right) dx = \frac{f^{(iv)}(\xi)}{12} \left(\frac{1}{5}x^5 - \frac{2}{9}x^3 + \frac{1}{9}x\right) \Big|_0^1 = \frac{1}{135} f^{(iv)}(\xi). \end{aligned}$$

For the 2-point Gauss-Laguerre quadrature

$$\int_0^\infty e^{-x} f(x) dx = \frac{2+\sqrt{2}}{4} f(2-\sqrt{2}) + \frac{2-\sqrt{2}}{4} f(2+\sqrt{2}) + R_2(f),$$

again, assuming  $f \in C^4[0, \infty)$ , there exists  $\xi > 0$  such that the remainder is expressed as

$$\begin{aligned} R_2(f) &= \frac{f^{(iv)}(\xi)}{4!} \int_0^\infty e^{-x} u^2(x) dx = \frac{f^{(iv)}(\xi)}{24} \int_0^\infty e^{-x} \left((x-2)^2 - 2\right)^2 dx \\ &= \frac{f^{(iv)}(\xi)}{24} \int_0^\infty e^{-x} (x^2 - 4x + 2)^2 dx \\ &= \frac{f^{(iv)}(\xi)}{24} \int_0^\infty e^{-x} (x^4 - 8x^3 + 20x^2 - 16x + 4) dx \\ &= \frac{f^{(iv)}(\xi)}{24} \left(\Gamma(5) - 8\Gamma(4) + 20\Gamma(3) - 16\Gamma(2) + 4\Gamma(1)\right) \\ &= \frac{f^{(iv)}(\xi)}{24} (4! - 8 \cdot 3! + 20 \cdot 2! - 16 \cdot 1! + 4 \cdot 0!) = \frac{1}{6} f^{(iv)}(\xi). \end{aligned}$$

■

Now we have better procedures for finding the nodes of a Gaussian quadrature formula (orthog-

onal polynomials are implemented in most mathematical software, for Matlab, see <https://www.mathworks.com/matlabcentral/fileexchange/69956-orthogonalpolynomials>). But system (3.23) is still a Vandermonde system in the coefficients. So, for an efficient implementation, we can still improve computations. For that, we will make use of the recurrence relation and the parameters in (3.26)–(3.27). With these, we define

$$J_m(w) = \begin{bmatrix} \alpha_0 & \sqrt{\beta_1} & & & 0 \\ \sqrt{\beta_1} & \alpha_1 & \sqrt{\beta_2} & & \\ & \sqrt{\beta_2} & \alpha_2 & \ddots & \\ & & \ddots & \ddots & \sqrt{\beta_{m-1}} \\ 0 & & & \sqrt{\beta_{m-1}} & \alpha_{m-1} \end{bmatrix}, \quad (3.29)$$

called the **Jacobi matrix** of order  $m$  for the weight function  $w$  on the interval  $[a, b]$ . The following holds:

**Theorem 3.18.** *The nodes  $\{x_k\}_{k=1}^m$  of the Gaussian formula (3.14) are the eigenvalues of  $J_m$ ,*

$$J_m v_k = x_k v_k, \quad v_k^T v_k = 1, \quad k = 1, \dots, m, \quad (3.30)$$

while the coefficients  $\{A_k\}_{k=1}^m$  are given by

$$A_k = \beta_0 v_{k,1}^2, \quad k = 1, \dots, m, \quad (3.31)$$

where  $v_{k,1}$  is the first component of the normalized ( $\|v_k\| = 1$ ) eigenvector associated with the eigenvalue  $x_k$ .

This is easily proved by writing the recurrence relation (3.26) in matrix (vector) form.

**Remark 3.19.** Thus, the problem of determining a Gauss numerical integration formula is now reduced to that of finding e-values and e-vectors for a *symmetric* and *tridiagonal* matrix. This problem has been studied extensively in linear algebra, there is a vast literature on it and there are many very efficient methods for solving it.

# Chapter 5. Numerical Solution of Nonlinear Equations

## 1 Introduction to Iterative Methods

Finding one or more roots of an equation

$$f(x) = 0, \quad (1.1)$$

is one of the most commonly occurring problems of Applied Mathematics. Even the simplest of nonlinear equations – e.g., algebraic equations – are known to not admit solutions that are expressible rationally in terms of the data. It is therefore impossible, in general, to compute roots of nonlinear equations in a finite numbers of arithmetic operations.

The function  $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$  is a nonlinear function, which will be assumed to have a certain degree of smoothness. If  $n > 1$ , then (1.1) represents a system of  $n$  equations (at least one nonlinear) with  $m$  unknowns.

For now, we will restrict our discussion to the case  $m = n = 1$ , although many of the procedures we describe can easily be generalized to the multidimensional case.

**Definition 1.1.** A number  $\alpha \in \mathbb{C}$  satisfying equation (1.1) is called a **zero** or a **root** of  $f$ .

As mentioned before, in most cases, explicit solutions of equation (1.1) are not available and we must try to find a root to any specified degree of accuracy. The numerical methods for finding the roots will be *iterative methods* and will require the knowledge of one (or more) initial value(s)  $x_0 (x_1, \dots)$ . Then the method will produce a sequence  $\{x_n\}_{n \in \mathbb{N}}$  of approximations of  $\alpha$ , such that  $\lim_{n \rightarrow \infty} x_n = \alpha$ . These initial values will be determined, in general, from the context of the problem or from the graph of the function.

The analysis of an iterative method will include

- the proof of convergence,  $x_n \rightarrow \alpha$ , as  $n \rightarrow \infty$ ;
- finding the *interval of convergence*, i.e. the set of values of the initial guess(es)  $x_0 (x_1, \dots)$  for which the method converges;
- determining the *speed* of convergence.

What makes an iterative method better than another is *how fast* it converges to the desired solution. Regarding the speed of convergence, we define the following:

**Definition 1.2.** We say that a sequence of iterates  $\{x_n\}_{n \in \mathbb{N}}$  converges to  $\alpha$  with **order of convergence**  $p \geq 1$ , if

$$|x_{n+1} - \alpha| \leq c |x_n - \alpha|^p, \text{ for all } n \in \mathbb{N}, \quad (1.2)$$

where  $c > 0$  is a constant independent of  $n$ .

If  $p = 1$ , the method is said to **converge linearly** to  $\alpha$ , in which case we also require that  $c < 1$ . Then the constant  $c$  is called the **rate of linear convergence** of  $x_n$  to  $\alpha$ .

For  $1 < p < 2$ , we say that the convergence is **superlinear**.

**Remark 1.3.** If  $p = 1$ , then

$$|x_n - \alpha| \leq c |x_{n-1} - \alpha| \leq \dots \leq c^n |x_0 - \alpha|,$$

which is why we require that  $c < 1$ .

## 2 Common Rootfinding Methods

We start with three simple methods and then give a general theory for one-point iteration methods. We recall some known results from Analysis, that will be used in the sequel.

**Theorem 2.1. [Intermediate Value Theorem]**

If  $f : [a, b] \rightarrow \mathbb{R}$  is a continuous function, then it takes on any given value between  $f(a)$  and  $f(b)$  at some point within the interval. As a consequence, if a continuous function has values of opposite sign inside an interval  $[a, b]$ , then it has at least one root in that interval.

**Theorem 2.2. [Rolle's Theorem]**

If a function  $f$  is continuous on  $[a, b]$  and differentiable on  $(a, b)$ , with  $f(a) = f(b)$ , then there exists a point  $c \in (a, b)$  such that  $f'(c) = 0$ . As a consequence, between any two distinct real roots of  $f$ , there is a root of the derivative.

So, combining the two, we can find the number of real zeros of a function (satisfying the conditions above) and locate them, by counting the number of *sign changes* of the function at the roots of the derivative and endpoints of the domain of definition.

### 2.1 Bisection Method

Assume that  $f : \mathbb{R} \rightarrow \mathbb{R}$  is continuous on an interval  $[a, b] \subset \mathbb{R}$  and that

$$f(a)f(b) < 0. \quad (2.1)$$

Then, by the Intermediate Value Theorem, there exists  $\alpha \in (a, b)$  such that  $f(\alpha) = 0$ .

The simplest numerical procedure for finding a root is to repeatedly *halve (bisection)* the interval  $[a, b]$ , keeping the half on which  $f(x)$  changes sign. This procedure is called the **bisection method**. Denoting by  $[a_1, b_1] = [a, b]$ , the method will produce a sequence of embedded intervals  $[a_n, b_n]$ , such that for every  $n \in \mathbb{N}$ ,  $\alpha \in [a_n, b_n]$ ,  $f(a_n)f(b_n) < 0$ , and a sequence of approximations

$$c_n = \frac{a_n + b_n}{2} \quad (2.2)$$

of the root  $\alpha$  (see Figure 1).

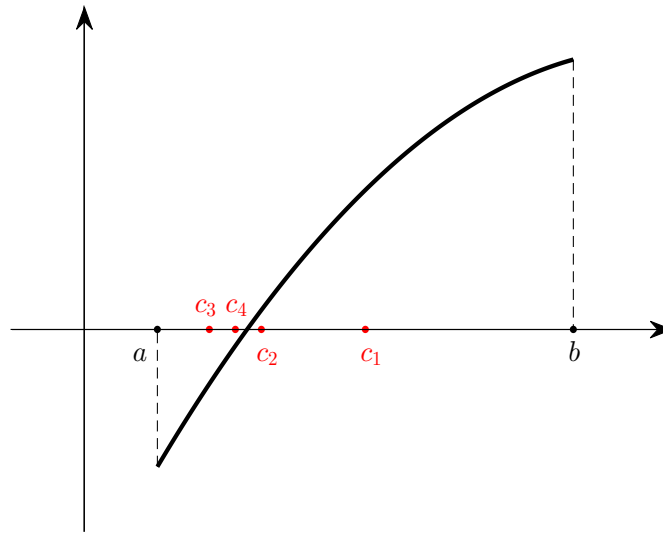


Fig. 1: Bisection method

Usually  $[a, b]$  is chosen to contain only one root  $\alpha$ , but the following algorithm for the bisection method will always converge to some root  $\alpha \in [a, b]$ , because of (2.1).

**Algorithm 2.3.** [Bisection method]

function  $\alpha = \text{Bisect}(f, a, b, \varepsilon)$

1. Define  $c = (a + b)/2$ .
2. If  $b - c \leq \varepsilon$ , then  $\alpha = c$  and exit.
3. If  $\text{sign}(f(b)) \cdot \text{sign}(f(c)) \leq 0$ , then  $a = c$ ; otherwise,  $b = c$ .
4. Return to step 1.

The sequence  $\{a_n\}_{n \in \mathbb{N}}$  is monotonely increasing, sequence  $\{b_n\}_{n \in \mathbb{N}}$  is monotonely decreasing and

$$\lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} b_n = \lim_{n \rightarrow \infty} c_n = \alpha.$$

Also, we have

$$\begin{aligned} |x_n - \alpha| &\leq b_n - a_n = \frac{b - a}{2^n}, \\ |x_{n+1} - \alpha| &\leq \frac{1}{2} |x_n - \alpha|, \end{aligned} \tag{2.3}$$

which shows that the bisection method converges *linearly* (order of convergence  $p = 1$ ) with a rate of convergence of  $\frac{1}{2}$ .

**Example 2.4.** Find the largest root of

$$f(x) \equiv x^6 - x - 1 = 0, \tag{2.4}$$

with an error of  $\varepsilon = 0.001$ .

**Solution.** First, let us see how many real roots are there and where they are (approximately) located. We have

$$\begin{aligned} f(x) &= x^6 - x - 1, \\ f'(x) &= 6x^5 - 1. \end{aligned}$$

The derivative  $f'$  has only one real root, namely  $\frac{1}{\sqrt[5]{6}}$ . Now,

$$f\left(\frac{1}{\sqrt[5]{6}}\right) = \frac{1}{6} \cdot \frac{1}{\sqrt[5]{6}} - \frac{1}{\sqrt[5]{6}} - 1 = -\frac{5}{6} \cdot \frac{1}{\sqrt[5]{6}} - 1 < 0,$$

so the table of variation of  $f$  is

$x$	$-\infty$	$\frac{1}{\sqrt[5]{6}}$	$\infty$
$f$	$+$	$-$	$+$

Thus,  $f$  has two real roots, one in  $\left(-\infty, \frac{1}{\sqrt[5]{6}}\right)$  and one,  $\alpha \in \left(\frac{1}{\sqrt[5]{6}}, \infty\right)$  (which we want to approximate).

In fact, since

$$\begin{aligned} f(-1) &= 1, f(0) = -1 \text{ and} \\ f(1) &= -1, f(2) = 61, \end{aligned}$$

we have a more precise location: a negative root between  $(-1, 0)$  and the positive root that we seek,  $\alpha \in (1, 2)$ . That also gives us the starting interval for the bisection method. Alternatively, we can see from the graph the approximate location of the two real roots (see Figure 2).

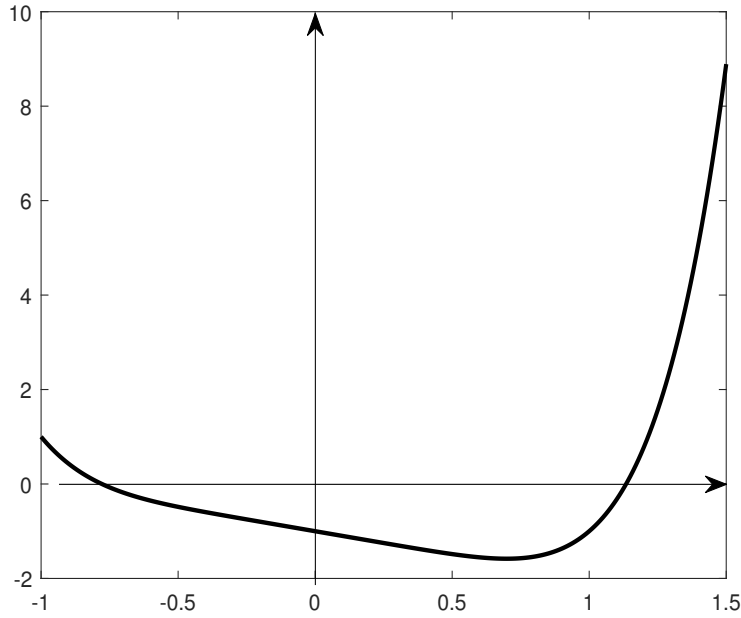


Fig. 2: Function  $f(x) = x^6 - x - 1$

So, we start with the interval  $[a_1, b_1] = [1, 2]$ . How many iterations are needed for precision  $\varepsilon = 0.001$ ? We find  $n$  from (2.3):

$$\begin{aligned} \frac{b-a}{2^n} &\leq \varepsilon, \text{ which means} \\ n &\geq \log_2 \left( \frac{b-a}{\varepsilon} \right), \text{ i.e., in our example,} \\ n &\geq \log_2 \left( \frac{1}{10^{-3}} \right) = 9.9658. \end{aligned}$$

The results of the bisection method are shown in Table 1. Indeed, after  $n = 10$  iterations, we obtain the desired precision. ■



$n$	$a$	$b$	$c$	$b - c$	$f(c)$
1	1.0000	2.0000	1.5000	0.5000	8.8906
2	1.0000	1.5000	1.2500	0.2500	1.5647
3	1.0000	1.2500	1.1250	0.1250	-0.0977
4	1.1250	1.2500	1.1875	0.0625	0.6167
5	1.1250	1.1875	1.1562	0.0312	0.2333
6	1.1250	1.1562	1.1406	0.0156	0.0616
7	1.1250	1.1406	1.1328	0.0078	-0.0196
8	1.1328	1.1406	1.1367	0.0039	0.0206
9	1.1328	1.1367	1.1348	0.0020	0.0004
10	1.1328	1.1348	1.1338	0.00098	-0.0096

Table 1: Bisection Method for  $x^6 - x - 1 = 0$

**Remark 2.5.** The bisection method is a *two-point method*, since two approximate values are needed to obtain an improved value. There are several advantages to the bisection method. The principal one is that the method is guaranteed to converge, as long as the function  $f$  is continuous and (2.1) is satisfied. In addition, the error bound given in (2.3) is guaranteed to decrease by one half with each iteration. This relation can also be used as a stopping criterion, as was done in the previous example. The principal disadvantage of the bisection method is that it generally converges slowly (only linearly), more slowly than most other methods. Also, it only approximates real roots.

The next two methods follow the same idea: approximate  $f$  by a linear interpolation polynomial and find the root of that polynomial. In other words, the graph of  $y = f(x)$  is approximated by a straight line and the  $x$ -intercept of that line is approximating the root of  $f$ .

## 2.2 Secant Method

Assume that two initial guesses to  $\alpha$  are known and denote them by  $x_0$  and  $x_1$ . We approximate  $f$  by its Lagrange polynomial at the nodes  $x_0$  and  $x_1$ . So the graph of  $y = f(x)$  is approximated by the *secant* line determined by the points  $(x_0, f(x_0))$  and  $(x_1, f(x_1))$ . The root  $\alpha$  of  $f$  is then approximated by  $x_2$ , the  $x$ -intercept of the secant line. We hope  $x_2$  will be an improved approximation of  $\alpha$ . This is illustrated in Figure 3.

Let us find the value of  $x_2$ . The equation of the secant line is

$$y - f(x_1) = \frac{f(x_1) - f(x_0)}{x_1 - x_0}(x - x_1).$$

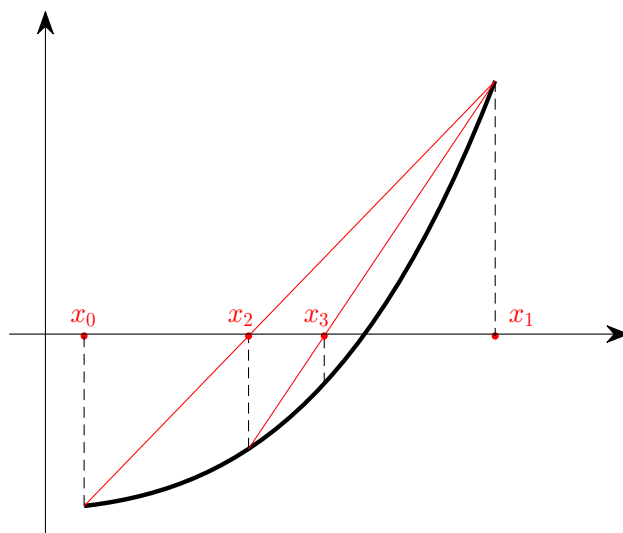


Fig. 3: Secant method

We find its point of intersection with the  $x$ -axis by letting  $y = 0$  and solving for  $x$ . We get

$$x_2 = x_1 - f(x_1) \frac{x_1 - x_0}{f(x_1) - f(x_0)}$$

Having found  $x_2$ , we use  $x_1$  and  $x_2$  as a new set of approximate values for  $\alpha$ . This leads to an improved value  $x_3$ . Recursively, we obtain a sequence of iterates given by

$$x_{n+1} = x_n - f(x_n) \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})}, \quad n = 1, 2, \dots, \quad (2.5)$$

called the **secant method**.

**Example 2.6.** We solve again the equation

$$f(x) \equiv x^6 - x - 1 = 0,$$

which was used previously as an example for the bisection method.

**Solution.** We start with

$$x_0 = 1, \quad x_1 = 2.$$

The results are given in Table 2, including the quantities  $x_n - x_{n-1}$  as an estimate of  $\alpha - x_{n-1}$ . The

iterate  $x_8$  equals  $\alpha$  rounded to nine significant digits.

$n$	$x_n$	$f(x_n)$	$x_n - x_{n-1}$	$\alpha - x_{n-1}$
0	2.0	61.0		
1	1.0	-1.0	-1.0	
2	1.01612903	$-9.15e-1$	$1.61e-2$	$1.35e-1$
3	1.19057777	$6.57e-1$	$1.74e-1$	$1.19e-1$
4	1.11765583	$-1.68e-1$	$-7.29e-2$	$-5.59e-2$
5	1.13253155	$-2.24e-2$	$1.49e-2$	$1.71e-2$
6	1.13481681	$9.54e-4$	$2.29e-3$	$2.19e-3$
7	1.13472365	$-5.07e-6$	$-9.32e-5$	$-9.27e-5$
8	1.13472414	$-1.13e-9$	$4.92e-7$	$4.92e-7$

Table 2: Secant Method for  $x^6 - x - 1 = 0$

■

The secant method is also a two-point iterative method. Unlike the bisection method, it *does not* always converge. For a convergence and error analysis, let us compute, from (2.5),

$$\begin{aligned}
x_{n+1} - \alpha &= x_n - \alpha - f(x_n) \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} \\
&= x_n - \alpha - \frac{f(x_n)}{\frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}} \\
&= x_n - \alpha - \frac{f(x_n) - f(\alpha)}{\frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}}, \text{ since } f(\alpha) = 0.
\end{aligned}$$

Further, we make use of divided differences and obtain

$$\begin{aligned}
x_{n+1} - \alpha &= x_n - \alpha - (x_n - \alpha) \frac{f[x_n, \alpha]}{f[x_{n-1}, x_n]} \\
&= (x_n - \alpha) \left[ 1 - \frac{f[x_n, \alpha]}{f[x_{n-1}, x_n]} \right] \\
&= (x_n - \alpha) \frac{f[x_{n-1}, x_n] - f[x_n, \alpha]}{f[x_{n-1}, x_n]} \\
&= (x_n - \alpha)(x_{n-1} - \alpha) \frac{f[x_{n-1}, x_n, \alpha]}{f[x_{n-1}, x_n]},
\end{aligned}$$

so, assuming  $f$  is smooth enough,

$$x_{n+1} - \alpha = (x_n - \alpha)(x_{n-1} - \alpha) \frac{f''(\xi_n)}{2f'(\zeta_n)}, \quad (2.6)$$

with  $\zeta_n$  between  $x_n$  and  $x_{n-1}$ , and  $\xi_n$  between the smallest and the largest of the numbers  $\alpha, x_n$  and  $x_{n-1}$ . Using (2.6) and a limiting argument, we have the following convergence result.

**Theorem 2.7.** Assume  $f, f'$  and  $f''$  are continuous on an interval  $I_\varepsilon = (\alpha - \varepsilon, \alpha + \varepsilon)$  containing the simple root  $\alpha$  ( $f'(\alpha) \neq 0$ ). Then, for starting values  $x_0$  and  $x_1$  sufficiently close to  $\alpha$ , the iterates in (2.5) converge to  $\alpha$ , with order of convergence

$$p = r = \frac{1 + \sqrt{5}}{2} \approx 1.618033 \dots, \quad (2.7)$$

known as the **golden ratio**.

Thus, the secant method converges *superlinearly*.

**Remark 2.8.**

1. To understand what “sufficiently close” means in the theorem above, let

$$M_\varepsilon = \frac{\max_{I_\varepsilon} |f''(x)|}{2 \min_{I_\varepsilon} |f'(x)|}, \quad e_0 = |x_0 - \alpha|, \quad e_1 = |x_1 - \alpha|. \quad (2.8)$$

Then the method above will converge if  $x_0, x_1 \in I_\varepsilon$ , with  $\varepsilon > 0$  chosen so that

$$\max\{M_\varepsilon e_0, M_\varepsilon e_1\} < 1. \quad (2.9)$$

2. It is clear now that the secant method does not always converge, but when it does, it does so *faster* than the bisection method (its order of convergence is higher). That was obvious in our example.

3. Another advantage is that the secant method can be used to approximate complex roots, as well, if the initial values  $x_0$  and  $x_1$  are taken to be complex numbers satisfying the conditions above.

4. It can be shown that

$$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - x_n|}{|x_n - \alpha|} = 1 \text{ and, thus,} \quad (2.10)$$

$$|x_n - \alpha| \approx |x_{n+1} - x_n|, \text{ for sufficiently large } n,$$

which can be used as a stopping criterion.

## 2.3 Newton's Method

In a similar fashion, now we start with one initial value,  $x_0$  and approximate  $f$  by its linear Taylor polynomial at the double node  $x_0$ . In other words, the graph of  $y = f(x)$  is approximated by the line *tangent* to the graph of  $f$  at the point  $(x_0, f(x_0))$ . The root  $\alpha$  of  $f$  is then approximated by  $x_1$ , the point of intersection of the tangent line with the  $x$ -axis. If  $x_0$  is close enough to  $\alpha$ , then the root of the Taylor polynomial should be close to  $\alpha$ .

The tangent line at  $x_0$  has equation

$$y - f(x_0) = f'(x_0)(x - x_0),$$

so, for  $x_1$ , we find

$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}.$$

Repeat the process to further improve the estimate of  $\alpha$ . Recursively, we get

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}, \quad n = 0, 1, \dots \quad (2.11)$$

This is called **Newton's (tangent) method** and it is illustrated in Figure 4.

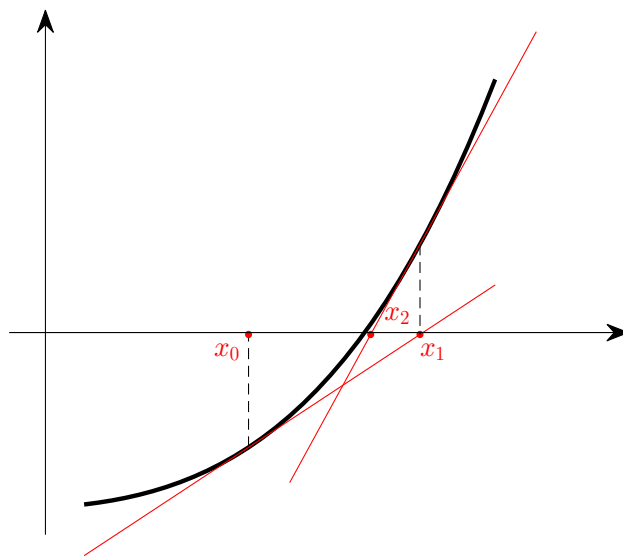


Fig. 4: Newton's method

**Example 2.9.** Let us approximate the positive solution of

$$f(x) \equiv x^6 - x - 1 = 0,$$

using Newton's method.

**Solution.** An initial guess  $x_0$  can be taken from the graph of  $y = f(x)$  in Figure 2. The iterative method is given by

$$x_{n+1} = x_n - \frac{x_n^6 - x_n - 1}{6x_n^5 - 1}, \quad n \geq 0.$$

Table 3 shows the results of Newton's method with initial value  $x_0 = 1.5$ .

$n$	$x_n$	$f(x_n)$	$x_n - x_{n-1}$	$\alpha - x_{n-1}$
0	1.5	$8.89e + 1$		
1	1.30049088	$2.54e + 1$	$-2.00e - 1$	$-3.65e - 1$
2	1.18148042	$5.38e - 1$	$-1.19e - 1$	$-1.66e - 1$
3	1.13945559	$4.92e - 2$	$-4.20e - 2$	$-4.68e - 2$
4	1.13477763	$5.50e - 4$	$-4.68e - 3$	$-4.73e - 3$
5	1.13472415	$7.11e - 8$	$-5.35e - 5$	$-5.35e - 5$
6	1.13472414	$1.55e - 15$	$-6.91e - 9$	$-6.91e - 9$

Table 3: Newton's Method for  $x^6 - x - 1 = 0$

As seen from the table, the convergence is very rapid. The iterate  $x_6$  is accurate (almost) to the machine precision of around 16 decimal digits. ■

As before, we can compute

$$\begin{aligned} x_{n+1} - \alpha &= (x_n - \alpha)^2 \frac{f[x_n, x_n, \alpha]}{f[x_n, x_n]} \\ &= (x_n - \alpha)^2 \frac{f''(\xi_n)}{2f'(x_n)}, \end{aligned} \tag{2.12}$$

with  $\xi_n$  between  $\alpha$  and  $x_n$ . Then we have the following convergence result.

**Theorem 2.10.** Assume  $f$ ,  $f'$  and  $f''$  are continuous on an interval  $I_\varepsilon = (\alpha - \varepsilon, \alpha + \varepsilon)$  containing the simple root  $\alpha$  ( $f'(\alpha) \neq 0$ ). Then, if the initial value  $x_0$  is sufficiently close to  $\alpha$ , the iterates in

(2.11) converge to  $\alpha$  and

$$\lim_{n \rightarrow \infty} \frac{x_{n+1} - \alpha}{(x_n - \alpha)^2} = \frac{f''(\alpha)}{2f'(\alpha)}, \quad (2.13)$$

which shows that the order of convergence of Newton's method is  $p = 2$ .

**Remark 2.11.**

1. Similarly with Remark 2.8, “sufficiently close” means  $x_0 \in I_\varepsilon$ , where  $\varepsilon$  is chosen so that  $M_\varepsilon e_0 < 1$ , with  $M_\varepsilon$  and  $e_0$  defined in (2.8).
2. Again, as before, Newton's method *does not* always converge, but when it does, it does so faster ( $p = 2$ ) than the bisection method ( $p = 1$ ) and the secant method ( $p = (1 + \sqrt{5})/2 \approx 1.618$ ).
3. Also, Newton's method can be used to approximate complex roots, as well, if the initial value  $x_0$  is a complex number satisfying the conditions above.
4. Again, for sufficiently large  $n$ ,

$$|x_n - \alpha| \approx |x_{n+1} - x_n|,$$

which can be used as a stopping criterion.

5. Unlike the bisection and secant methods, Newton's method is a *one-step* iterative method, as it only requires one initial value. Later on, we will give a more comprehensive analysis of one-step iterative methods.

## 2.4 Comparison Between Newton's and Secant Methods

As we have seen, Newton's method and the secant method are closely related. If the approximation

$$f'(x_n) \approx \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}$$

is used in Newton's formula (2.11), we obtain the secant formula (2.5).

The conditions for convergence are almost identical and the error formulas are similar. Nonetheless, there are two major differences. Newton's method requires two function evaluations per iterate, those of  $f(x_n)$  and  $f'(x_n)$ , whereas the secant method requires only one function evaluation per iterate, that of  $f(x_n)$  (provided that the value of  $f(x_{n-1})$  is retained from the last iteration). So, Newton's method is generally more expensive per iteration. On the other hand, it converges more rapidly (order  $p = 2$  versus  $p = r \approx 1.62$ ) and consequently, it will require fewer iterations to attain a given desired accuracy. A comparison of the expenditure of computational time needed to approximate a root  $\alpha$  within a desired tolerance, can be made.

To simplify the analysis, we assume that the initial guesses are quite close to the desired root, so both methods converge. Let  $t$  be the time needed to evaluate  $f(x)$ , and  $s \cdot t$  the time required to evaluate  $f'(x)$ . By writing the operations involved in the two methods, it can be then shown that the minimum time to obtain the desired accuracy with Newton's method is

$$T_N = \frac{(1+s)tK}{\log 2},$$

while, for the secant method, a similar calculation shows that the minimum time necessary to obtain the desired accuracy is

$$T_S = \frac{tK}{\log r},$$

where  $K$  is a positive constant that depends on  $\varepsilon$ ,  $x_0$  and  $c = \left| \frac{f''(\alpha)}{2f'(\alpha)} \right|$ . Thus,

$$\frac{T_S}{T_N} = \frac{\log 2}{(1+s) \log r}.$$

The secant method is faster than Newton's method if the ratio is less than one, i.e.

$$\begin{aligned} \frac{T_S}{T_N} &< 1 \\ s &> \frac{\log 2}{\log r} - 1 \approx 0.44. \end{aligned}$$

In conclusion, if the time needed to evaluate  $f'(x)$  is more than 44% of that necessary to evaluate  $f(x)$ , then the secant method is more efficient. In practice, many other factors will affect the relative costs of the two methods, so that the .44 factor should be used with caution.

### 3 One-Point Iteration Methods – General Theory

#### 3.1 Fixed Point Iteration

A classical approach is to reformulate equation  $f(x) = 0$  as

$$x = g(x) \tag{3.1}$$



and find a *fixed point* for  $g$ . Let us first note that the form (3.1) is *not* restrictive in any way. In fact, any equation can be written in the form (3.1) in a multitude of ways.

**Example 3.1.** Consider the equation

$$x^2 - 3 = 0.$$

It can be rewritten, for instance, as

$$\begin{aligned} \text{(a)} \quad x &= x^2 + x - 3, \\ \text{(b)} \quad x &= \frac{3}{x}, \\ \text{(c)} \quad x &= \frac{1}{2} \left( x + \frac{3}{x} \right), \\ \text{(d)} \quad x &= x + c(x^2 - 3), \text{ for some constant } c \in \mathbb{R}, \end{aligned}$$

and many other ways.

Now we can employ fixed point theory to discuss the solvability of equation (3.1). In what follows, the notation

$$g([a, b]) \subseteq [a, b]$$

means

$$x \in [a, b] \implies g(x) \in [a, b].$$

**Lemma 3.2.** *Let  $g \in C[a, b]$ , such that  $g([a, b]) \subseteq [a, b]$ . Then  $g$  has at least one fixed point in  $[a, b]$ .*

*Proof.* This follows immediately from the Intermediate Value Theorem applied to the function

$$G(x) = g(x) - x.$$

Since  $G$  is continuous and  $G(a) \geq 0$ ,  $G(b) \leq 0$ ,  $G$  must have at least one zero in  $[a, b]$ , which is, obviously, a fixed point of  $g$ . □

**Theorem 3.3. [Banach]** *Let  $g \in C[a, b]$ , with  $g([a, b]) \subseteq [a, b]$ . Assume that there exists  $0 < \lambda < 1$  such that*

$$|g(x) - g(y)| \leq \lambda |x - y|, \quad \forall x, y \in [a, b] \quad (\text{i.e., } g \text{ is a } \mathbf{contraction}). \quad (3.2)$$

Then  $g$  has a unique fixed point  $\alpha \in [a, b]$ . Furthermore, the iterates

$$x_{n+1} = g(x_n), \quad n \geq 0, \quad (3.3)$$

converge to  $\alpha$ , for any choice of  $x_0 \in [a, b]$  and the following error estimates hold:

$$\begin{aligned} |x_n - \alpha| &\leq \lambda |x_{n-1} - \alpha|, \quad n \geq 1, \\ |x_n - \alpha| &\leq \frac{\lambda^n}{1 - \lambda} |x_1 - x_0|. \end{aligned} \quad (3.4)$$

*Proof.* The existence of the fixed point is guaranteed by Lemma 3.2.

To prove its uniqueness, assume there are two fixed points,  $\alpha = g(\alpha)$ ,  $\beta = g(\beta)$ ,  $\alpha \neq \beta$ . Then

$$|\alpha - \beta| = |g(\alpha) - g(\beta)| \stackrel{(3.2)}{\leq} \lambda |\alpha - \beta|$$

and, so,

$$(1 - \lambda)|\alpha - \beta| \leq 0,$$

which is a contradiction. Thus,  $\alpha = \beta$ .

To prove the convergence, let us note that

$$x_0 \in [a, b] \implies x_1 = g(x_0) \in [a, b] \implies \dots \implies x_n \in [a, b], \quad \forall n \geq 0.$$

Then,

$$\begin{aligned} |x_n - \alpha| &= |g(x_{n-1}) - g(\alpha)| \leq \lambda |x_{n-1} - \alpha| = \lambda |g(x_{n-2}) - g(\alpha)| \\ &\leq \lambda^2 |x_{n-2} - \alpha| \leq \dots \leq \lambda^n |x_0 - \alpha|. \end{aligned}$$

Letting  $n \rightarrow \infty$ , since  $\lambda^n \rightarrow 0$ , it follows that  $x_n \rightarrow \alpha$  and the first bound in (3.4) holds.

For the second bound, we write

$$\begin{aligned} |x_0 - \alpha| &\leq |x_0 - x_1| + |x_1 - \alpha| \leq |x_0 - x_1| + \lambda |x_0 - \alpha|, \\ |x_0 - \alpha| &\leq \frac{1}{1 - \lambda} |x_1 - x_0|. \end{aligned}$$

Combining this with the previous relation, we get the second bound in (3.4).  $\square$

#### **Remark 3.4.**

**1.** The first bound shows that  $\{x_n\}_{n \in \mathbb{N}}$  converges *linearly*, with a rate of convergence bounded by

the contraction constant  $\lambda$ .

2. From the proof of Theorem 3.3, we can also show that

$$\begin{aligned} |x_n - \alpha| &\leq \frac{1}{1 - \lambda} |x_{n+1} - x_n| \text{ and, hence,} \\ |x_{n+1} - \alpha| &\leq \lambda |x_n - \alpha| \leq \frac{\lambda}{1 - \lambda} |x_{n+1} - x_n|, \end{aligned}$$

which gives a stopping criterion

$$|x_{n+1} - x_n| \leq \frac{1 - \lambda}{\lambda} \varepsilon. \quad (3.5)$$

3. If  $g$  is also differentiable on  $(a, b)$ , then, by the MVT, there exists  $c \in (a, b)$  such that

$$g(x) - g(y) = g'(c)(x - y), \forall x, y \in [a, b].$$

Letting  $\lambda = \max_{x \in [a, b]} |g'(x)|$ , it follows that

$$|g(x) - g(y)| \leq \lambda |x - y|, \forall x, y \in [a, b].$$

Then, we can restate the convergence result.

**Theorem 3.5.** *Let  $g \in C^1[a, b]$ , such that  $g([a, b]) \subseteq [a, b]$  and*

$$\lambda := \max_{x \in [a, b]} |g'(x)| < 1. \quad (3.6)$$

*Then:*

*a) Function  $g$  has a unique fixed point  $\alpha \in [a, b]$ .*

*b) For any initial choice  $x_0 \in [a, b]$ , the sequence  $x_{n+1} = g(x_n)$  converges to  $\alpha$ , as  $n \rightarrow \infty$ .*

*c)  $|x_n - \alpha| \leq \lambda^n |x_0 - \alpha| \leq \frac{\lambda^n}{1 - \lambda} |x_1 - x_0|$ ,  $n \geq 1$ .*

*d)*

$$\lim_{n \rightarrow \infty} \frac{x_{n+1} - \alpha}{x_n - \alpha} = g'(\alpha), \quad (3.7)$$

*so, if  $g'(\alpha) \neq 0$ , the iterative method  $x_{n+1} = g(x_n)$  is linearly convergent to the root  $\alpha$  with rate of convergence bounded by  $\lambda$ .*

The conditions of Theorem 3.5 can be relaxed and imposed only *locally*, near the root  $\alpha$ .

**Theorem 3.6.** Assume  $\alpha$  is a fixed point of  $g$  and that  $g$  is continuously differentiable in a neighborhood of  $\alpha$ , with

$$|g'(\alpha)| < 1. \quad (3.8)$$

Then the conclusions of Theorem 3.5 still hold, provided that  $x_0$  is chosen sufficiently close to  $\alpha$ .

**Example 3.7.** Refer back to the equation  $x^2 - 3 = 0$  in Example 3.1, with  $\alpha = \sqrt{3}$ . Let us see which of the four iterative methods are convergent, for  $x_0$  sufficiently close to  $\alpha$ .

**Solution.**

(a)

$$g(x) = x^2 + x - 3, \quad g'(x) = 2x + 1, \quad g'(\alpha) = 2\sqrt{3} + 1 > 1,$$

so this method *does not converge*.

(b)

$$g(x) = \frac{3}{x}, \quad g'(x) = -\frac{3}{x^2}, \quad g'(\alpha) = -1,$$

so this method does not converge, either.

(c)

$$g(x) = \frac{1}{2}\left(x + \frac{3}{x}\right), \quad g'(x) = \frac{1}{2}\left(1 - \frac{3}{x^2}\right), \quad g'(\alpha) = 0.$$

This method *will converge* at least linearly.

(d)

$$g(x) = x + c(x^2 - 3), \quad g'(x) = 1 + 2cx, \quad g'(\alpha) = 1 + 2c\sqrt{3}.$$

For convergence, pick  $c$  such that  $|g'(\alpha)| < 1$ , i.e.,  $-\frac{1}{\sqrt{3}} < c < 0$ .

For a good rate of linear convergence, pick  $c$  such that  $1 + 2c\sqrt{3} \approx 0$ , or  $c \approx -\frac{1}{2\sqrt{3}}$ , for example,

$$c = -\frac{1}{4}.$$

■

**Example 3.8.** How many real roots does the equation

$$x - 1 - \arctan x = 0 \quad (3.9)$$

have? Will the iterative method

$$x_{n+1} = 1 + \arctan x_n \quad (3.10)$$

converge? For what starting values  $x_0$ ? Find a bound for the error.

**Solution.** First, let us recall that the function  $\arctan x$  is defined on the entire  $\mathbb{R}$ , but takes values *only* in the interval  $\left(-\frac{\pi}{2}, \frac{\pi}{2}\right)$  (see its graph below). That means that

$$-\frac{\pi}{2} < \arctan x < \frac{\pi}{2}, \quad \forall x \in \mathbb{R}.$$

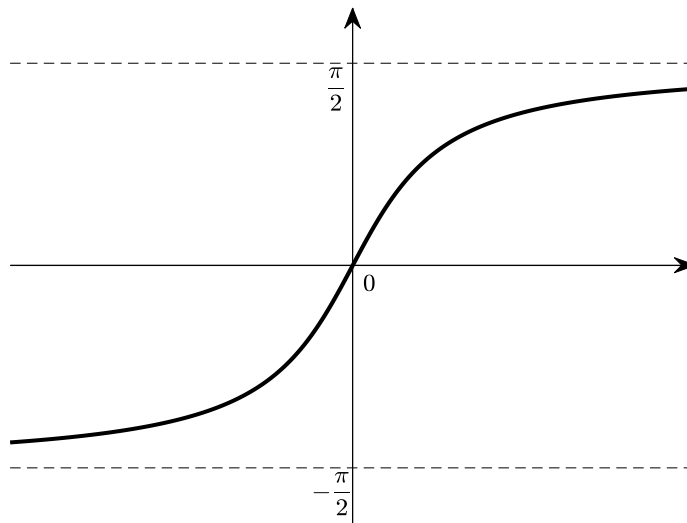


Fig. 5: Function  $\arctan x$

Now, to find the number of real roots, let

$$f(x) = x - 1 - \arctan x.$$

Then

$$f'(x) = 1 - \frac{1}{1+x^2} = \frac{x^2}{1+x^2} \geq 0$$

and is 0 only for  $x = 0$ . The table of variation of  $f$  is

$x$	$-\infty$	$0$	$\infty$
$f$	$-$	$-$	$+$

So there is only one real root  $\alpha > 0$ . To locate it better, compute a few more values:

$$\begin{aligned} f(1) &= 1 - 1 - \frac{\pi}{4} = -\frac{\pi}{4} < 0, \\ f\left(1 + \frac{\pi}{2}\right) &= 1 + \frac{\pi}{2} - 1 - \arctan\left(1 + \frac{\pi}{2}\right) = \frac{\pi}{2} - \arctan\left(1 + \frac{\pi}{2}\right) > 0, \end{aligned}$$

so  $\alpha \in \left(1, 1 + \frac{\pi}{2}\right)$ .

To study the iterative method (3.10), let

$$g(x) = 1 + \arctan x.$$

Now equation (3.9) can be written in the fixed-point form  $x = g(x)$  and the iteration (3.10) is given by  $x_{n+1} = g(x_n)$ .

Let us see if we can use Theorem 3.5, a *global* result, i.e., find an interval  $[a, b]$  such that  $g([a, b]) \subseteq [a, b]$ . Since  $\arctan x \leq \frac{\pi}{2}$ , it follows that  $g(x) = 1 + \arctan x \leq 1 + \frac{\pi}{2}$ , for all  $x \in \mathbb{R}$ . Also,

$$g'(x) = \frac{1}{1+x^2},$$

which is strictly positive for all  $x \in \mathbb{R}$ . That means that  $g$  is strictly increasing on  $\mathbb{R}$ . So, for  $x \in \left[1, 1 + \frac{\pi}{2}\right]$ , we have

$$g(1) \leq g(x) \leq g\left(1 + \frac{\pi}{2}\right).$$

But

$$\begin{aligned} g(1) &= 1 + \frac{\pi}{4} > 1 \quad \text{and} \\ g\left(1 + \frac{\pi}{2}\right) &= 1 + \arctan\left(1 + \frac{\pi}{2}\right) < 1 + \frac{\pi}{2}. \end{aligned}$$

Thus,

$$g\left(\left[1, 1 + \frac{\pi}{2}\right]\right) \subseteq \left[1, 1 + \frac{\pi}{2}\right].$$

Now,  $g''(x) = -\frac{2x}{(1+x^2)^2}$ , which is strictly negative on  $\left[1, 1 + \frac{\pi}{2}\right]$ , so  $g'$  is strictly decreasing on that interval. Then, for all  $x \in \left[1, 1 + \frac{\pi}{2}\right]$ ,

$$g'(x) \leq g'(1) = \frac{1}{2} < 1.$$

So, by Theorem 3.5 with  $\lambda = \frac{1}{2}$ , the iteration (3.10),  $x_{n+1} = g(x_n) = 1 + \arctan x_n$  converges to  $\alpha$ , for any starting value  $x_0 \in \left[1, 1 + \frac{\pi}{2}\right]$  and we have the error estimate

$$|x_n - \alpha| \leq \frac{1}{2^{n-1}} |x_1 - x_0|.$$

The exact solution with 10 correct decimals is  $\alpha = 2.1322679602$ . Indeed, the convergence of the iteration (3.10) is quite fast, as seen in Table 4, for various values of  $x_0 \in \left[1, 1 + \frac{\pi}{2}\right]$ .

$n$	$x_0 = 1$		$x_0 = 1 + \pi/4$		$x_0 = 1 + \pi/2$	
	$x_n$	$ x_n - \alpha $	$x_n$	$ x_n - \alpha $	$x_n$	$ x_n - \alpha $
1	1.78540	$3.47e-1$	2.06023	$7.20e-2$	2.19982	$6.76e-2$
2	2.06023	$7.20e-2$	2.11891	$1.34e-2$	2.14414	$1.19e-2$
3	2.11891	$1.34e-2$	2.12985	$2.42e-3$	2.13440	$2.13e-3$
4	2.12985	$2.42e-3$	2.13183	$4.37e-4$	2.13265	$3.84e-4$
5	2.13183	$4.37e-4$	2.13219	$7.90e-5$	2.13234	$6.89e-5$
6	2.13219	$7.90e-5$	2.13225	$1.44e-5$	2.13228	$1.22e-5$
7	2.13225	$1.44e-5$	2.13227	$2.80e-6$	2.13227	$2.01e-6$
8	2.13227	$2.80e-6$	2.13227	$6.97e-7$	2.13227	$1.70e-7$

Table 4: Example 3.8

■

## 3.2 Higher Order One-Point Iteration Methods

Let us recall the main fixed-point iteration results from last time.

**Theorem 3.1.** Let  $g \in C^1[a, b]$ , such that  $g([a, b]) \subseteq [a, b]$  and

$$\lambda := \max_{x \in [a, b]} |g'(x)| < 1. \quad (3.1)$$

Then:

- a) Function  $g$  has a unique fixed point  $\alpha \in [a, b]$ .
- b) For any initial choice  $x_0 \in [a, b]$ , the sequence  $x_{n+1} = g(x_n)$  converges to  $\alpha$ , as  $n \rightarrow \infty$ .
- c)  $|x_n - \alpha| \leq \lambda^n |x_0 - \alpha| \leq \frac{\lambda^n}{1 - \lambda} |x_1 - x_0|$ ,  $n \geq 1$ .
- d)

$$\lim_{n \rightarrow \infty} \frac{x_{n+1} - \alpha}{x_n - \alpha} = g'(\alpha), \quad (3.2)$$

so, if  $g'(\alpha) \neq 0$ , the iterative method  $x_{n+1} = g(x_n)$  is linearly convergent to the root  $\alpha$  with rate of convergence bounded by  $\lambda$ .

**Theorem 3.2.** Assume  $\alpha$  is a fixed point of  $g$  and that  $g$  is continuously differentiable in a neighborhood of  $\alpha$ , with

$$|g'(\alpha)| < 1. \quad (3.3)$$

Then the conclusions of Theorem 3.1 still hold, provided that  $x_0$  is chosen sufficiently close to  $\alpha$ .

So far, there isn't much information in the case  $g'(\alpha) = 0$ , although the convergence is clearly quite good. Moreover, what happens if the derivatives of  $g$  of up to some order are all 0 at  $\alpha$ ? Can we expect a *faster* convergence? The answer is in the following result.

**Theorem 3.3.** Assume  $\alpha$  is a fixed point of  $g$  and that  $g$  is  $p$  times continuously differentiable for all  $x$  near  $\alpha$ , for some  $p \geq 2$ . Furthermore, assume that

$$g'(\alpha) = \dots = g^{(p-1)}(\alpha) = 0, \quad g^{(p)}(\alpha) \neq 0. \quad (3.4)$$

Then, if the initial value  $x_0$  is chosen sufficiently close to  $\alpha$ , the iteration  $x_{n+1} = g(x_n)$  converges to  $\alpha$  with order of convergence  $p$ , and

$$\lim_{n \rightarrow \infty} \frac{x_{n+1} - \alpha}{(x_n - \alpha)^p} = \frac{1}{p!} g^{(p)}(\alpha). \quad (3.5)$$



*Proof.* Since  $g'(\alpha) = 0$ , by Theorem 3.2, it follows that the iterative method  $x_{n+1} = g(x_n)$  converges to  $\alpha$ , if  $x_0$  is sufficiently close to  $\alpha$ .

For the order of convergence, we use the Taylor series expansion of  $g$  around  $\alpha$ :

$$x_{n+1} = g(x_n) = g(\alpha) + (x_n - \alpha)g'(\alpha) + \cdots + \frac{(x_n - \alpha)^{p-1}}{(p-1)!}g^{(p-1)}(\alpha) + \frac{(x_n - \alpha)^p}{p!}g^{(p)}(\xi_n),$$

for some  $\xi_n$  between  $x_n$  and  $\alpha$ . Using (3.4) and the fact that  $g(\alpha) = \alpha$ , we get

$$\begin{aligned} x_{n+1} - \alpha &= \frac{(x_n - \alpha)^p}{p!}g^{(p)}(\xi_n), \\ \frac{x_{n+1} - \alpha}{(x_n - \alpha)^p} &= \frac{1}{p!}g^{(p)}(\xi_n). \end{aligned}$$

Letting  $n \rightarrow \infty$ , both  $x_n, \xi_n \rightarrow \alpha$  and, hence, (3.5) follows. □

**Example 3.4.** Recall Newton's iterative method

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}, \quad n \geq 0. \quad (3.6)$$

Let us analyze it by this new result.

**Solution.** We have

$$g(x) = x - \frac{f(x)}{f'(x)},$$

for a simple root of  $f$ ,  $\alpha$ , which means  $f(\alpha) = 0$  and  $f'(\alpha) \neq 0$ . We have

$$g'(x) = 1 - \frac{(f'(x))^2 - f(x)f''(x)}{(f'(x))^2} = \frac{f(x)f''(x)}{(f'(x))^2}.$$

We compute the second derivative, but discard the argument  $x$ :

$$g'' = \frac{(f'f'' + ff''')(f')^2 - 2f'f'' \cdot ff''}{(f')^4}.$$

Then

$$\begin{aligned}g(\alpha) &= \alpha, \\g'(\alpha) &= 0, \\g''(\alpha) &= \frac{(f'(\alpha)f''(\alpha))(f'(\alpha))^2}{(f'(\alpha))^4} = \frac{f''(\alpha)}{f'(\alpha)}\end{aligned}$$

and Theorem 3.3 gives the previously found quadratic convergence ( $p = 2$ ) and error estimate

$$\lim_{n \rightarrow \infty} \frac{x_{n+1} - \alpha}{(x_n - \alpha)^2} = \frac{1}{2}g''(\alpha) = \frac{f''(\alpha)}{2f'(\alpha)}.$$

■

**Example 3.5.** Let us revisit the problem from Example 3.7 (Lecture 11): equation  $x^2 - 3 = 0$ , with  $\alpha = \sqrt{3}$ .

**Solution.** Last time we saw several ways of rewriting the equation in the form  $g(x) = x$ , some “better” than others, from the convergence point of view.

Let us use Newton’s iteration. We have

$$f(x) = x^2 - 3, \quad f'(x) = 2x,$$

so,

$$g(x) = x - \frac{f(x)}{f'(x)} = x - \frac{x^2 - 3}{2x} = \frac{1}{2}\left(2x - x + \frac{3}{x}\right) = \frac{1}{2}\left(x + \frac{3}{x}\right), \quad g(\alpha) = \alpha,$$

which was one of the methods discussed (part (c)). Further, we have

$$\begin{aligned}g'(x) &= \frac{1}{2}\left(1 - \frac{3}{x^2}\right), \quad g'(\alpha) = 0, \\g''(x) &= \frac{3}{x^3}, \quad g''(\alpha) = \frac{1}{\sqrt{3}} \neq 0.\end{aligned}$$

So, indeed, the iteration  $x_{n+1} = \frac{1}{2}\left(x_n + \frac{3}{x_n}\right)$  converges quadratically ( $p = 2$ ) to  $\alpha$ .

As a side note, rewriting the equation  $x^2 - 3 = 0$  as  $x = \frac{1}{2}\left(x + \frac{3}{x}\right)$  (which leads to faster convergence) was *not* a “lucky guess”, it is actually Newton’s method.

■

**Example 3.6.** Consider the equation

$$x = g(x) = cx(1 - x), \quad (3.7)$$

with  $c \neq 0$ . This is called a *logistic equation* and is of great interest in the *mathematical theory of chaos*. This equation has one nonzero solution, denoted by  $\alpha_c$ . For what values of  $c$  will the iteration  $x_{n+1} = g(x_n)$  converge to  $\alpha_c$  (provided that  $x_0$  is chosen sufficiently close to  $\alpha_c$ )? Determine the convergence order.

**Solution.** The nonzero solution of equation (3.7) is given by

$$c(1 - x) = 1, \quad x = \alpha_c = 1 - \frac{1}{c} = \frac{c - 1}{c}.$$

We have

$$\begin{aligned} g(x) &= c(x - x^2), \quad g(\alpha_c) = \alpha_c, \\ g'(x) &= c(1 - 2x), \quad g'(\alpha_c) = c\left(1 - 2 + \frac{2}{c}\right) = 2 - c, \\ g''(x) &= -2c. \end{aligned}$$

In order to have the iteration  $x_{n+1} = g(x_n)$  converge to  $\alpha_c$ , we impose the condition

$$|g'(\alpha_c)| < 1 \iff |2 - c| < 1 \iff c \in (1, 3).$$

So, for any  $1 < c < 3$ , the method converges (at least) linearly, with rate of convergence  $|c - 2|$ .

Now, if  $g'(\alpha_c) = 0$ , i.e.,  $c = 2$ , then the convergence is quadratic ( $p = 2$ ) and

$$\lim_{n \rightarrow \infty} \frac{x_{n+1} - \alpha_c}{(x_n - \alpha_c)^2} = \frac{1}{2}g''(\alpha_c) = \frac{-2c}{2} = -c = -2.$$

■

**Example 3.7.** Does the iteration

$$x_{n+1} = x_n^5 - 10x_n^3 - 20x_n^2 - 15x_n - 5$$

converge to  $\alpha = -1$ , provided that  $x_0$  is chosen sufficiently close to  $\alpha$ ? If so, determine the convergence order and bound the error.

**Solution.** We have  $g(x) = x^5 - 10x^3 - 20x^2 - 15x - 5$  and

$$g(-1) = -1 + 10 - 20 + 15 - 5 = -1,$$

so  $\alpha = -1$  is a fixed point of  $g$ .

For convergence (and the order of convergence), we compute successively

$$\begin{aligned} g'(x) &= 5x^4 - 30x^2 - 40x - 15, \quad g'(-1) = 5 - 30 + 40 - 15 = 0, \\ g''(x) &= 20x^3 - 60x - 40 = 20(x^3 - 3x - 2), \quad g''(-1) = 20(-1 + 3 - 2) = 0, \\ g'''(x) &= 20(3x^2 - 3) = 60(x^2 - 1), \quad g'''(-1) = 0, \\ g^{(4)}(x) &= 120x, \quad g^{(4)}(-1) = -120 \neq 0. \end{aligned}$$

So, for  $x_0$  sufficiently close to  $-1$ , we have

$$\lim_{n \rightarrow \infty} \frac{x_{n+1} + 1}{(x_n + 1)^4} = \frac{1}{4!} g^{(4)}(-1) = -\frac{1}{24} \cdot 120 = -5.$$

Then, it follows that

$$|x_{n+1} + 1| \leq c_1 |x_n + 1|^4 \leq \dots \leq c |x_0 + 1|^{4^{n+1}}.$$

So, for  $|x_0 + 1| < 1$ , the method converges with order of convergence  $p = 4$  and the error estimate above. ■

## 4 Numerical Approximation of Multiple Roots

**Definition 4.1.** We say that a function  $f$  has a **root**  $\alpha$  **of multiplicity**  $m > 1$  if

$$f(x) = (x - \alpha)^m h(x), \quad h \text{ continuous at } x = \alpha \text{ and } h(\alpha) \neq 0. \quad (4.1)$$

We restrict our discussion to the case where  $m$  is a positive integer, although some of our considerations are equally valid for non-integer values. If  $h$  is smooth enough at  $x = \alpha$ , then (4.1) is equivalent to

$$f(\alpha) = f'(\alpha) = \dots = f^{(m-1)}(\alpha) = 0, \quad f^{(m)}(\alpha) \neq 0. \quad (4.2)$$

There are several challenges in approximating multiple roots.

## Uncertainty

When finding a root of any function on a computer, there is always an *interval of uncertainty* about the root, due to measuring/rounding/truncation errors, and this is made worse when the root is multiple.

**Example 4.2.** Consider evaluating the two functions

$$\begin{aligned}f_1(x) &= x^2 - 3, \\f_2(x) &= x^2(x^2 - 6) + 9.\end{aligned}$$

Notice that  $\alpha = \sqrt{3}$  has multiplicity 1 as a root of  $f_1$  and multiplicity 2 as a root of  $f_2$ , since

$$f_2'(x) = 4x(x^2 - 3).$$

Using four-digit decimal arithmetic, we have

$$\begin{aligned}f_1(x) &< 0, \text{ for } x \leq 1.731, \\f_1(1.732) &= 0, \text{ and} \\f_1(x) &> 0, \text{ for } x > 1.733,\end{aligned}$$

so  $\alpha \in (1.731, 1.733)$ . But for  $f_2$ ,

$$f_2(x) = 0, \text{ for } 1.726 \leq x \leq 1.738,$$

implying that  $\alpha \in [1.726, 1.738]$ , thus limiting the amount of accuracy that can be attained in finding a root of  $f_2$ .

## Loss of Precision

Another problem with multiple roots is that the earlier rootfinding methods will not perform as well when the root being sought is multiple. Let us investigate this for Newton's method. We consider Newton's method as a fixed-point iteration

$$x_{n+1} = g(x_n), \quad g(x) = x - \frac{f(x)}{f'(x)}, \quad x \neq \alpha, \quad f(x) = (x - \alpha)^m h(x), \quad m > 1.$$

We have

$$f'(x) = m(x - \alpha)^{m-1} h(x) + (x - \alpha)^m h'(x) = (x - \alpha)^{m-1} [m h(x) + (x - \alpha) h'(x)],$$

so,

$$\begin{aligned} g(x) &= x - \frac{(x - \alpha)^m h(x)}{(x - \alpha)^{m-1} [m h(x) + (x - \alpha) h'(x)]} \\ &= x - (x - \alpha) \frac{h(x)}{m h(x) + (x - \alpha) h'(x)} = x - (x - \alpha) \varphi(x), \end{aligned}$$

where

$$\varphi(x) \stackrel{\text{not}}{=} \frac{h(x)}{m h(x) + (x - \alpha) h'(x)}, \quad \varphi(\alpha) = \frac{1}{m}.$$

Then,

$$\begin{aligned} g'(x) &= 1 - \varphi(x) - (x - \alpha) \varphi'(x), \\ g'(\alpha) &= 1 - \varphi(\alpha) = 1 - \frac{1}{m} \neq 0. \end{aligned}$$

Thus, in this case, Newton's method converges only *linearly*, with rate of convergence  $1 - \frac{1}{m} < 1$ .

One way to fix this loss of accuracy would be to change the problem into an equivalent one: instead of solving  $f(x) = 0$  which has  $\alpha$  as a multiple root, consider the equation

$$u(x) := \frac{f(x)}{f'(x)} = 0, \tag{4.3}$$

for which  $\alpha$  is a *simple* root. Then Newton's method is defined by

$$x_{n+1} = x_n - \frac{u(x_n)}{u'(x_n)}. \tag{4.4}$$

We have

$$u' = \frac{(f')^2 - f f''}{(f')^2}, \quad \frac{u}{u'} = \frac{\frac{f}{f'}}{\frac{(f')^2 - f f''}{(f')^2}} = \frac{f f'}{(f')^2 - f f''},$$

so Newton's method is given by

$$x_{n+1} = x_n - \frac{f(x_n)f'(x_n)}{(f'(x_n))^2 - f(x_n)f''(x_n)}, \quad n \geq 0. \quad (4.5)$$

Although this method restores the order of convergence  $p = 2$ , it has several disadvantages: it requires the computation of the second derivative  $f''$ , it involves more complex computations than the original method, and the denominator in (4.5) can take very small values, as  $x_n \rightarrow \alpha$ .

A better alternative is to modify the *method*, instead of the *function*.

## Newton's Method for Multiple Roots

To improve Newton's method, we would like a function  $g$  for which  $g'(\alpha) = 0$  (as before), even for multiple roots. Consider the following idea: if  $\alpha$  is a root of  $f$ , then in the vicinity of  $\alpha$ , we have

$$f(x) = (x - \alpha)^m h(x) \approx (x - \alpha)^m c,$$

for some constant  $c$ . Then

$$\begin{aligned} f'(x) &\approx m(x - \alpha)^{m-1} c, \quad \frac{f(x)}{f'(x)} \approx \frac{x - \alpha}{m}, \\ x - \alpha &\approx m \frac{f(x)}{f'(x)}, \quad \alpha \approx x - m \frac{f(x)}{f'(x)}. \end{aligned}$$

Thus, we define **Newton's method for multiple roots** by

$$x_{n+1} = x_n - m \frac{f(x_n)}{f'(x_n)}. \quad (4.6)$$

Now, indeed, it is easy to check (by our earlier computations) that for  $g(x) = x - m \frac{f(x)}{f'(x)}$ , we have

$$g(\alpha) = \alpha, \quad g'(\alpha) = 1 - m \varphi(\alpha) = 0,$$

so the method converges quadratically again and

$$\lim_{n \rightarrow \infty} \frac{x_{n+1} - \alpha}{(x_n - \alpha)^2} = \frac{1}{2} g''(\alpha).$$

## 5 Newton's Method for Nonlinear Systems

Many of the methods considered in the previous sections can be generalized to the multidimensional case, i.e. to *systems* of nonlinear equations. These problems are widespread in applications, and they are varied in form. There is a great variety of methods for the solution of such systems. We only consider the two-dimensional case

$$\begin{aligned} f_1(x_1, x_2) &= 0 \\ f_2(x_1, x_2) &= 0, \end{aligned} \tag{5.1}$$

or, in vector notation,

$$\mathbf{f}(\mathbf{x}) = \mathbf{0}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \mathbf{f}(\mathbf{x}) = \begin{bmatrix} f_1(x_1, x_2) \\ f_2(x_1, x_2) \end{bmatrix}. \tag{5.2}$$

The one-point iteration theory discussed in the previous sections still stands, with appropriate adjustments (norm instead of absolute value and *Jacobian matrix* instead of derivative). Newton's method is derived similarly with the one-dimensional case, considering Taylor series expansions of each  $f_i, i = 1, 2$  and expanding  $f_i(\mathbf{x})$  about  $\mathbf{x}_0 = [x_{1,0} \ x_{2,0}]^T$ . We get

$$\mathbf{x}_{n+1} = \mathbf{x}_n - (\mathbf{J}_f(\mathbf{x}_n))^{-1} \mathbf{f}(\mathbf{x}_n), \quad n \geq 0, \tag{5.3}$$

where  $\mathbf{J}_f(\mathbf{x}_n)$  is the Jacobian matrix of  $\mathbf{f}$  at  $\mathbf{x}_n$

$$\mathbf{J}_f(\mathbf{x}_n) = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_1(\mathbf{x})}{\partial x_2} \\ \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} \end{bmatrix} \tag{5.4}$$

This is **Newton's method for nonlinear systems**.

In actual practice, we *do not invert*  $\mathbf{J}_f(\mathbf{x}_n)$ , particularly for systems of more than two equations. Instead we solve a linear system for a correction term to  $\mathbf{x}_n$ :

$$\begin{aligned} \mathbf{J}_f(\mathbf{x}_n) \delta_{n+1} &= -\mathbf{f}(\mathbf{x}_n), \\ \mathbf{x}_{n+1} &= \mathbf{x}_n + \delta_{n+1}. \end{aligned} \tag{5.5}$$

This is more efficient in computation time, requiring only about one-third as many operations as inverting  $\mathbf{J}_f(\mathbf{x}_n)$ .



**Example 5.1.** Solve the system

$$\begin{aligned} f_1(x_1, x_2) &\equiv 4x_1^2 + x_2^2 - 4 = 0 \\ f_2(x_1, x_2) &\equiv x_1 + x_2 - \sin(x_1 - x_2) = 0. \end{aligned}$$

**Solution.** There are only two roots, one near  $(1, 0)$  and its reflection about the origin near  $(-1, 0)$  (see Figure 1).

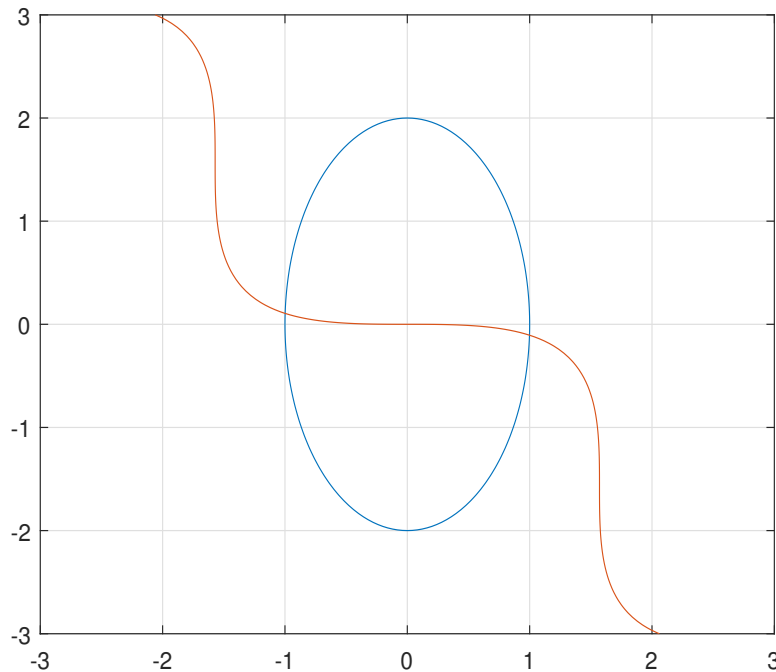


Fig. 1: Example 5.1

Using Newton's method with  $\mathbf{x}_0 = [1 \ 0]^T$  we obtain the results in Table 1.

$n$	$x_{1,n}$	$x_{2,n}$	$f_1(\mathbf{x}_n)$	$f_2(\mathbf{x}_n)$
0	1.0	0.0	0.0	$1.59e - 1$
1	1.0	$-0.1029207154$	$1.06e - 2$	$4.55e - 3$
2	$0.9986087598$	$-0.1055307239$	$1.46e - 5$	$6.63e - 7$
3	$0.9986069441$	$-0.1055304923$	$1.32e - 11$	$1.87e - 12$

Table 1: Newton's Method for Example 5.1

■