

MA637 - Numerical Analysis and Computing Lecture Notes

Abhinav Jha
Indian Institute of Technology, Gandhinagar
Winter Semester 2024-2025

Preface

These notes are designed to provide a structured and comprehensive understanding of the course content. They will cover key topics, concepts, and computational techniques that are fundamental to numerical analysis. Please note that this is the first iteration (Version 0.1.0) of the notes and hence there is a chance that some of the content is incorrect. If you find some flaws, please email me at abhinav.jha@iitgn.ac.in.

Contents

1	Interpolation	7
1.1	Polynomial Interpolation	9
1.1.1	Drawbacks	10
1.2	Lagrange Interpolation	11
1.2.1	Drawbacks	14
1.2.2	Runge Phenomena	15
1.3	Newton Divided Difference Interpolation	17
1.3.1	Computational Complexity	21
1.3.2	Forward Difference Formula	21
1.4	Hermite Interpolation	22
1.4.1	Hermite Polynomials using Divided Difference	25
1.5	Spline Interpolation	27
1.5.1	Cubic Splines	27
1.5.2	B-Splines	31
2	System of Equations	35
2.1	Gaussian Elimination	36
2.1.1	Computational Complexity	39
2.1.2	Gauss-Jordan Algorithm	42
2.2	Matrix Factorisation	42
2.2.1	LU Decomposition	42
2.2.2	LDL ^T Decomposition	48
2.2.3	Cholesky Decomposition	54
2.3	Iterative Methods	56
2.3.1	Jacobi Method	61
2.3.2	Gauss-Seidel Method	63
2.3.3	Successive Over Relaxation	69
2.3.4	Condition Number	72
2.4	Least Square Methods	74
2.4.1	QR Decomposition	76
3	Computing	79
3.1	Good Practices in Coding	80
3.1.1	Variable Initialization and Naming	80
3.1.2	Reusability and Modularity	81
3.2	Testing and Continuous Integration	81

3.3	Introduction to Computing Using Python	81
3.3.1	Variables	81
3.3.2	Arithmetic Operations	82
3.3.3	Compound Assignment	82
3.3.4	Logical Operations	83
3.4	Conditional Statements	83
3.5	Recursive Statements	84
3.5.1	For Loop	84
3.5.2	Custom Step Size	85
3.5.3	Break and Continue	85
3.5.4	Nested Loops	85
3.6	Functions	86
3.7	NumPy Library	87
3.7.1	Arrays and Matrices	87
3.7.2	Linspace	88
3.7.3	Mathematical Functions	88

Chapter 1

Interpolation

Interpolation has various definitions depending on the search engine. For example, *Wikipedia* states,

“Interpolation is a type of estimation, a method of constructing (finding) new data points based on the range of a discrete set of known data points.”

Blackphoto says,

“It is a technique used by digital scanners, cameras, and printers to increase the size of an image in pixels by averaging the colour and brightness values of surrounding pixels.”

One can see such an example in image processing. A rather famous (or infamous) example is the *Ecce Homo* painting (see Fig. 1 (left)). This is a fresco painting painted in 1330 by the Spanish painter Elías García Martínez depicting Jesus Christ. With wear and tear, the painting got degraded, and in 2012, an 81-year-old lady, Cecilia Giménez “tried” to restore it (see Fig. 1 (right)); as we can see, it is not very good, and hence it was named *Ecce Mono*. We can get much better results with modern image processing techniques (which inherently use a form of interpolation).



Figure 1.1: Elías García Martínez, *Ecce Homo*: The leftmost photograph, taken in 2010, shows some initial flaking of the paintwork. The central photograph was taken in July 2012, just a month before the attempted restoration, showing the extent of damage and deterioration. The rightmost photograph documents the artwork following Giménez’s efforts to repair it.

In interpolation, we try to approximate general functions by a “simple” class of functions. In analysis, a powerful result connects the continuous functions and polynomial approximation:

the Weierstrass Approximation theorem given by Karl Weierstrass.



Figure 1.2: Karl Weierstrass: 31 October 1815-19 February 1897

Theorem 1.1. [1, Theorem 5.4.14] (**Weierstrass Approximation Theorem**) Let $f \in \mathcal{C}[a, b]$. Then for each $\varepsilon > 0$ there exists a polynomial $p(x)$ with the property that

$$|f(x) - p(x)| < \varepsilon \quad \text{for all } x \in [a, b].$$

This theorem is important because polynomials have excellent differentiation and integration properties as their derivatives and integrals are polynomials. Another interpretation of Theorem 1 is that given a continuous function on a closed and bounded interval, there exists a polynomial, i.e. , as “close” to the given function as desired.

But in analysis, there exists one more kind of polynomial approximation, and that is the Taylor’s theorem

Theorem 1.2. [1, Theorem 6.4.1] (**Taylor’s Theorem**) Suppose $f \in \mathcal{C}^n[a, b]$ and $f^{(n+1)}$ exists on $[a, b]$ and $x_0 \in [a, b]$. For every $x \in [a, b]$ there exists a number $\xi(x) \in [x_0, x]$ with

$$f(x) = P_n(x) + R_n(x),$$

where $P_n(x) = \sum_{k=0}^n \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k$ and $R_n(x) = \frac{f^{(n+1)}(\xi(x))}{(n+1)!} (x - x_0)^{n+1}$.



Figure 1.3: Brook Taylor: 18 August 1685-29 December 1731

There are two issues here:

1. We need to know the higher derivatives of $f(x)$.

2. This is a *local* approximation, i.e., the approximation is excellent near x_0 but we need certain global approximation. For example, if we do the Taylor series expansion for $\exp(x)$ around zero, then it becomes worse as we move away from zero (see Fig. 1.4).

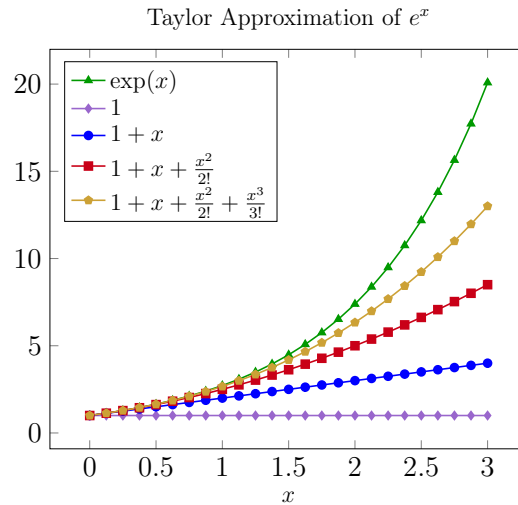


Figure 1.4: Taylor polynomials for exponential function approximated at $x = 0$.

However, it should be noted that the Taylor theorem is still a powerful result whose main purpose is the derivation of numerical techniques and error estimation.

1.1 Polynomial Interpolation

Suppose we have a finite set of data points f_i associated with parameters x_i . We want to depict these data points as a function $f(x)$ with the property that $f(x_i) = f_i$. This is clearly not well-defined since there are many such functions. But if we restrict to finite-dimensional spaces (such as polynomials), then we can define such functions, or to be more precise, the process is well-defined.

We first start with the idea of polynomial interpolation. Polynomials representing an unknown functional dependence of the discrete set of data points are called *interpolants*. The main problem that we want to tackle with interpolation is:

Problem: Given a set of $(n + 1)$ data points say $\{(x_i, f_i)\}_{i=0}^n$ find a polynomial $p_n(x)$ of degree n satisfying

$$p_n^{\forall}(x_i) = f_i \quad \text{for all } i = 0, 1, \dots, n.$$

Now the general form of a polynomial $p_n^{\forall}(x)$ is given by

$$p_n^{\forall}(x) = \sum_{i=0}^n c_i x^{n-i} := c_0 x^n + c_1 x^{n-1} + \dots + c_n,$$

for coefficients $c_i \in \mathbb{R}$. Since each polynomial of degree n can be determined by $(n + 1)$ coefficients, we can re-write the above problem as solving the following system of equations:

$$c_0 x_i^n + c_1 x_i^{n-1} + \dots + c_{n-1} x_i + c_n = f_i \quad i = 0, 1, 2, \dots, n,$$

or in the matrix form

$$\mathbf{V}\mathbf{c} = \mathbf{f} \quad (1.1)$$

where

$$\mathbf{V} = \begin{bmatrix} x_0^n & x_0^{n-1} & \dots & x_0 & 1 \\ x_1^n & x_1^{n-1} & \dots & x_1 & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_n^n & x_n^{n-1} & \dots & x_n & 1 \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} c_0 \\ c_1 \\ \vdots \\ c_n \end{bmatrix}, \quad \text{and} \quad \mathbf{f} = \begin{bmatrix} f_0 \\ f_1 \\ \vdots \\ f_n \end{bmatrix}.$$

This system has a unique solution if \mathbf{V} is invertible [5, Theorem 3.10], which is equivalent to saying that $\det(\mathbf{V}) \neq 0$. This matrix \mathbf{V} is called as the *Vandermonde matrix* and its determinant is given by

$$\det(\mathbf{V}) = \prod_{i=0}^{n-1} \prod_{j=i+1}^n (x_i - x_j).$$



Figure 1.5: Alexandre-Théophile Vandermonde: 28 February 1735 – 1 January 1796

This determinant is non-zero if we have distinct points. Hence, from now on, we assume we have $(n + 1)$ distinct points.

The algorithm for using the polynomial interpolation using Vandermonde matrix for finding solution at a given point x_{eval} is given in Algorithm 1.

Note that we have introduced the notation $p_n^{\mathbf{V}}(x)$ only to denote the polynomial $p_n(x)$ computed using the Vandermonde matrix.

1.1.1 Drawbacks

Even though Eq. (1.1) has a perfect mathematical solution, computationally, it is not that good. The reason being Vandermonde matrices are ill-conditioned¹ (we will do conditioning of a system in the following chapters). The matrix \mathbf{V} has a large condition number leading to inaccurate solutions.

To understand why the Vandermonde matrix is ill-conditioned for large n , we can plot x^k for $0 \leq k \leq n$ in $[0, 1]$ (see Fig. 1.6). Even though x^k are distinct for larger k , they tend to look the same. As a result, it is harder to identify projections of a particular polynomial $p_n^{\mathbf{V}}(x)$ into the nearly collinear basis of monomials x^k for large k .

¹**Ill-Conditioned System:** In numerical analysis, the condition number of a function quantifies the extent to which the output can change in response to small variations in the input. It measures a function's sensitivity to input changes or errors, indicating how much an input error can propagate into the output. A problem with a low condition number is said to be *well-conditioned* while a problem with a high condition number is said to be *ill-conditioned*.

Algorithm 1 Vandermonde Interpolation

Given: Data sets $\{(x_i, f_i)\}_{i=0}^n$, Evaluation point x_{eval} .

Find: Interpolated polynomial $p_n^{\mathbb{V}}(x_{\text{eval}})$.

Step 1: Compute Vandermonde Matrix

Initialize an empty Vandermonde matrix \mathbf{V} of size $(n + 1) \times (n + 1)$

for $i = 0$ **to** n **do**

for $j = 0$ **to** n **do**

$\mathbf{V}_{i,j} = x_i^{(n-j)}$

end for

end for

Step 2: Solve the System of Linear Equations

Solve the system $\mathbf{V} \cdot \mathbf{c} = \mathbf{f}$ to get coefficient vector \mathbf{c}

Step 3: Evaluate the Vandermonde Polynomial $p_n^{\mathbb{V}}(x)$ at x_{eval}

Initialize $p_n^{\mathbb{V}}(x_{\text{eval}}) = 0$

for $i = 0$ **to** n **do**

$p_n^{\mathbb{V}}(x_{\text{eval}}) = p_n^{\mathbb{V}}(x_{\text{eval}}) + c_i \cdot x_{\text{eval}}^{(n-i)}$

end for

return $p_n^{\mathbb{V}}(x_{\text{eval}})$

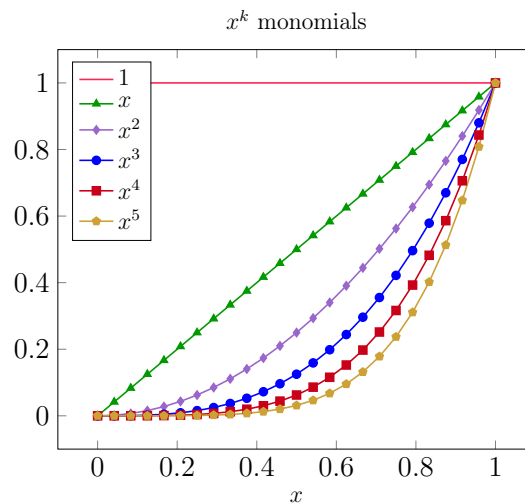


Figure 1.6: Monomial basis x^k for $k = 0, 1, \dots, 5$.

1.2 Lagrange Interpolation

After examining how unstable polynomial interpolation is, we need to develop more stable methods. One of the most known methods is the Lagrange interpolation. The formula was first published by Waring in 1779, rediscovered by Euler in 1783, and published by Lagrange in 1795 (Jeffreys & Jeffreys, 1988).

Let us start with a basic example of two points (x_0, f_0) and (x_1, f_1) , then we define

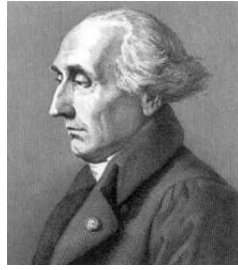


Figure 1.7: Joseph-Louis Lagrange: 25 January 1736-10 April 1813

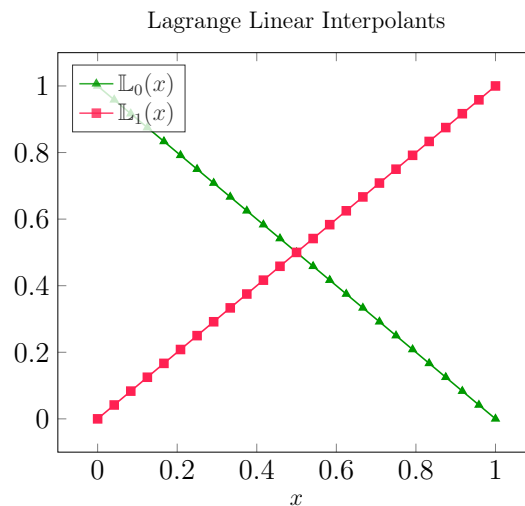
functions:

$$\mathbb{L}_0(x) = \frac{x - x_1}{x_0 - x_1} \quad \text{and} \quad \mathbb{L}_1(x) = \frac{x - x_0}{x_1 - x_0}. \quad (1.2)$$

Then, a linear interpolating polynomial passing through the above points is given by:

$$p_1(x) = \mathbb{L}_0(x)f_0 + \mathbb{L}_1(x)f_1 = \frac{x - x_1}{x_0 - x_1}f_0 + \frac{x - x_0}{x_1 - x_0}f_1,$$

as $\mathbb{L}_0(x_0) = 1$; $\mathbb{L}_0(x_1) = 0$; $\mathbb{L}_1(x_0) = 0$; and $\mathbb{L}_1(x_1) = 1$, we have $p_1(x_0) = f_0$ and $p_1(x_1) = f_1$. This polynomial is called the *Lagrange linear interpolating polynomial*. In fact, this is a unique polynomial. Fig. 1.8 shows $\mathbb{L}_0(x)$ and $\mathbb{L}_1(x)$ for $x_0 = 0$ and $x_1 = 1$.

Figure 1.8: Lagrange linear interpolating polynomials for $x_i = \{0, 1\}$.

What happens if we generalize this concept, i.e., we have $\{(x_i, f_i)\}_{i=0}^n$? In this case we first need to construct for each $i = 0, 1, \dots, n$ a function $\mathbb{L}_{n,i}(x)$ with the property that

$$\mathbb{L}_{n,i}(x_k) = \delta_{ik} \quad \text{for } k = 0, \dots, n.$$

Based on Eq. (1.2) the general form should look like:

$$\mathbb{L}_{n,i}(x) = \prod_{j=0, j \neq i}^n \left(\frac{x - x_j}{x_i - x_j} \right).$$

Then, we can define the polynomial as

$$p_n^{\mathbb{L}}(x) = \sum_{i=0}^n f_i \mathbb{L}_{n,i}(x), \quad (1.3)$$

where we have used the notation $p_n^{\mathbb{L}}(x)$ to denote the interpolating polynomial obtained by the Lagrange interpolation. If the degree of the polynomial is clear we can write $\mathbb{L}_{n,i}(x)$ as $\mathbb{L}_i(x)$. We call $\mathbb{L}_{n,i}(x)$ as the n^{th} Lagrange interpolating polynomial (see Fig. 1.9). One can compute the Lagrange interpolating polynomial using algorithm 2.

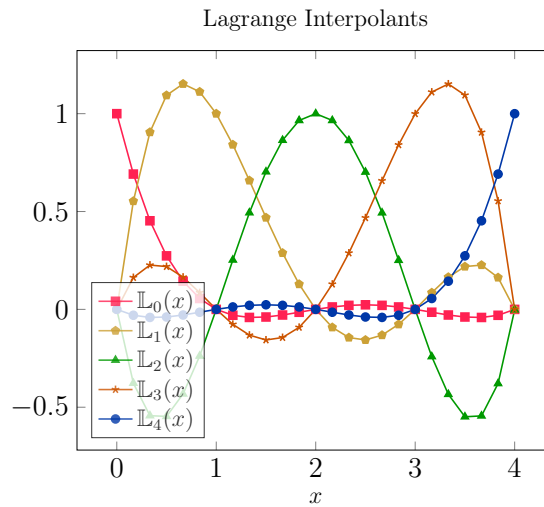


Figure 1.9: Lagrange interpolating polynomials defined over $x_i = 0, 1, 2, 3, 4$.

We have certain remarks for the Lagrange interpolation:

1. We note that in Eq. (1.3) $p_n^{\mathbb{L}}(x)$ maps the linear space \mathbb{R}^{n+1} to the space of polynomials \mathbb{P}_n which is a linear map.
2. We can extend the Lagrange interpolant to any continuous function $f(x)$ by

$$p_n^{\mathbb{L}}f(x) = \sum_{i=0}^n f(x_i)\mathbb{L}_i(x).$$

3. The operator $p_n^{\mathbb{L}}(x)$ is a projection, i.e., $p_n^{\mathbb{L}}q = q$ for all $q \in \mathbb{P}_n$.

Now we present a theorem that tells us about the error obtained using Lagrange interpolation.

Theorem 1.3. Suppose $\{x_0, x_1, \dots, x_n\}$ are distinct numbers in the interval $[a, b]$ and $f \in C^{n+1}[a, b]$. Then for each $x \in [a, b]$ there exists a number $\xi(x) \in (a, b)$ with

$$f(x) = p_n^{\mathbb{L}}(x) + \frac{f^{(n+1)}(\xi(x))}{(n+1)!} \prod_{i=0}^n (x - x_i), \tag{1.4}$$

where $p_n^{\mathbb{L}}(x)$ is given by Eq. (1.3).

Proof. Note that if $x = x_k$ then $f(x_k) = p_n^{\mathbb{L}}(x_k)$ for any $k = 0, 1, \dots, n$. Hence Eq. (1.4) is trivial for any $\xi(x) \in (a, b)$.

Suppose $x \neq x_k$ for any $k = 0, 1, \dots, n$ then define a function g for t in $[a, b]$ as

$$g(t) = f(t) - p_n^{\mathbb{L}}(t) - [f(x) - p_n^{\mathbb{L}}(x)] \prod_{i=0}^n \frac{(t - x_i)}{(x - x_i)}.$$

Since $f \in \mathcal{C}^{n+1}[a, b]$ and $p_n^{\mathbb{L}} \in \mathcal{C}^{\infty}[a, b]$ we have $g \in \mathcal{C}^{n+1}[a, b]$.

Theorem 1.4. [3, Theorem 1.10]/(Generalized Rolle's Theorem)

Suppose $f \in \mathcal{C}[a, b]$ is n -times differentiable on (a, b) . If $f(x) = 0$ at $(n+1)$ distinct points $a \leq x_0 < x_1 < \dots < x_n \leq b$ then there exists a number $c \in (x_0, x_n) (\subset (a, b))$ such that $f^{(n)}(c) = 0$.

For $t = x_k$ for any k , we have

$$g(x_k) = f(x_k) - p_n^{\mathbb{L}}(x_k) = 0.$$

Moreover $g(x) = 0$. Thus $g \in \mathcal{C}^{n+1}[a, b]$ with $(n+2)$ distinct zeros. By Generalized Rolle's theorem 1.2 there exists a $\xi \in (a, b)$ for which $g^{(n+1)}(\xi) = 0$. So,

$$0 = g^{(n+1)}(\xi) = f^{(n+1)}(\xi) - p_n^{\mathbb{L}(n+1)}(\xi) - [f(x) - p_n^{\mathbb{L}}(x)] \frac{d^{n+1}}{dt^{n+1}} \left[\prod_{i=0}^n \frac{(t - x_i)}{(x - x_i)} \right]_{t=\xi}. \quad (1.5)$$

Now, $p_n^{\mathbb{L}}(x)$ is a polynomial of degree at most n . Hence, $p_n^{\mathbb{L}(n+1)}(x) = 0$. Also, $\prod_{i=0}^n \frac{(t-x_i)}{(x-x_i)}$ is a polynomial of degree $(n+1)$ with leading coefficient being $\frac{1}{\prod_{i=0}^n (x-x_i)}$. Hence,

$$\frac{d^{n+1}}{dt^{n+1}} \left(\prod_{i=0}^n \frac{(t - x_i)}{(x - x_i)} \right) = \frac{(n+1)!}{\prod_{i=0}^n (x - x_i)}.$$

Hence, Eq. (1.5) becomes

$$f^{(n+1)}(\xi) - [f(x) - p_n^{\mathbb{L}}(x)] \frac{(n+1)!}{\prod_{i=0}^n (x - x_i)} = 0 \quad \Rightarrow \quad f(x) = p_n^{\mathbb{L}}(x) + \frac{f^{(n+1)}(\xi)}{(n+1)!} \prod_{i=0}^n (x - x_i).$$

□

Note that this error term is similar to Taylor's theorem, but it has information on all the points instead of the error being concentrated along one point.

1.2.1 Drawbacks

Lagrange interpolant suffers from certain drawbacks. The first one is regarding its *computational complexity*². For the evaluation of an unknown point x , we will check the computational

²**Computational Complexity:** Computational complexity measures how hard it is for a computer to solve a problem as the size of the problem increases. It tells us how much time and resources are needed to find a solution.

Algorithm 2 Lagrange Interpolation

Given: Data sets $\{(x_i, f_i)\}_{i=0}^n$, Evaluation point x_{eval} .
Find: Interpolated polynomial $p_n^{\mathbb{L}}(x_{\text{eval}})$.

Step 1: Compute Lagrange Basis Polynomials $\mathbb{L}_i(x)$
for $i = 0$ **to** n **do**
 $\mathbb{L}_i(x_{\text{eval}}) = 1$
 for $j = 0$ **to** n **do**
 if $j \neq i$ **then**
 $\mathbb{L}_i(x_{\text{eval}}) = \mathbb{L}_i(x_{\text{eval}}) \times \frac{x_{\text{eval}} - x_j}{x_i - x_j}$
 end if
 end for
end for

Step 2: Compute Lagrange Polynomial $p_n^{\mathbb{L}}(x)$ at x_{eval}
Initialize $p_n^{\mathbb{L}}(x_{\text{eval}}) = 0$
for $i = 0$ **to** n **do**
 $p_n^{\mathbb{L}}(x_{\text{eval}}) = p_n^{\mathbb{L}}(x_{\text{eval}}) + f_i \times \mathbb{L}_i(x_{\text{eval}})$
end for

return $p_n^{\mathbb{L}}(x_{\text{eval}})$

complexity. An individual Lagrange interpolating polynomial of degree n looks like

$$\mathbb{L}_i(x) = \prod_{j=0, j \neq i}^n \frac{(x - x_j)}{(x_i - x_j)},$$

and then $p_n^{\mathbb{L}}(x) = f_0 + \mathbb{L}_0(x) + f_1 \mathbb{L}_1(x) + \dots + f_n \mathbb{L}_n(x)$. For the computation of each $\mathbb{L}_i(x)$ we need $\mathcal{O}(n)$ multiplications. As we have $(n + 1)$ points, we need $\mathcal{O}(n^2 + n)$ operations. The final operation for computing of $p_n^{\mathbb{L}}(x)$ is of multiplication and addition and hence a total of $\mathcal{O}(n)$ operations. Therefore, in totality, we need $\mathcal{O}(n^2)$ operations, which is not very nice as, generally, we prefer to have linear ($\mathcal{O}(n)$) complexity.

Apart from the above drawback, another major drawback is that if we want to add a new point, say (x_{n+1}, f_{n+1}) , then we need to perform new computations from scratch.

But there are advantages as well; for example, the computation of $\{\mathbb{L}_i(x)\}_{i=0}^n$ is independent of $f(x_k)$. Another one is that it does not depend on the arrangement of nodes.

1.2.2 Runge Phenomena

In 1901, Carl David Tolmé Runge observed that while approximating

$$f(x) = \frac{1}{1 + 25x^2}, \quad x \in [-1, 1],$$

using polynomial approximation, there are large errors at the endpoints of the interval while using equally spaced points (see Fig. 1.11). This is what is called as the *Runge phenomena* and the above function is called the *Runge function*.



Figure 1.10: Carl David Tolmé Runge: 30 August 1856-3 January 1927

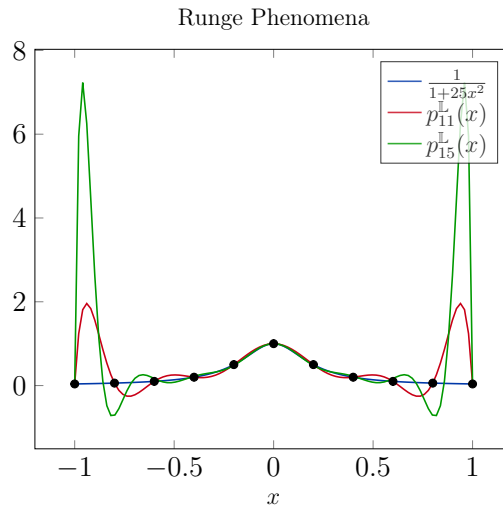


Figure 1.11: Runge phenomena for the function $1/(1+25x^2)$. $p_{11}^{\mathbb{L}}(x)$ refers to an approximation computed using 11 points (the dots refer to $\{x_i\}_{i=0}^{10}$), $p_{15}^{\mathbb{L}}(x)$ refers to approximation using 15 points.

Let us look at the interpolation error and try to understand this phenomenon. In Theorem 1.3 it was seen that

$$f(x) - p_n^{\mathbb{L}}(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \prod_{i=0}^n (x - x_i) \quad \text{for } \xi \in (-1, 1).$$

Thus, we have

$$\max_{-1 \leq x \leq 1} |f(x) - p_n^{\mathbb{L}}(x)| \leq \max_{-1 \leq x \leq 1} \left| \frac{f^{(n+1)}(\xi)}{(n+1)!} \right| \max_{-1 \leq x \leq 1} \prod_{i=0}^n |x - x_i|.$$

Now, it can be shown (although not very easily) that $\max_{-1 \leq x \leq 1} \prod_{i=0}^n |x - x_i| \leq h^{n+1} n!$ where $h = 2/n$ and we suppose that the $(n+1)^{\text{th}}$ derivative of $f(x)$ can be bounded by M_{n+1} which in turn can be bounded by $5^{n+1}(n+1)!$ (see this PDF). Hence in total

$$\lim_{n \rightarrow \infty} \left(\max_{-1 \leq x \leq 1} |f(x) - p_n^{\mathbb{L}}(x)| \right) \leq \lim_{n \rightarrow \infty} \left(\left(\frac{10}{n} \right)^{n+1} n! \right) = \infty.$$

To mitigate this problem, one idea is to use a non-uniform grid with points accumulated at the endpoints. If one is interested, I suggest this excellent review paper by Berrut and Trefethen [2].

1.3 Newton Divided Difference Interpolation

We noticed in Sec. 1.2 that Lagrange interpolation suffers from $\mathcal{O}(n^2)$ evaluation computational complexity. Now, we have another interpolating method that overcomes this and is referred to as the *Newton Divided Differences*.

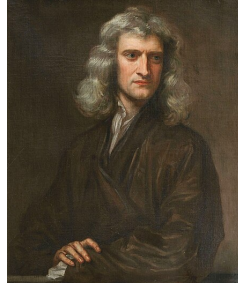


Figure 1.12: Isaac Newton: 4 January 1643-31 March 1727

Let $p_n(x)$ be a polynomial interpolating the data points $\{(x_i, f_i)\}_{i=0}^n$. Another way of expressing such a polynomial is

$$p_n(x) = a_0 + a_1(x - x_0) + a_2(x - x_0)(x - x_1) + \cdots + a_n \prod_{i=0}^{n-1} (x - x_i), \quad (1.6)$$

for appropriate constants $\{a_i\}_{i=0}^n$. Now, the question is, how do we determine these coefficients? At $x = x_0$ we have $p_n(x_0) = f_0$. Hence, $y_0 = a_0$. Similarly at $x = x_1$, $p_n(x_1) = f_1$ which implies

$$a_1 = \frac{f_1 - a_0}{x_1 - x_0} = \frac{f_1 - f_0}{x_1 - x_0}.$$

Now, we can continue in this manner and compute each a_i . For this, we introduce the divided difference (DD) notation. The zeroth divided difference of a function $f(x)$ with respect to x_i is denoted by $f[x_i] = f(x_i) = f_i$. For the rest, we define them in a recursive way.

- 1st DD of $f(x)$ with respect to x_i and x_{i+1} is

$$f[x_i, x_{i+1}] = \frac{f[x_{i+1}] - f[x_i]}{x_{i+1} - x_i}. \quad (1.7)$$

- 2nd DD of $f(x)$ with respect to x_i , x_{i+1} and x_{i+2} is

$$f[x_i, x_{i+1}, x_{i+2}] = \frac{f[x_{i+1}, x_{i+2}] - f[x_i, x_{i+1}]}{x_{i+2} - x_i}. \quad (1.8)$$

- k^{th} DD of $f(x)$ with respect to x_i , x_{i+1}, \dots, x_{i+k} is

$$f[x_i, x_{i+1}, \dots, x_{i+k-1}, x_{i+k}] = \frac{f[x_{i+1}, \dots, x_{i+k}] - f[x_i, \dots, x_{i+k-1}]}{x_{i+k} - x_i}. \quad (1.9)$$

- n^{th} DD of $f(x)$ with respect to x_0, x_1, \dots, x_n is

$$f[x_0, x_1, \dots, x_n] = \frac{f[x_1, \dots, x_n] - f[x_0, \dots, x_{n-1}]}{x_n - x_0}. \quad (1.10)$$

Hence, we can rewrite the polynomial $p_n(x)$ defined in Eq. (1.6) as

$$p_n^{\mathbb{N}}(x) = f[x_0] + \sum_{k=1}^n f[x_0, x_1, \dots, x_k](x - x_0)(x - x_1) \dots (x - x_{k-1}). \quad (1.11)$$

We have introduced the notation $p_n^{\mathbb{N}}(x)$ to identify the Newton DD polynomial.

For simplicity, let us look at the DD table we obtain for 4 points (see Table 1.1)

x	$f(x) = 0^{\text{th}}$ DD	1 st DD	2 nd DD	3 rd DD
x_0	$f[x_0]$			
x_1	$f[x_1]$	$f[x_0, x_1] = \frac{f[x_1]-f[x_0]}{x_1-x_0}$	$f[x_0, x_1, x_2] = \frac{f[x_1, x_2]-f[x_0, x_1]}{x_2-x_0}$	
x_2	$f[x_2]$	$f[x_1, x_2] = \frac{f[x_2]-f[x_1]}{x_2-x_1}$	$f[x_1, x_2, x_3] = \frac{f[x_2, x_3]-f[x_1, x_2]}{x_3-x_1}$	$f[x_0, x_1, x_2, x_3] = \frac{f[x_1, x_2, x_3]-f[x_0, x_1, x_2]}{x_3-x_0}$
x_3	$f[x_3]$	$f[x_2, x_3] = \frac{f[x_3]-f[x_2]}{x_3-x_2}$		

Table 1.1: Divided difference table for four points x_0, x_1, x_2, x_3 .

The Lagrange interpolating polynomial has a polynomial basis as $\mathbb{L}_{n,i}(x)$, we can consider the Newton DD as another method with a basis defined by $\omega_i(x) = \prod_{k=0}^{i-1} (x - x_k)$ for $i \geq 1$ and $\omega_0(x) = 1$ (see Fig. 1.13).

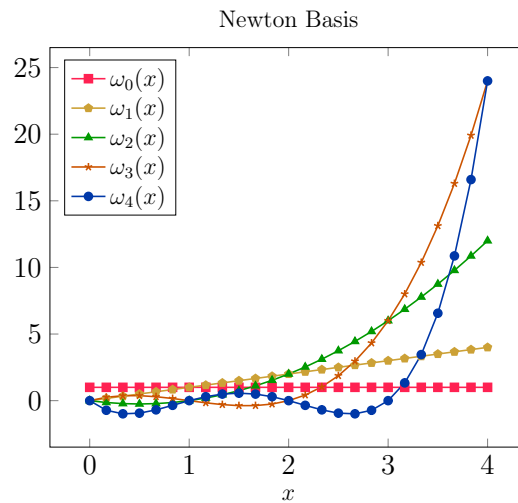


Figure 1.13: Newton basis polynomials defined over $x_i = 0, 1, 2, 3, 4$.

Now, we try to establish a relation between the DD and the derivatives of f . First, we recall the mean value theorem

Theorem 1.5. [1, Theorem 6.2.4] **(Mean Value Theorem)** If $f \in \mathcal{C}[a, b]$ and $f(x)$ is differentiable in (a, b) then there exists a $c \in (a, b)$ such that

$$f'(c) = \frac{f(b) - f(a)}{b - a}.$$

Now, if we apply the MVT on the interval $[x_i, x_{i+1}]$ then there exists a $\xi \in (x_i, x_{i+1})$ such that

$$f'(\xi) = \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i} = f[x_i, x_{i+1}].$$

In fact, we can generalize this concept.

Theorem 1.6. *Suppose that $f \in \mathcal{C}^n[a, b]$ and x_0, x_1, \dots, x_n are distinct numbers in $[a, b]$. Then there exists a $\xi \in (a, b)$ with*

$$f[x_0, x_1, \dots, x_n] = \frac{f^{(n)}(\xi)}{n!}.$$

Proof. Let $g(x) = f(x) - p_n^{\mathbb{N}}(x)$. Since, $f(x_i) = p_n^{\mathbb{N}}(x_i)$ at $i = 0, 1, \dots, n$. Then $g(x)$ has $(n+1)$ distinct zeros in $[a, b]$. So by generalized Rolle's theorem 1.2 there exists a $\xi \in (a, b)$ with $g^{(n)}(\xi) = 0$, so

$$0 = f^{(n)}(\xi) - p_n^{\mathbb{N}(n)}(\xi).$$

Since, $p_n^{\mathbb{N}}(x)$ is polynomial of degree n with leading coefficient $f[x_0, x_1, \dots, x_n]$, we have

$$p_n^{\mathbb{N}(n)}(x) = n!f[x_0, x_1, \dots, x_n] \quad \text{for all } x.$$

Hence,

$$f[x_0, x_1, \dots, x_n] = \frac{f^{(n)}(\xi)}{n!}.$$

□

Next, we give a result which gives an explicit representation of the Newton DD formula.

Theorem 1.7. *For distinct points x_0, \dots, x_n , the n^{th} coefficient of the Newton interpolation satisfies*

$$f[x_0, x_1, \dots, x_n] = \sum_{k=0}^n \frac{f(x_k)}{\prod_{i \neq k} (x_k - x_i)},$$

where $f[x_0] = f(x_0)$ in the case $n = 0$.

Proof. Using the representation Eq. (1.11) of the interpolant the n^{th} derivative of $p_n^{\mathbb{N}}(x)$ is given by to $f[x_0, \dots, x_n]\omega_n^{(n)}(x)$ where $\omega_n(x) = \prod_{i=0}^{n-1} (x - x_i)$.

Now, the polynomial $p_n^{\mathbb{N}}(x)$ is just another representation of $p_n^{\mathbb{L}}(x)$. Hence their n^{th} derivatives must match.

For the Lagrange polynomial the n^{th} derivative is $\sum_{k=0}^n f(x_k)\mathbb{L}_k^{(n)}(x)$ (see Eq. (1.3)). Hence,

$$f[x_0, x_1, \dots, x_n]\omega_n^{(n)}(x) = \sum_{k=0}^n f(x_k)\mathbb{L}_k^{(n)}(x).$$

Now, the k^{th} Lagrange interpolating polynomial is given by

$$\mathbb{L}_k(x) = \prod_{j \neq k} \left(\frac{x - x_j}{x_k - x_j} \right),$$

and its n^{th} derivative $\mathbb{L}_k^{(n)}(x) = \frac{n!}{\prod_{j \neq k} (x_k - x_j)}$ since x^n is the leading term. Since the leading term in $\omega_n(x)$ is x^n , we get $\omega_n^{(n)}(x) = n!$. Cancelling out these factorial term we get the expression. \square

The algorithm for computing the Newton DD interpolation is provided in Algorithm 3.

Algorithm 3 Newton Interpolation

Given: Data sets $\{(x_i, f_i)\}_{i=0}^n$, Evaluation point x_{eval} .

Find: Interpolated polynomial $p_n^{\mathbb{N}}(x_{\text{eval}})$.

Step 1: Construct Divided Difference Table

Initialize DD as a zero matrix of size $(n + 1) \times (n + 1)$.

for $i = 0$ **to** n **do**

$DD_{i,0} = f_i$

end for

for $j = 1$ **to** n **do**

for $i = 0$ **to** $n - j$ **do**

Compute:

$$DD_{i,j} = \frac{DD_{i+1,j-1} - DD_{i,j-1}}{x_{i+j} - x_i}.$$

end for

end for

Step 2: Evaluate Newton Polynomial $p_n^{\mathbb{N}}(x)$ at x_{eval}

Initialize $p_n^{\mathbb{N}}(x_{\text{eval}}) = DD_{0,0}$.

for $k = 1$ **to** n **do**

Initialize $\omega = 1$.

for $j = 0$ **to** $k - 1$ **do**

$\omega = \omega \times (x_{\text{eval}} - x_j)$

end for

Update the interpolated value:

$$p_n^{\mathbb{N}}(x_{\text{eval}}) = p_n^{\mathbb{N}}(x_{\text{eval}}) + DD_{0,k} \cdot \omega.$$

end for

return $p_n^{\mathbb{N}}(x_{\text{eval}})$.

1.3.1 Computational Complexity

We can rewrite the Newton interpolation as

$$\begin{aligned} p_n^{\mathbb{N}}(x) &= a_0 + a_1(x - x_0) + a_2(x - x_0)(x - x_1) + \cdots + a_n(x - x_0)(x - x_1) \cdots (x - x_{n-1}) \\ &= a_0 + (x - x_0) [a_1 + (x - x_1) \{a_2 + \cdots + (x - x_{n-2}) \{a_{n-1} + a_n(x - x_{n-1})\}\}]. \end{aligned}$$

We notice that each term requires one multiplication and one addition for evaluation, and we have n points. Hence, we require $2n$ operations, which is of $\mathcal{O}(n)$, whereas for Lagrange, we have $\mathcal{O}(n^2)$.

Another advantage of the Newton interpolation over Lagrange interpolation is that it is easy to update the DD table whenever we have a new data set as it does not require new computation only a modification of the DD table.

1.3.2 Forward Difference Formula

Suppose we have an equal spacing of points; then we can rewrite Newton's formula in a better way. Let $h = x_{i+1} - x_i$ for all $i = 0, 1, \dots, n-1$ and $x = x_0 + sh$. Then we can rewrite Eq. (1.11) as

$$\begin{aligned} p_n^{\mathbb{N}}(x) &= f[x_0] + (x - x_0)f[x_0, x_1] + (x - x_0)(x - x_1)f[x_0, x_1, x_2] + \cdots \\ &\quad + (x - x_0) \cdots (x - x_{n-1})f[x_0, x_1, \dots, x_n] \\ &= f[x_0] + shf[x_0, x_1] + s(s-1)h^2f[x_0, x_1, x_2] + \cdots \\ &\quad + s(s-1) \cdots (s-n+1)h^nf[x_0, x_1, \dots, x_n] \\ &= f[x_0] + \sum_{k=1}^n s(s-1) \cdots (s-k+1)h^k f[x_0, x_1, \dots, x_k]. \end{aligned}$$

Using the binomial coefficient notation

$${}^sC_k = \frac{s(s-1) \cdots (s-k+1)}{k!},$$

we can express

$$p_n^{\mathbb{N}}(x) = p_n^{\mathbb{N}}(x_0 + sh) = f[x_0] + \sum_{k=1}^n {}^sC_k k! h^k f[x_0, x_1, \dots, x_k].$$

Let us use the Δ notation for forward difference, i.e, $\Delta f(x_0) = f(x_1) - f(x_0)$. Similarly for higher differences we use the notation $\Delta^2 f(x_0) = \Delta f(x_1) - \Delta f(x_0)$, then we can rewrite the divided differences as

$$\begin{aligned} f[x_0, x_1] &= \frac{f[x_1] - f[x_0]}{x_1 - x_0} = \frac{\Delta f(x_0)}{h} \\ f[x_0, x_1, x_2] &= \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_0} = \frac{1}{2h} \frac{\Delta f(x_1) - \Delta f(x_0)}{h} = \frac{\Delta^2 f(x_0)}{2!h^2} \\ &\vdots \\ f[x_0, x_1, \dots, x_n] &= \frac{\Delta^n f(x_0)}{n!h^n}. \end{aligned}$$

Hence,

$$p_n^{\mathbb{N}}(x) = f(x_0) + \sum_{k=1}^n {}^sC_k \Delta^k f(x_0).$$

1.4 Hermite Interpolation

Definition 1.8. Let $\{x_0, x_1, \dots, x_n\}$ be $(n + 1)$ distinct points in $[a, b]$ and for $i = 0, 1, \dots, n$ let m_i be a non-negative integer. Suppose that $f \in \mathcal{C}^m[a, b]$ where $m = \max_{0 \leq i \leq n} m_i$. The *osculating polynomial* approximating $f(x)$ is the polynomial $p(x)$ of least degree such that

$$\frac{d^k p(x_i)}{dx^k} = \frac{d^k f(x_i)}{dx^k}, \quad \text{for } i = 0, 1, \dots, n \quad \text{and } k = 0, 1, \dots, m_i.$$

Not when $n = 0$ the osculating polynomial approximating f is the m_0^{th} Taylor polynomial for f at x_0 . When $m_i = 0$ for all i then the osculating polynomial is the n^{th} Lagrange polynomial interpolating f at x_0, x_1, \dots, x_n .

Hermite Polynomials

If $m_i = 1$ for all $i = 0, 1, \dots, n$ then we get the Hermite polynomials. For a given function f these polynomials agree with f at x_0, x_1, \dots, x_n . In addition they agree with their derivatives as well.



Figure 1.14: Charles Hermite: 24 December 1822- 14 January 1901

Theorem 1.9. If $f \in \mathcal{C}^1[a, b]$ and $x_0, x_1, \dots, x_n \in [a, b]$ are distinct, the unique polynomial of least degree agreeing with f and f' at x_0, x_1, \dots, x_n is the Hermite polynomial of degree at most $2n + 1$ given by

$$p_{2n+1}^{\mathbb{H}}(x) = \sum_{j=0}^n f(x_j) \mathbb{H}_{n,j}(x) + \sum_{j=0}^n f'(x_j) \hat{\mathbb{H}}_{n,j}(x),$$

where for $\mathbb{L}_{n,j}(x)$ denoting the j^{th} Lagrange coefficient polynomial of degree n we have

$$\mathbb{H}_{n,j}(x) = [1 - 2(x - x_j)\mathbb{L}_{n,j}'(x)] \mathbb{L}_{n,j}^2(x) \quad \text{and} \quad \hat{\mathbb{H}}_{n,j}(x) = (x - x_j)\mathbb{L}_{n,j}^2(x).$$

Moreover, if $f \in \mathcal{C}^{2n+2}[a, b]$ then

$$f(x) = p_{2n+1}^{\mathbb{H}}(x) + \frac{(x - x_0)^2(x - x_1)^2 \dots (x - x_n)^2}{(2n + 2)!} f^{(2n+2)}(\xi(x)),$$

for some unknown $\xi(x) \in (a, b)$.

Proof. We know that $\mathbb{L}_{n,j}(x_i) = \delta_{ij}$. Hence when $i \neq j$ $\mathbb{H}_{n,j}(x_i) = 0$ and $\hat{\mathbb{H}}_{n,j}(x_i) = 0$ whereas for each i

$$\mathbb{H}_{n,i}(x_i) = [1 - 2(x_i - x_i)\mathbb{L}_{n,i}'(x_i)] \mathbb{L}_{n,i}^2(x_i) = 1 \quad \text{and} \quad \hat{\mathbb{H}}_{n,i}(x_i) = (x_i - x_i)\mathbb{L}_{n,i}^2(x_i) = 0.$$

Hence, we can say $\mathbb{H}_{n,j}(x_i) = \delta_{ij}$ and $\hat{\mathbb{H}}_{n,j}(x_i) = 0$ for all i, j . As a consequence

$$p_{2n+1}^{\mathbb{H}}(x_i) = \sum_{j=0}^n f(x_j) \mathbb{H}_{n,j}(x_i) + \sum_{j=0}^n f'(x_j) \hat{\mathbb{H}}_{n,j}(x_i) = f(x_i),$$

so $p_{2n+1}^{\mathbb{H}}$ agrees with f at x_0, x_1, \dots, x_n .

Now we need to show that they match at the derivatives as well, i.e., $p_{2n+1}^{\mathbb{H}'}$ and f' match at x_i . We will tackle this by differentiating both the terms $\mathbb{H}_{n,j}$ and $\hat{\mathbb{H}}_{n,j}$.

The derivative of $\mathbb{H}_{n,j}(x)$ is given by

$$\begin{aligned} \mathbb{H}'_{n,j}(x) &= [1 - 2(x - x_j)\mathbb{L}'_{n,j}(x_j)] 2\mathbb{L}_{n,j}(x)\mathbb{L}'_{n,j}(x) + \mathbb{L}_{n,j}^2(x) [-2\mathbb{L}'_{n,j}(x_j)] \\ &= 2\mathbb{L}_{n,j}(x) [\{1 - 2(x - x_j)\mathbb{L}'_{n,j}(x_j)\} \mathbb{L}'_{n,j}(x) - \mathbb{L}_{n,j}(x)\mathbb{L}'_{n,j}(x_j)]. \end{aligned}$$

As $\mathbb{L}_{n,j}(x_i) = \delta_{ij}$ we get that at $i \neq j$, $\mathbb{H}'_{n,j}(x_i) = 0$. At $i = j$ we have

$$\begin{aligned} \mathbb{H}'_{n,i}(x_i) &= [1 - 2(x_i - x_i)\mathbb{L}'_{n,i}(x_i)] 2\mathbb{L}_{n,i}(x_i)\mathbb{L}'_{n,i}(x_i) + \mathbb{L}_{n,i}^2(x_i) [-2\mathbb{L}'_{n,i}(x_i)] \\ &= [2\mathbb{L}'_{n,i}(x_i) - 2\mathbb{L}'_{n,i}(x_i)] = 0. \end{aligned}$$

Now for the second term the derivative is given by

$$\begin{aligned} \hat{\mathbb{H}}'_{n,j}(x) &= (x - x_j)2\mathbb{L}_{n,j}(x)\mathbb{L}'_{n,j}(x) + \mathbb{L}_{n,j}^2(x) \\ &= \mathbb{L}_{n,j}(x) [2(x - x_j)\mathbb{L}'_{n,j}(x) + \mathbb{L}_{n,j}(x)]. \end{aligned}$$

At $x = x_i$ we have $\mathbb{L}_{n,j}(x_i) = \delta_{ij}$. Hence $\hat{\mathbb{H}}'_{n,j}(x_i) = 0$ if $i \neq j$ and at $i = j$

$$\begin{aligned}\hat{\mathbb{H}}'_{n,i}(x_i) &= (x_i - x_i)2\mathbb{L}_{n,i}(x_i)\mathbb{L}'_{n,i}(x_i) + \mathbb{L}_{n,i}^2(x_i) \\ &= 1.\end{aligned}$$

Hence $\hat{\mathbb{H}}'_{n,j}(x_i) = \delta_{ij}$. Therefore

$$p_{2n+1}^{\mathbb{H}'}(x_i) = \sum_{j=0}^n f(x_j)\hat{\mathbb{H}}'_{n,j}(x_i) + \sum_{j=0}^n f'(x_j)\hat{\mathbb{H}}'_{n,j}(x_i) = f'(x_i).$$

Therefore $p_{2n+1}^{\mathbb{H}}$ agrees with f and $p_{2n+1}^{\mathbb{H}'}$ with f' at x_0, x_1, \dots, x_n . So, we have existence of a polynomial that agrees with f and f' at $\{x_i\}_{i=0}^n$.

For the uniqueness we will use the method of contradiction. Suppose there exists another polynomial of least degree say $q(x)$ such that

$$q(x_i) = f(x_i) \quad \text{and} \quad q'(x_i) = f'(x_i) \quad \forall i.$$

Now consider the polynomial $D(x) = p_{2n+1}^{\mathbb{H}}(x) - q(x)$ of degree at most $(2n+1)$. Obviously

$$D(x_i) = 0 \quad \text{and} \quad D'(x_i) = 0 \quad \forall i$$

Hence x_i are distinct roots of multiplicity two. Therefore we have $2n+2$ roots, which is only possible if $D(x) = 0$. Hence, we get $p_{2n+1}^{\mathbb{H}}(x) = q(x)$ leading to a contradiction.

For showing the error term we will use the same strategy as in theorem 1.3., if $x = x_i$ for some i then we can choose $\xi(x)$ arbitrary.

Suppose $x \neq x_i$ for any i , then define

$$g(t) = f(t) - p_{2n+1}^{\mathbb{H}}(t) - [f(x) - p_{2n+1}^{\mathbb{H}}(x)] \prod_{i=0}^n \frac{(t - x_i)^2}{(x - x_i)^2}.$$

Now $g(x) = 0$ and $g(x_i) = 0$ for all i . Hence $g(t)$ has distinct $n+2$ roots in $[a, b]$. Hence, by Rolle's theorem $g'(t)$ has $n+1$ distinct roots between x_0, x_1, \dots, x_n and x , say $\xi_0, \xi_1, \dots, \xi_n$.

Now taking the derivative of $g(t)$ with respect to t we get

$$\begin{aligned}g'(t) &= f'(t) - p_{2n+1}^{\mathbb{H}'}(t) - \frac{[f(x) - p_{2n+1}^{\mathbb{H}}(x)]}{\prod_{i=0}^n (x - x_i)^2} \frac{d}{dt} \left(\prod_{i=0}^n (t - x_i)^2 \right) \\ &= f'(t) - p_{2n+1}^{\mathbb{H}'}(t) - \frac{[f(x) - p_{2n+1}^{\mathbb{H}}(x)]}{\prod_{i=0}^n (x - x_i)^2} \frac{d}{dt} ((t - x_0)^2 (t - x_1)^2 \dots (t - x_n)^2) \\ &= f'(t) - p_{2n+1}^{\mathbb{H}'}(t) - \frac{[f(x) - p_{2n+1}^{\mathbb{H}}(x)]}{\prod_{i=0}^n (x - x_i)^2} \left(2(t - x_0)(t - x_1)^2 \dots (t - x_n)^2 \right. \\ &\quad \left. + (t - x_0)^2 2(t - x_1) \dots (t - x_n)^2 + \dots + (t - x_0)^2 (t - x_1)^2 \dots 2(t - x_n) \right) \\ &= f'(t) - p_{2n+1}^{\mathbb{H}'}(t) - 2 \frac{[f(x) - p_{2n+1}^{\mathbb{H}}(x)]}{\prod_{i=0}^n (x - x_i)^2} \sum_{k=0}^n (t - x_k) \prod_{j=0, j \neq k}^n (t - x_j)^2\end{aligned}$$

At $t = x_i$ for any i we have $g'(x_i) = 0$ for $i = 0, 1, \dots, n$. Hence, $g'(t)$ has $2n+2$ roots. Using the generalized Rolle's theorem on $g'(t)$ and then following the same pattern as in Theorem 1.3 we get the result. \square

Theorem 1.9 gives all the details about the Hermite polynomials but it is computationally expensive as we need to compute the Lagrange polynomials and its derivatives.

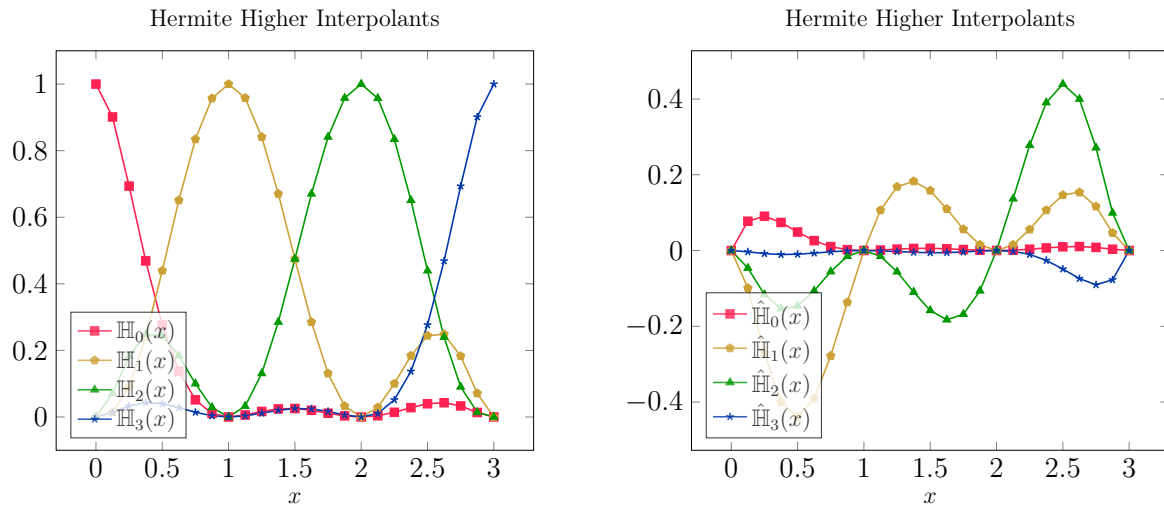


Figure 1.15: Hermite interpolating polynomials defined over $x_i = 0, 1, 2, 3$.

1.4.1 Hermite Polynomials using Divided Difference

For computing the Hermite polynomials using the Newton DD. We will use the relation between the n^{th} DD and the n^{th} derivative of $f(x)$ as in Theorem 1.6.

Suppose we have $(n + 1)$ distinct points $\{x_i\}_{i=0}^n$, we define a new sequence $\{z_i\}_{i=0}^{2n+1}$ by

$$z_{2i} = z_{2i+1} = x_i \quad \text{for each } i = 0, 1, \dots, n,$$

i.e., $z_0 = z_1 = x_0$, $z_2 = z_3 = x_1$, and so on. Then we can construct the DD table using these values. Since $z_{2i} = z_{2i+1}$ we cannot define $f[z_{2i}, z_{2i+1}]$. However from Theorem 1.6 we can make a reasonable substitution that

$$f[z_{2i}, z_{2i+1}] = f'(z_{2i}) = f'(x_i).$$

Hence we can use the derivative entries for the undefined DD.

The remaining entries of the DD are defined in the same manner and we get the Hermite polynomial as

$$p_{2n+1}^{\mathbb{H}}(x) = f[z_0] + \sum_{k=1}^{2n+1} f[z_0, z_1, \dots, z_k](x - z_0)(x - z_1) \dots (x - z_{k-1}).$$

For an example let us consider a data set of two points x_0 and x_1 . Then the DD table is given by table 1.2.

The algorithm for Hermite interpolation is given in algorithm 4.

z	$f(z)$	1 st DD	2 nd DD	3 rd DD
$z_0 = x_0$	$f[z_0] = f(x_0)$	$f[z_0, z_1] = f'(x_0)$		
$z_1 = x_0$	$f[z_1] = f(x_0)$	$f[z_1, z_2] = \frac{f[z_2] - f[z_1]}{z_2 - z_1}$	$f[z_0, z_1, z_2] = \frac{f[z_1, z_2] - f[z_0, z_1]}{z_2 - z_0}$	$f[z_0, z_1, z_2, z_3] = \frac{f[z_1, z_2, z_3] - f[z_0, z_1, z_2]}{z_3 - z_0}$
$z_2 = x_1$	$f[z_2] = f(x_1)$	$f[z_2, z_3] = f'(x_1)$	$f[z_1, z_2, z_3] = \frac{f[z_2, z_3] - f[z_1, z_2]}{z_3 - z_1}$	
$z_3 = x_1$	$f[z_3] = f(x_1)$			

Table 1.2: Divided difference table for two points and the Hermite polynomial

Algorithm 4 Hermite Interpolation

Given: Data sets $\{(x_i, f_i, f'_i)\}_{i=0}^n$, Evaluation point x_{eval} .

Find: Interpolated polynomial $p_{2n+1}^{\mathbb{H}}(x_{\text{eval}})$.

Step 1: Create z_i and $f(z_i)$ arrays

Construct $\{z_i\}_{i=0}^{2n+1}$ and $\{f(z_i)\}_{i=0}^{2n+1}$

for $i = 0$ **to** n **do**

$$z_{2i} = z_{2i+1} = x_i, \quad f(z_{2i}) = f(z_{2i+1}) = f_i$$

end for

Step 2: Construct Divided Difference Table

Initialize DD as a zero matrix of size $(2n + 2) \times (2n + 2)$.

for $i = 0$ **to** $2n + 1$ **do**

$$\text{DD}_{i,0} = f(z_i)$$

end for

for $j = 1$ **to** $2n + 1$ **do**

for $i = 0$ **to** $2n + 1 - j$ **do**

if $j = 1$ **and** $i \% 2 = 0$ **then**

$$\text{DD}_{i,j} = f'_i$$

else

Compute:

$$\text{DD}_{i,j} = \frac{\text{DD}_{i+1,j-1} - \text{DD}_{i,j-1}}{z_{i+j} - z_i}.$$

end if

end for

end for

Step 3: Evaluate Hermite Polynomial $p_{2n+1}^{\mathbb{H}}(x)$ at x_{eval}

Initialize $p_{2n+1}^{\mathbb{H}}(x) = \text{DD}_{0,0}$.

for $k = 1$ **to** $2n + 1$ **do**

Initialize $\omega = 1$.

for $j = 0$ **to** $k - 1$ **do**

$$\omega = \omega \times (x_{\text{eval}} - z_j)$$

end for

Update the interpolated value:

$$p_{2n+1}^{\mathbb{H}}(x_{\text{eval}}) = p_{2n+1}^{\mathbb{H}}(x_{\text{eval}}) + \text{DD}_{0,k} \cdot \omega.$$

end for

return $p_{2n+1}^{\mathbb{H}}(x_{\text{eval}})$.

1.5 Spline Interpolation

Both Lagrange and Newton interpolation methods suffer from the Runge phenomenon, where oscillations occur at the edges of the interval, especially with high-degree polynomials. This issue arises because these methods rely on a single global polynomial, meaning that every data point influences the entire approximation. This “global approximation” can lead to poor performance for non-uniform or large datasets.

An alternative approach is to divide the interval into smaller sub-intervals and use piecewise polynomial approximation. This strategy, known as local interpolation, reduces the influence of distant data points, resulting in more stable and accurate approximations.

Given a set of points $\{(x_i, f_i)\}_{i=0}^n$ we can use piecewise-linear interpolation that consists of joining set of data points using straight lines (see Fig. 1.16). An immediate disadvantage of such an interpolation is that the approximating polynomial is not differentiable at the nodal points which geometrically mean the function is not “smooth”.

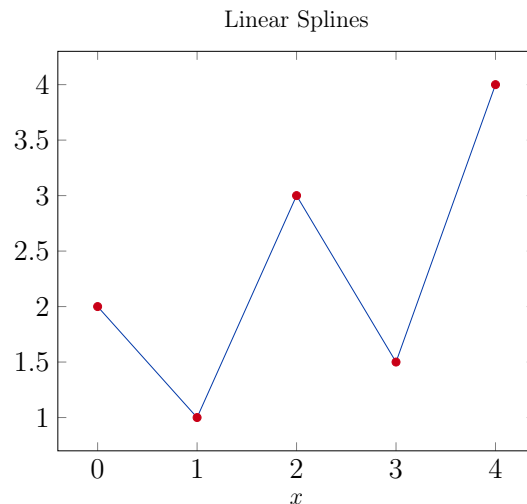


Figure 1.16: Linear spline defined over $x_i = 0, 1, 2, 3, 4$.

To address the limitations of linear interpolation, we use *splines*, which are piecewise polynomials of higher degree. The term “spline” was introduced by Isaac Jacob Schoenberg in the 1930s, inspired by drafting tools called “flat splines”. These tools were used to draw smooth curves on paper before the advent of computer-aided design. A spline curve behaves like a flexible beam, ensuring continuity in both slope and curvature.

1.5.1 Cubic Splines

The most common piecewise polynomial approximation uses the cubic polynomials between each successive pair of nodes and is called *cubic spline interpolation* (see Fig. 1.18) .



Figure 1.17: Isaac Jacob Schoenberg: 21 April 1903-21 February 1990

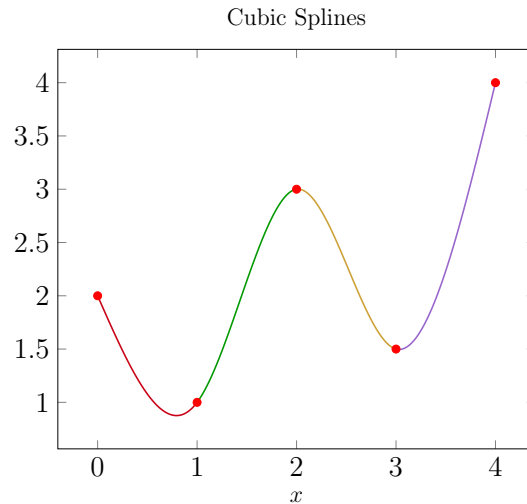


Figure 1.18: Cubic spline defined over $x_i = 0, 1, 2, 3, 4$.

Definition 1.10. Given a function f defined on $[a, b]$ and a set of nodes $a = x_0 < x_1 < \dots < x_n = b$ (called *knots*), a cubic spline interpolant $p^{\mathbb{S}}$ for f is a function that satisfies the following conditions:

- a) $p^{\mathbb{S}}(x)$ is a cubic polynomial, whose restriction on the interval $[x_j, x_{j+1}]$ is denoted by $p_j^{\mathbb{S}}(x)$ for each $j = 0, 1, \dots, n - 1$.
- b) $p_j^{\mathbb{S}}(x_j) = f(x_j)$ and $p_j^{\mathbb{S}}(x_{j+1}) = f(x_{j+1})$ for $j = 0, 1, \dots, n - 1$.
- c) $p_{j+1}^{\mathbb{S}}(x_{j+1}) = p_j^{\mathbb{S}}(x_{j+1})$ for $j = 0, 1, \dots, n - 2$ (implied by **b**).
- d) $p_{j+1}^{\mathbb{S}'}(x_{j+1}) = p_j^{\mathbb{S}'}(x_{j+1})$ for $j = 0, 1, \dots, n - 2$.
- e) $p_{j+1}^{\mathbb{S}''}(x_{j+1}) = p_j^{\mathbb{S}''}(x_{j+1})$ for $j = 0, 1, \dots, n - 2$.
- f) One of the following sets of boundary conditions is satisfied:
 - i) $p^{\mathbb{S}''}(x_0) = p^{\mathbb{S}''}(x_n) = 0$ (*natural* (or *free*) boundary).
 - ii) $p^{\mathbb{S}'}(x_0) = f'(x_0)$ and $p^{\mathbb{S}'}(x_n) = f'(x_n)$ (*clamped* boundary).

When the free boundary condition occurs the spline is called *natural spline*. In general clamped boundary conditions lead to more accurate results but it includes the information about the derivative of the function which is not easily available.

Notice that we have n intervals and on each interval we have 4 unknowns. Hence we have a total of $4n$ unknowns.

Construction of Cubic Splines

Let $[a, b]$ be an interval. We divide this interval into n subintervals denoted $[x_j, x_{j+1}]$ for any $j = 0, 1, 2, \dots, n-1$, then for each subinterval we define a cubic polynomial as

$$p_j^{\mathbb{S}}(x) = a_j + b_j(x - x_j) + c_j(x - x_j)^2 + d_j(x - x_j)^3.$$

Since, $p_j^{\mathbb{S}}(x) = f(x_j)$ we get $a_j = f(x_j)$.

Now from condition **c**) we have

$$\begin{aligned} p_{j+1}^{\mathbb{S}}(x_{j+1}) &= p_j^{\mathbb{S}}(x_{j+1}) \\ a_{j+1} &= a_j + b_j(x_{j+1} - x_j) + c_j(x_{j+1} - x_j)^2 + d_j(x_{j+1} - x_j)^3, \end{aligned}$$

for $j = 0, 1, 2, \dots, n-2$. Let us denote $x_{j+1} - x_j$ by h_j for $j = 1, 2, \dots, n-1$. If we define $a_n = f(x_n)$, then we get the relation

$$a_{j+1} = a_j + b_j h_j + c_j h_j^2 + d_j h_j^3, \quad \text{for } j = 0, 1, \dots, n-1. \quad (1.12)$$

We also note that

$$p_j^{\mathbb{S}'}(x) = b_j + 2c_j(x - x_j) + 3d_j(x - x_j)^2.$$

Substituting $x = x_j$, we get $p_j^{\mathbb{S}'}(x_j) = b_j$ for each $j = 0, 1, \dots, n-2$. Defining $b_n = p^{\mathbb{S}'}(x_n)$ we get the relation

$$b_{j+1} = b_j + 2c_j h_j + 3d_j h_j^2 \quad \text{for } j = 0, 1, \dots, n-1. \quad (1.13)$$

Now $p_j^{\mathbb{S}''}(x) = 2c_j + 6d_j(x - x_j)$ and hence $p_j^{\mathbb{S}''}(x_j) = 2c_j$. Defining $p^{\mathbb{S}''}(x_n) = 2c_n$, from condition **e**) we get

$$2c_{j+1} = 2c_j + 6d_j h_j \quad \text{for } j = 0, 1, \dots, n-1. \quad (1.14)$$

Solving for d_j in Eq. (1.14) we get $d_j = \frac{c_{j+1} - c_j}{3h_j}$ and substituting this back in Eq. (1.12) and Eq. (1.13) we get

$$\begin{aligned} a_{j+1} &= a_j + b_j h_j + c_j h_j^2 + \frac{(c_{j+1} - c_j)}{3h_j} h_j^3 \\ &= a_j + b_j h_j + \frac{(2c_j + c_{j+1})}{3} h_j^2. \end{aligned} \quad (1.15)$$

$$b_{j+1} = b_j + 2c_j h_j + h_j(c_{j+1} - c_j). \quad (1.16)$$

From Eq. (1.15) we get for b_j

$$b_j = \frac{a_{j+1} - a_j}{h_j} - \frac{h_j}{3}(2c_j + c_{j+1}). \quad (1.17)$$

Substituting Eq. (1.17) into Eq. (1.16) we get

$$\begin{aligned} \frac{a_{j+2} - a_{j+1}}{h_{j+1}} - \frac{h_{j+1}}{3}(2c_{j+1} + c_{j+2}) &= \frac{a_{j+1} - a_j}{h_j} - \frac{h_j}{3}(2c_j + c_{j+1}) + h_j(c_j + c_{j+1}) \\ \frac{a_{j+2} - a_{j+1}}{h_{j+1}} - \frac{a_{j+1} - a_j}{h_j} &= \frac{c_j h_j}{3} + \frac{2c_{j+1} h_j}{3} + \frac{2c_{j+1} h_{j+1}}{3} + \frac{c_{j+2} h_{j+1}}{3} \\ c_j h_j + 2c_{j+1}(h_j + h_{j+1}) + c_{j+2} h_{j+1} &= \frac{3(a_{j+2} - a_{j+1})}{h_{j+1}} - \frac{3(a_{j+1} - a_j)}{h_j}, \end{aligned}$$

for $j = 0, 1, 2, \dots, n - 2$.

For simplicity we do the shifting of the index by 1. Hence, finally we get the system of equation as

$$c_{j-1}h_{j-1} + 2c_j(h_{j-1} + h_j) + c_{j+1}h_j = \frac{3}{h_j}(a_{j+1} - a_j) - \frac{3}{h_{j-1}}(a_j - a_{j-1}), \quad (1.18)$$

for $j = 1, 2, \dots, n - 1$. The system of equations given by Eq. (1.18) involves the unknown $\{c_i\}_{i=0}^n$ as the values of $\{h_j\}_{j=0}^n$ and $\{a_j\}_{j=0}^n$ are known. Hence, if we can compute $\{c_j\}$ then we can compute $\{b_j\}_{j=0}^{n-1}$ using Eq. (1.17) and $\{d_j\}_{j=0}^{n-1}$ from Eq. (1.14). Hence after these computations we can compute $\{p_j^{\mathbb{S}}(x)\}_{j=0}^{n-1}$. So if Eq. (1.18) has an unique solution then we are done.

Theorem 1.11. *If f is defined at $a = x_0 < x_1 < \dots < x_n = b$, then f has a unique natural spline interpolant $p^{\mathbb{S}}(x)$ on the nodes x_0, x_1, \dots, x_n , i.e., a spline interpolant that satisfies the natural boundary condition $p^{\mathbb{S}''}(a) = p^{\mathbb{S}''}(b) = 0$.*

Proof. Let us consider $p_0^{\mathbb{S}}(x)$, which is given by $p_0^{\mathbb{S}}(x) = a_0 + b_0(x - x_0) + c_0(x - x_0)^2 + d_0(x - x_0)^3$. Now, $p_0^{\mathbb{S}''}(x) = 2c_0 + 6d_0(x - x_0)$ and at $x = x_0$, we have $p_0^{\mathbb{S}''}(x_0) = 2c_0 = 0$ which implies $c_0 = 0$. Similarly $c_n = 0$.

Let us look at Eq. (1.18)

$$h_{j-1}c_{j-1} + 2(h_{j-1} + h_j)c_j + h_jc_{j+1} = \frac{3}{h_j}(a_{j+1} - a_j) - \frac{3}{h_{j-1}}(a_j - a_{j-1}),$$

for $j = 1, 2, \dots, n - 1$. If we substitute for each j we get a system of equation

$$\mathbf{S}\mathbf{c} = \mathbf{v},$$

where

$$\mathbf{S} = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & \dots & 0 & 0 \\ h_0 & 2(h_0 + h_1) & h_1 & 0 & \dots & \dots & 0 & 0 \\ 0 & h_1 & 2(h_1 + h_2) & h_2 & \dots & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \ddots & \ddots & \vdots & \vdots \\ \vdots & \vdots & \ddots & \ddots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & \dots & \dots & h_{n-2} & 2(h_{n-2} + h_{n-1}) & h_{n-1} \\ 0 & 0 & \dots & \dots & \dots & 0 & 0 & 1 \end{bmatrix},$$

$$\mathbf{v} = \begin{bmatrix} 0 \\ \frac{3}{h_1}(a_2 - a_1) - \frac{3}{h_0}(a_1 - a_0) \\ \vdots \\ \frac{3}{h_{n-1}}(a_n - a_{n-1}) - \frac{3}{h_{n-2}}(a_{n-1} - a_{n-2}) \\ 0 \end{bmatrix}, \quad \text{and} \quad \mathbf{c} = \begin{bmatrix} c_0 \\ c_1 \\ \vdots \\ c_n \end{bmatrix}.$$

We use the following theorem to show that the matrix \mathbf{S} is invertible.

Theorem 1.12. (Strictly Diagonal Dominant Matrix)

A strictly diagonally dominant matrix \mathbf{A} is nonsingular.

We notice that our matrix \mathbf{S} is strictly diagonally dominant³ and hence it is invertible, which leads to a unique solution. \square

The algorithm for natural spline interpolation can be found in Algorithm 5 We have a similar result for the clamped spline interpolation.

Theorem 1.13. *If f is defined at $a = x_0 < x_1 < \dots < x_n = b$ and differentiable at a and b then f has a unique clamped spline interpolant $p^{\mathbb{S}}(x)$ on the nodes x_0, x_1, \dots, x_n , i.e., a spline interpolant that satisfies the clamped boundary condition $p^{\mathbb{S}'}(a) = f'(a)$ and $p^{\mathbb{S}'}(b) = f'(b)$.*

Now we present a result regarding the error bound of the spline interpolation but we will not delve into its proof as the proof requires a lot of technicalities from Numerical Analysis which is out of scope of this lecture.

Theorem 1.14. *Let $f \in C^4[a, b]$ with $M = \max_{a \leq x \leq b} |f^{(4)}(x)|$. If $p^{\mathbb{S}}(x)$ is the unique clamped cubic spline interpolant to f with respect to the nodes $a = x_0 < x_1 < \dots < x_n = b$ then for all $x \in [a, b]$*

$$|f(x) - p^{\mathbb{S}}(x)| \leq \frac{5M}{384} \max_{0 \leq j \leq n-1} (x_{j+1} - x_j)^4.$$

A fourth order error bound also exist for the case of natural boundary splines, but they are more difficult to express. An alternative to the natural boundary condition is the *not-a-knot* condition, it states that $p^{\mathbb{S}''''}(x)$ has to be continuous at x_1 and x_{n-1} .

1.5.2 B-Splines

So far, we have focused on a specific type of spline function called cubic splines. A natural question arises: can we generalize this to splines of other degrees? The answer to this is yes. A generalization of the cubic splines is the basis splines or B-splines.

Let $\{x_i\}_{i=0}^n$ be the data points (or knots); then we define the zeroth degree B-spline as

$$B_{j,0}(x) = \begin{cases} 1, & x \in [x_j, x_{j+1}), \\ 0, & \text{else,} \end{cases}$$

³**Strictly Diagonally Dominant Matrix:** A matrix $\mathbf{A} = \{a_{ij}\}_{i=1, j=1}^n$ is said to be strictly diagonally dominant if

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}| \quad \forall i.$$

Algorithm 5 Cubic Natural Spline Interpolation

Given: Data sets $\{(x_i, f_i)\}_{i=0}^n$, Evaluation point x_{eval} .

Find: Interpolated polynomial $p^{\mathbb{S}}(x_{\text{eval}})$.

Step 1: Compute h_i arrays

Construct $\{h_i\}_{i=0}^{n-1}$,

for $i = 0$ **to** $n - 1$ **do**

$h_i = x_{i+1} - x_i$

end for

Step 2: Construct \mathbf{S} and \mathbf{v}

Initialize \mathbf{S} as a zero matrix of size $(n + 1) \times (n + 1)$ and \mathbf{v} as $(n + 1)$.

for $i = 0$ **to** n **do**

if $i = 0$ **or** $i = n$ **then**

$\mathbf{S}_{ii} = 1$ and $\mathbf{v}_i = 0$

else

$\mathbf{S}_{ii} = 2(h_{i-1} + h_i)$

$\mathbf{S}_{i,i-1} = h_{i-1}$

$\mathbf{S}_{i,i+1} = h_i$

$\mathbf{v}_i = \frac{3}{h_i}(f_{i+1} - f_i) - \frac{3}{h_{i-1}}(f_i - f_{i-1})$

end if

end for

Step 3: Solve $\mathbf{S}\mathbf{c} = \mathbf{v}$

Solve the system $\mathbf{S}\mathbf{c} = \mathbf{v}$ to get coefficient vector \mathbf{c} .

Step 4: Locate x_{eval}

$\text{loc} = 0$

for $i = 0$ **to** n **do**

if $x_{\text{eval}} \leq x_i$ **then**

$\text{loc} = i - 1$

break

end if

end for

Step 5: Compute b_{loc} and d_{loc}

$$b_{\text{loc}} = \frac{f_{\text{loc}+1} - f_{\text{loc}}}{h_{\text{loc}}} - \frac{h_{\text{loc}}}{3}(2\mathbf{c}_{\text{loc}} + \mathbf{c}_{\text{loc}+1})$$

$$d_{\text{loc}} = \frac{\mathbf{c}_{\text{loc}+1} - \mathbf{c}_{\text{loc}}}{3h_{\text{loc}}}$$

Step 6: Evaluate Spline Polynomial $p^{\mathbb{S}}(x)$ at x_{eval}

$$p^{\mathbb{S}}(x_{\text{eval}}) = f_{\text{loc}} + b_{\text{loc}}(x_{\text{eval}} - x_{\text{loc}}) + \mathbf{c}_{\text{loc}}(x_{\text{eval}} - x_{\text{loc}})^2 + d_{\text{loc}}(x_{\text{eval}} - x_{\text{loc}})^3$$

return $p^{\mathbb{S}}(x_{\text{eval}})$

for $j = 0, 1, \dots, n - 1$. In Fig. 1.19 we have the zero degree spline $B_{0,0}(x)$ for $x_j = 0$. Higher-

degree splines are constructed recursively using lower-degree splines as follows:

$$B_{j,k}(x) = \frac{x - x_j}{x_{j+k} - x_j} B_{j,k-1}(x) + \frac{x_{j+k+1} - x}{x_{j+k+1} - x_{j+1}} B_{j+1,k-1}(x) \quad k \geq 1.$$

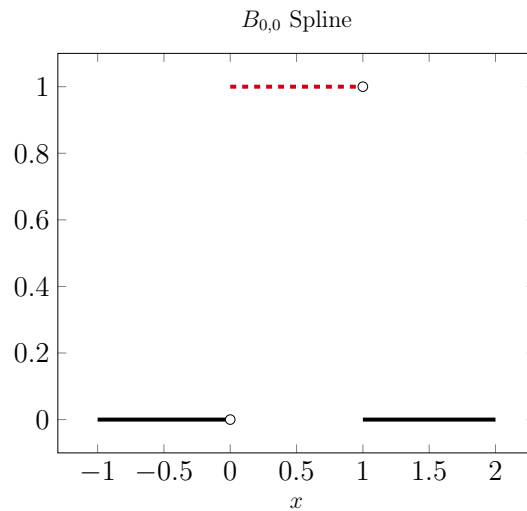


Figure 1.19: Zeroth degree B-spline $B_{0,0}(x)$ for $x_i = 0$.

Although not obvious, one can see that $B_{j,k}(x)$ has one more continuous derivative than $B_{j,k-1}(x)$. Thus while $B_{j,0}(x)$ is discontinuous, $B_{j,1}(x)$ is continuous, $B_{j,2}(x) \in \mathcal{C}^1(\mathbb{R})$, and $B_{j,3}(x) \in \mathcal{C}^2(\mathbb{R})$.

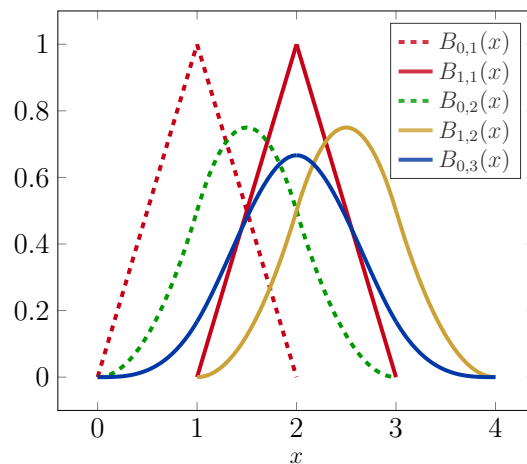


Figure 1.20: Higher degree B-spline polynomial for $x_i = 0$.

As the degree of the B-splines increases, they become more smooth, but the support of $B_{j,k}(x)$ also increases. Based on these results, we can make the following observations:

1. $B_{j,k}(x) \in \mathcal{C}^{k-1}(\mathbb{R})$ (Continuity).
2. $B_{j,k}(x) = 0$ if $x \notin (x_j, x_{j+k+1})$ (Compact Support) ⁴.

⁴**Compact Support:** Let $f : X \rightarrow \mathbb{R}$ be a real-valued function whose domain is an arbitrary set X . The support of f written as $\text{supp}(f)$, is the set of points in X where f is non-zero, i.e., $\text{supp}(f) = \{x \in X : f(x) \neq 0\}$. If $\text{supp}(f)$ is a compact set, then the support is referred to as compact support

3. $B_{j,k}(x) > 0$ for $x \in (x_j, x_{j+k+1})$ (Positivity).

Note: We can notice from Fig. 1.20 that as the degree of the function increases, the support of the function increases as well. Hence, we might get points outside $[x_0, x_n]$. To develop the method, we include additional points beyond the original domain as follows:

$$\cdots < x_{-2} < x_{-1} < x_0 < x_1 < \cdots < x_n < x_{n+1} < \cdots$$

Let $p_k^{\mathbb{S}}(x)$ denote the spline of piecewise polynomial in \mathbb{P}_k . Then we have the following two conditions:

1. $p_k^{\mathbb{S}}(x_i) = f_i$ for $i = 0, 1, \dots, n$.
2. $p_k^{\mathbb{S}} \in \mathcal{C}^{k-1}[x_0, x_n]$ for $k \geq 1$.

Notice that we have an abuse of notation here. In the previous section we used $p_j^{\mathbb{S}}$ to denote the restriction to $[x_j, x_{j+1}]$, whereas here $p_k^{\mathbb{S}}$ denote a spline of degree k .

Let $c_{j,k}$ denote the unknown coefficients, then

$$p_k^{\mathbb{S}}(x) = \sum_j c_{j,k} B_{j,k}(x).$$

Now the question remains on what values of j the summation applies. For the greatest flexibility, we take j for which

$$B_{j,k}(x) \neq 0 \quad \text{for some } x \in [x_0, x_n].$$

Now, for $k \geq 1$, $B_{j,k}(x)$ has support of (x_j, x_{j+k+1}) and hence

$$p_k^{\mathbb{S}}(x) = \sum_{j=-k}^{n-1} c_{j,k} B_{j,k}(x), \quad k > 1.$$

The inclusion of negative indices for j arises due to the support of the B-spline at boundary knots, particularly at x_0 . For the B-spline $B_{j,k}(x)$ to contribute at x_0 , its support must include x_0 . Since the support of $B_{j,k}(x)$ spans from x_j to x_{j+k+1} , and the last point of this support is x_1 when considering x_0 , we require $j + k + 1 = 1$. Solving for j , this gives $j = -k$, which explains the inclusion of “ghost points” $x_{-1}, x_{-2}, \dots, x_{-k}$ in the extended knot sequence. At the other boundary, x_n , the support extends back to x_{n-1} , ensuring that $B_{j,k}(x)$ contributes only within the domain of the spline. To include all valid intervals in the original knot sequence, the upper bound for j is $j \leq n - 1$. Thus, the range of j is determined as $-k \leq j \leq n - 1$, ensuring that the spline remains well-defined and accounts for boundary contributions at x_0 and x_n .

Hence we have $n + k$ (include $k = 0$) coefficients (unknowns) that satisfy the $n + 1$ interpolation condition,

$$p_k^{\mathbb{S}}(x_i) = f_i = \sum_{j=-k}^{n-1} c_{j,k} B_{j,k}(x_i) \quad i = 0, 1, \dots, n.$$

The system becomes underdetermined for higher degrees k , meaning there are more unknowns than equations. As we observed in the cubic interpolation, we might need to impose more conditions to get a system of equations.

Chapter 2

System of Equations

In many applications of science and engineering, solving a system of equations is essential. One prominent example arises in Operations Research, where traffic flow modelling involves solving such systems. The foundational work of Ford and Fulkerson [4] introduced the maximum flow problem, which significantly advanced the theory and applications of system-solving techniques.

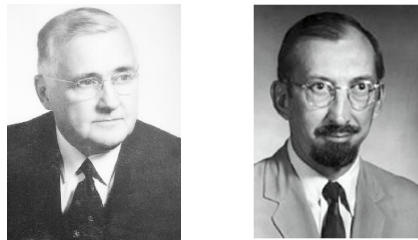


Figure 2.1: Lester Randolph Ford Jr. (23 September 1927–26 February 2017, left) and Delbert Ray Fulkerson (14 August 1924–10 January 1976, right).

Another important application arises in the discretisation of differential equations, a technique widely used in civil, mechanical, and electrical engineering. When differential equations are discretised, the resulting system of equations often takes the form of a *band matrix*¹. Depending on the choice of polynomial approximation used in the discretisation, the resulting band matrix can be tridiagonal, pentadiagonal, or a more general band matrix. Efficiently solving these systems is crucial to obtaining solutions to the differential equations.

In this chapter, we first introduce *direct methods* for solving systems of equations, followed by iterative methods. Direct methods aim to find the exact solution theoretically in a finite number of steps, though practical computations are subject to round-off errors, which must be carefully managed to ensure accuracy.

¹**Band Matrix:** A matrix $\{a_{ij}\}_{i,j=1}^n$ is called a band matrix if all elements outside a certain diagonal band are zero. The band is determined by:

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

where k_1 and k_2 are the lower and upper bandwidths, respectively. Special cases include diagonal matrices ($k_1 = k_2 = 0$) and tridiagonal matrices ($k_1 = k_2 = 1$).

Let us for the uniformity of the notation denote b_i by $a_{i,n+1}$ for $i = 1, 2, \dots, n$ then

$$\tilde{\mathbf{A}} = [\mathbf{A}, \mathbf{b}] = \left[\begin{array}{cccc|c} a_{11} & a_{12} & \dots & a_{1n} & a_{1,n+1} \\ a_{21} & a_{22} & \dots & a_{2n} & a_{2,n+1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} & a_{n,n+1} \end{array} \right]. \quad (2.3)$$

Provided $a_{11} \neq 0$ we perform the operations corresponding to

$$R_j \mapsto R_j - \frac{a_{j1}}{a_{11}} R_1 \quad \text{for } j = 2, 3, \dots, n,$$

to eliminate the coefficient x_1 in each of the rows. Once the coefficients of x_1 are cancelled, we do the same for other rows and follow a sequential procedure for $i = 2, 3, \dots, n-1$ and perform the operation

$$R_j \mapsto R_j - \frac{a_{ji}}{a_{ii}} R_i \quad \text{for } j = i+1, i+2, \dots, n.$$

The resulting matrix has the form

$$\tilde{\tilde{\mathbf{A}}} = [\mathbf{A}, \mathbf{b}] = \left[\begin{array}{cccc|c} a_{11} & a_{12} & \dots & a_{1n} & a_{1,n+1} \\ 0 & \tilde{a}_{22} & \dots & \tilde{a}_{2n} & \tilde{a}_{2,n+1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & \tilde{a}_{nn} & \tilde{a}_{n,n+1} \end{array} \right].$$

This system of equation has the same solution set as Eq. (2.1). But the new system of equation has the form:

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n &= a_{1,n+1} \\ \tilde{a}_{22}x_2 + \dots + \tilde{a}_{2n}x_n &= \tilde{a}_{2,n+1} \\ &\vdots \\ \tilde{a}_{nn}x_n &= \tilde{a}_{n,n+1}. \end{aligned}$$

By backward substitution we get

$$x_n = \frac{\tilde{a}_{n,n+1}}{\tilde{a}_{nn}}.$$

Solving the $(n-1)^{\text{th}}$ equation for x_{n-1} and using the value of x_n we get

$$\begin{aligned} \tilde{a}_{n-1,n-1}x_{n-1} + \tilde{a}_{n-1,n}x_n &= \tilde{a}_{n-1,n+1} \\ x_{n-1} &= \frac{\tilde{a}_{n-1,n+1} - \tilde{a}_{n-1,n}x_n}{\tilde{a}_{n-1,n-1}}. \end{aligned}$$

Continuing this process we get

$$x_i = \frac{\tilde{a}_{i,n+1} - \sum_{j=i+1}^n \tilde{a}_{ij}x_j}{\tilde{a}_{ii}},$$

for $i = n-1, n-2, \dots, 2, 1$ where for $i = 1$, $\tilde{a}_{1,n+1} = a_{1,n+1}$ and $\tilde{a}_{11} = a_{11}$.

Gaussian elimination can also be seen more precisely by forming a sequence of augmented matrices $\tilde{\mathbf{A}}^{(1)}, \tilde{\mathbf{A}}^{(2)}, \dots, \tilde{\mathbf{A}}^{(n)}$ where $\tilde{\mathbf{A}}^{(1)}$ is the matrix given in Eq. (2.3) and a general $\tilde{\mathbf{A}}^{(k)}$ matrix for $k = 2, 3, \dots, n$ is given by

$$\tilde{\mathbf{A}}^{(k)} = \left[\begin{array}{cccc|cccc|c} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} & \cdots & a_{1,k-1}^{(1)} & a_{1,k}^{(1)} & \cdots & a_{1,n}^{(1)} & a_{1,n+1}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} & \cdots & a_{2,k-1}^{(2)} & a_{2,k}^{(2)} & \cdots & a_{2,n}^{(2)} & a_{2,n+1}^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & a_{k-1,k-1}^{(k-1)} & a_{k-1,k}^{(k-1)} & \cdots & a_{k-1,n}^{(k-1)} & a_{k-1,n+1}^{(k-1)} \\ 0 & 0 & 0 & \cdots & 0 & a_{k,k}^{(k)} & \cdots & a_{k,n}^{(k)} & a_{k,n+1}^{(k)} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & a_{n,k}^{(k)} & \cdots & a_{n,n}^{(k)} & a_{n,n+1}^{(k)} \end{array} \right]. \quad (2.4)$$

where x_{k-1} has been eliminated from R_k, \dots, R_n .

In general the matrix entries are given by

$$a_{ij}^{(k)} = \begin{cases} a_{ij}^{(k-1)} & \text{if } i = 1, 2, \dots, k-1 \text{ and } j = 1, 2, \dots, n+1, \\ 0 & \text{if } i = k, k+1, \dots, n \text{ and } j = 1, 2, \dots, k-1, \\ a_{ij}^{(k-1)} - \frac{a_{i,k-1}^{(k-1)}}{a_{k-1,k-1}^{(k-1)}} a_{k-1,j}^{(k-1)} & \text{if } i = k, k+1, \dots, n \text{ and } j = k, k+1, \dots, n+1. \end{cases}$$

This procedure will fail if any of the elements $\{a_{ii}^{(i)}\}$ for $i = 1, 2, \dots, n$ is zero as

$$R_i \mapsto R_i - \frac{a_{i,k}^{(k)}}{a_{k,k}^{(k)}} R_k$$

cannot be performed or the backward substitution fails.

The system may still have solution but the technique might be altered. For example, consider the augmented matrix

$$\tilde{\mathbf{A}} = \tilde{\mathbf{A}}^{(1)} = \left[\begin{array}{cccc|c} 1 & -1 & 2 & -1 & -8 \\ 2 & -2 & 3 & -3 & -20 \\ 1 & 1 & 1 & 0 & -2 \\ 1 & -1 & 4 & 3 & 4 \end{array} \right].$$

Performing the operations, $R_3 \mapsto R_3 - R_1$, $R_4 \mapsto R_4 - R_1$, and $R_2 \mapsto R_2 - 2R_1$,

$$\tilde{\mathbf{A}} = \tilde{\mathbf{A}}^{(2)} = \left[\begin{array}{cccc|c} 1 & -1 & 2 & -1 & -8 \\ 0 & 0 & -1 & -1 & -4 \\ 0 & 2 & -1 & 1 & 6 \\ 0 & 0 & 2 & 4 & 12 \end{array} \right].$$

Here $a_{22}^{(2)}$ is zero and is called the *pivot* element. Hence the procedure cannot proceed. So we search the second column for first non-zero entry after 2nd row. Since, $a_{32}^{(2)} \neq 0$, we perform $R_2 \leftrightarrow R_3$ and then proceed.

The above example shows what happens if $a_{kk}^{(k)} = 0$ for some $k = 1, 2, \dots, n-1$. In this case we follow:

1. The k^{th} column of $\tilde{\mathbf{A}}^{(k-1)}$ is searched from the k^{th} row to the n^{th} row for first non zero entry, $a_{pk}^{(k)} \neq 0$ for $k + 1 \leq p \leq n$.
2. Then $R_p \leftrightarrow R_k$ is performed to get a temporary matrix $\tilde{\mathbf{A}}^{(k-1)'}$ and then the usual elimination follows.

In the case $a_{pk}^{(k)} = 0$ for each p , then the system does not have a unique solution as two columns are the linearly dependent. Finally, if $a_{nn}^{(n)} = 0$ then the system does not have a unique solution.

The algorithm for the Gaussian elimination is provided in Algorithm 6. Although the algorithm looks like we are creating new matrices $\tilde{\mathbf{A}}^{(i)}$ for $i = 1, 2, \dots, n$ but we can perform all the computation using only one $n \times (n + 1)$ matrix for storage.

2.1.1 Computational Complexity

Now we look at the computational complexity of the Gaussian elimination. Generally time taken to perform a multiplication or division is generally more than addition or subtraction. Hence, we count these operations separately.

The arithmetic operations happens in Step 2.3:

1. **Computation of m_{ki} :** Requires division and $(n - i)$ operations.
2. **Multiplication of $m_{ki}R_i$:** This multiplication happens with the non-zero entries of R_i which is $(n - i) \times (n - i + 1)$ as non zero entries is given by $(n - i + 1)$ in the R_i^{th} row.
3. **Subtraction for $R_k - m_{ki}R_i$:** These will also $(n - i + 1) \times (n - i)$ as we subtract the non-zero entries.

The first two are multiplication and division and the last one is addition/subtraction.

Multiplication/Division Complexity

Now, $(n - i) + (n - i) \times (n - i + 1) = (n - i) \times (n - i + 2) = (n - i)^2 + 2(n - i)$. Summing i from 1 to $n - 1$ we get

$$\begin{aligned} \sum_{i=1}^{n-1} (n - i)(n - i + 2) &= \sum_{i=1}^{n-1} (n - i)^2 + 2 \sum_{i=1}^{n-1} (n - i) \\ &= \sum_{i=1}^{n-1} i^2 + 2 \sum_{i=1}^{n-1} i \\ &= \frac{(n - 1)n(2n - 1)}{6} + \frac{2n(n - 1)}{2} = \frac{2n^3 + 3n^2 - 5n}{6}. \end{aligned}$$

In the above equation we have used the basic identities of summation, namely $\sum_{i=1}^n i^2$ and $\sum_{i=1}^n i$.

Algorithm 6 Gauss Elimination

Given: Matrix \mathbf{A} , right hand side \mathbf{b} and dimension n .

Find: Solution \mathbf{x} .

Step 1: Create Augmented Matrix $\tilde{\mathbf{A}}$

Initialize $\tilde{\mathbf{A}}$ as a zero matrix of size $n \times (n + 1)$

```

for  $i = 1$  to  $n$  do
  for  $j = 1$  to  $n + 1$  do
    if  $j \leq n$  then
       $\tilde{\mathbf{A}}_{ij} = \mathbf{A}_{ij}$ 
    else
       $\tilde{\mathbf{A}}_{ij} = \mathbf{b}_i$ 
    end if
  end for
end for
end for

```

Step 2: Reduce the matrix to Row-Echelon form

for $i = 1$ to $n - 1$ do

Step 2.1: Check Pivot

Initialize $p = -1$

for $q = i$ to n do

if $\tilde{\mathbf{A}}_{qi} \neq 0$ then

$p = q$

break

end if

end for

if $p = -1$ then

Output("No Unique Solution")

exit()

end if

Step 2.2: Exchange Rows $R_i \leftrightarrow R_p$

if $p \neq i$ then

temp = 0

for $j = 1$ to $n + 1$ do

temp = $\tilde{\mathbf{A}}_{ij}$

$\tilde{\mathbf{A}}_{ij} = \tilde{\mathbf{A}}_{pj}$

$\tilde{\mathbf{A}}_{pj} =$ temp

end for

end if

Step 2.3: Matrix Reduction

$m_{ki} = 0$

for $k = i + 1$ to n do

$m_{ki} = \tilde{\mathbf{A}}_{ki} / \tilde{\mathbf{A}}_{ii}$

for $j = i$ to $n + 1$ do

$\tilde{\mathbf{A}}_{kj} = \tilde{\mathbf{A}}_{kj} - m_{ki} \tilde{\mathbf{A}}_{ij}$

end for

end for

end for

Step 3: Check for no Solution

if $\tilde{\mathbf{A}}_{nn} = 0$ then

Output("No Unique Solution")

exit()

end if

Step 4: Backward Substitution

Initialize x as a vector of size n

$x_n = \frac{\tilde{\mathbf{A}}_{n,n+1}}{\tilde{\mathbf{A}}_{nn}}$

for $i = n - 1$ to 1 do

sum = 0

for $j = i + 1$ to n do

sum = sum + $\tilde{\mathbf{A}}_{ij} x_j$

end for

$x_i = \frac{\tilde{\mathbf{A}}_{i,n+1} - \text{sum}}{\tilde{\mathbf{A}}_{ii}}$

end for

return $\{x_i\}_{i=1}^n$

Addition/Subtraction Complexity

$$\begin{aligned}
\sum_{i=1}^{n-1} (n-i)(n-i+1) &= \sum_{i=1}^{n-1} (n-i)^2 + \sum_{i=1}^{n-1} (n-i) \\
&= \sum_{i=1}^{n-1} i^2 + \sum_{i=1}^{n-1} i \\
&= \frac{n(n-1)(2n-1)}{6} + \frac{n(n-1)}{2} = \frac{n^3 - n}{3}.
\end{aligned}$$

Hence we notice that in Step 2.3 we require $\mathcal{O}(n^3)$ operations.

The next step that require arithmetic operations are the ones in backward substitution, i.e, Step 4. First is in the computation of x_n which requires one division. For the computation of rest of the $\{x_i\}$ we need $(n-i)$ multiplications and one division for each i and $(n-i-1)$ addition for each summation followed by one subtraction.

Multiplication/Division Complexity

$$\begin{aligned}
1 + \sum_{i=1}^{n-1} ((n-i) + 1) &= 1 + \left(\sum_{i=1}^{n-1} (n-i) \right) + n - 1 \\
&= n + \sum_{i=1}^{n-1} (n-i) \\
&= n + \sum_{i=1}^{n-1} i \\
&= n + \frac{n(n-1)}{2} = \frac{n^2 + n}{2}.
\end{aligned}$$

Addition/Subtraction Complexity

$$\sum_{i=1}^{n-1} ((n-i-1) + 1) = \sum_{i=1}^{n-1} (n-i) = \frac{n^2 - n}{2}.$$

Hence in total we require

$$\frac{2n^3 + 3n^2 - 5n}{6} + \frac{n^2 + n}{2} = \frac{n^3}{3} + n^2 - \frac{n}{3},$$

operations for multiplication and division; and

$$\frac{n^3 - n}{3} + \frac{n^2 - n}{2} = \frac{2n^3 + 3n^2 - 5n}{6},$$

for addition and subtraction. Hence we have $\mathcal{O}(n^3/3)$ computational complexity.

2.1.2 Gauss-Jordan Algorithm

Wilhelm Jordan was a geodesist who extended the basic Gaussian elimination to achieve a full row-reduced echelon form of a matrix. Do not confuse Wilhelm Jordan with Camille Jordan (who gave us Jordan Curve theorem and Jordan Canonical form).



Figure 2.3: Wilhelm Jordan: 1 March 1842-17 April 1899.

This method is a variation of Gaussian elimination where the variable x_i is not only removed from $R_{i+1}, R_{i+2}, \dots, R_n$ but also from R_1, R_2, \dots, R_{i-1} . Upon this reduction the augmented matrix looks like

$$[\mathbf{A}, \mathbf{b}] = \left[\begin{array}{cccc|c} a_{11}^{(1)} & 0 & \dots & 0 & a_{1,n+1}^{(1)} \\ 0 & a_{22}^{(2)} & \dots & 0 & a_{2,n+1}^{(2)} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & a_{nn}^{(n)} & a_{n,n+1}^{(n)} \end{array} \right].$$

Then the solution can be easily obtained using

$$x_i = \frac{a_{i,n+1}^{(i)}}{a_{ii}^{(i)}}, \quad \text{for } i = 1, 2, \dots, n.$$

The Gauss-Jordan algorithm is presented in Algorithm 7.

2.2 Matrix Factorisation

Like polynomial interpolation, which was the basis for developing more efficient algorithms such as Lagrange and Newton divided differences, Gaussian elimination is the foundation for more advanced topics.

Gaussian elimination consists of two steps: the row-reduction step and backward substitution. The former has a computational complexity of $\mathcal{O}(n^3)$, while the latter requires only $\mathcal{O}(n^2)$. This means that if we have a triangular matrix, solving the system requires only $\mathcal{O}(n^2)$ operations.

2.2.1 LU Decomposition

Suppose that we have $\mathbf{A} = \mathbf{L}\mathbf{U}$, meaning that \mathbf{A} has been factored into a lower triangular matrix (\mathbf{L})² and an upper triangular matrix (\mathbf{U})³. Then, solving $\mathbf{A}\mathbf{x} = \mathbf{b}$ can be done in two

²**Lower Triangular Matrix:** A matrix $\mathbf{A} = \{a_{ij}\}_{i,j=1}^n$ is said to be lower triangular if $a_{ij} = 0$ for $i > j$.

³**Upper Triangular Matrix:** A matrix $\mathbf{A} = \{a_{ij}\}_{i,j=1}^n$ is said to be upper triangular if $a_{ij} = 0$ for $j > i$.

Algorithm 7 Gauss Jordan

Given: Matrix \mathbf{A} , right hand side \mathbf{b} and dimension n .

Find: Solution \mathbf{x} .

Step 1: Create Augmented Matrix $\tilde{\mathbf{A}}$

Initialize $\tilde{\mathbf{A}}$ as a zero matrix of size $n \times (n + 1)$

```

for  $i = 1$  to  $n$  do
  for  $j = 1$  to  $n + 1$  do
    if  $j \leq n$  then
       $\tilde{\mathbf{A}}_{ij} = \mathbf{A}_{ij}$ 
    else
       $\tilde{\mathbf{A}}_{ij} = \mathbf{b}_i$ 
    end if
  end for
end for

```

Step 2: Reduce the matrix to Row-Echelon form

for $i = 1$ to n do

Step 2.1: Check Pivot

Initialize $p = -1$

```

for  $q = i$  to  $n$  do
  if  $\tilde{\mathbf{A}}_{qi} \neq 0$  then
     $p = q$ 
    break
  end if
end for
if  $p = -1$  then
  Output("No Unique Solution")
  exit()
end if

```

Step 2.2: Exchange Rows $R_i \leftrightarrow R_p$

```

if  $p \neq i$  then
  temp = 0
  for  $j = 1$  to  $n + 1$  do
    temp =  $\tilde{\mathbf{A}}_{ij}$ 
     $\tilde{\mathbf{A}}_{ij} = \tilde{\mathbf{A}}_{pj}$ 
     $\tilde{\mathbf{A}}_{pj} = temp$ 
  end for
end if

```

Step 2.3: Matrix Reduction

$m_{ki} = 0$

```

for  $k = 1$  to  $n$  do
  if  $k = i$  then
    continue
  else
     $m_{ki} = \tilde{\mathbf{A}}_{ki} / \tilde{\mathbf{A}}_{ii}$ 
    for  $j = i$  to  $n + 1$  do
       $\tilde{\mathbf{A}}_{kj} = \tilde{\mathbf{A}}_{kj} - m_{ki} \tilde{\mathbf{A}}_{ij}$ 
    end for
  end if
end for
end for

```

Step 3: Check for no Solution

```

if  $\tilde{\mathbf{A}}_{nn} = 0$  then
  Output("No Unique Solution")
  exit()
end if

```

Step 4: Backward Substitution

Initialize \mathbf{x} as a zero vector of size n

```

for  $i = 1$  to  $n$  do
   $x_i = \frac{\tilde{\mathbf{A}}_{i,n+1}}{\tilde{\mathbf{A}}_{ii}}$ 
end for

```

return $\{x_i\}_{i=1}^n$

steps:

- Solve $\mathbf{L}\mathbf{y} = \mathbf{b}$ for \mathbf{y} .
- Solve $\mathbf{U}\mathbf{x} = \mathbf{y}$ for \mathbf{x} .

Both steps require only $\mathcal{O}(n^2)$ operations.

Although several mathematicians introduced LU decomposition, the Polish mathematician Tadeusz Banachiewicz is credited with generalizing the method for arbitrary matrices.



Figure 2.4: Tadeusz Banachiewicz: 13 February 1882 - 17 November 1954.

LU decomposition reduces an $\mathcal{O}(n^3/3)$ problem to an $\mathcal{O}(2n^2)$ problem. This reduction is useful but comes at a cost: the factorization of \mathbf{A} into \mathbf{L} and \mathbf{U} itself requires $\mathcal{O}(n^3/3)$ operations. However, once computed, the factorization can be stored and reused for multiple right-hand-side vectors \mathbf{b} .

To proceed with LU decomposition, we assume that $\mathbf{A}\mathbf{x} = \mathbf{b}$ can be solved using Gaussian elimination without row pivoting, i.e., $a_{ii}^{(i)} \neq 0$ for $i = 1, 2, \dots, n$.

The first step in Gaussian elimination consists of performing, for each $j = 2, 3, \dots, n$,

$$R_j \mapsto R_j - m_{j1}R_1, \quad \text{where } m_{j1} = \frac{a_{j1}^{(1)}}{a_{11}^{(1)}}.$$

An equivalent way of viewing this is by multiplying \mathbf{A} on the left by the matrix $\mathbf{M}^{(1)}$, where

$$\mathbf{M}^{(1)} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ -m_{21} & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ -m_{n1} & 0 & \cdots & 1 \end{bmatrix}.$$

This is called the *first Gaussian transformation matrix*. The product of this matrix with \mathbf{A} is denoted by $\mathbf{A}^{(2)}$, so that

$$\mathbf{A}^{(2)} = \mathbf{M}^{(1)}\mathbf{A}.$$

Similarly, the right-hand side vector is updated as

$$\mathbf{b}^{(2)} = \mathbf{M}^{(1)}\mathbf{b}.$$

Next, we construct $\mathbf{M}^{(2)}$ by replacing the subdiagonal entries in the second column of the identity matrix with the negative of the multipliers

$$m_{j2} = \frac{a_{j2}^{(2)}}{a_{22}^{(2)}}.$$

This process continues until we obtain an upper triangular matrix $\mathbf{A}^{(n)}$, given by

$$\mathbf{A}^{(n)} = \mathbf{M}^{(n-1)}\mathbf{M}^{(n-2)} \dots \mathbf{M}^{(1)}\mathbf{A}.$$

At this point, we define $\mathbf{U} = \mathbf{A}^{(n)}$ as the upper triangular matrix in the LU factorization.

To compute the lower triangular matrix \mathbf{L} , we note that the inverse of each $\mathbf{M}^{(k)}$ matrix is given by

$$\mathbf{L}^{(k)} = [\mathbf{M}^{(k)}]^{-1} = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 & 0 & \dots & 0 \\ 0 & 0 & \dots & m_{k+1,k} & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & m_{n,k} & 0 & \dots & 1 \end{bmatrix}.$$

The lower triangular matrix \mathbf{L} is then obtained as

$$\mathbf{L} = \mathbf{L}^{(1)}\mathbf{L}^{(2)} \dots \mathbf{L}^{(n-1)}.$$

Since each $\mathbf{L}^{(k)}$ is the inverse of $\mathbf{M}^{(k)}$, we confirm that

$$\mathbf{LU} = \mathbf{A}.$$

Theorem 2.1. (Doolittle LU Decomposition) *If Gaussian elimination can be performed on the system $\mathbf{Ax} = \mathbf{b}$ without row interchanges, then the matrix \mathbf{A} can be factored as $\mathbf{A} = \mathbf{LU}$, where*

$$m_{ji} = \frac{a_{ji}^{(i)}}{a_{ii}^{(i)}},$$

and

$$\mathbf{U} = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & a_{nn}^{(n)} \end{bmatrix}, \quad \mathbf{L} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ m_{21} & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ m_{n1} & m_{n2} & \dots & 1 \end{bmatrix}.$$

The above factorization is the Doolittle method, where \mathbf{L} has ones on its diagonal. Alternatively, if the ones are placed on the diagonal of \mathbf{U} , the technique is called Crout's LU decomposition.

Once the LU factorization is obtained, the system $\mathbf{Ax} = \mathbf{LUx}$ is solved efficiently by first computing \mathbf{y} from $\mathbf{Ly} = \mathbf{b}$ using forward substitution and then solving $\mathbf{Ux} = \mathbf{y}$ using backward substitution.

It is important to note that not all square matrices have an LU factorization. For example, the matrix

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

has no LU factorization. Suppose it did; then there would exist \mathbf{L} and \mathbf{U} such that $\mathbf{A} = \mathbf{LU}$. However, this would lead to a contradiction, as one of the factors would necessarily be singular while \mathbf{A} is not.

Next, we note that the LU decomposition is not unique.

Theorem 2.2. *If a matrix has an LU decomposition, then it is not unique.*

Proof. Let \mathbf{A} have an LU decomposition, i.e., $\mathbf{A} = \mathbf{LU}$. Then, we can write

$$\begin{aligned}\mathbf{A} &= \mathbf{LU} \\ &= \mathbf{LDD}^{-1}\mathbf{U} \\ &= (\mathbf{LD})(\mathbf{D}^{-1}\mathbf{U}),\end{aligned}$$

where \mathbf{D} is any diagonal matrix. Since \mathbf{LD} remains lower triangular and $\mathbf{D}^{-1}\mathbf{U}$ is still upper triangular, we obtain infinitely many LU decompositions of \mathbf{A} by varying \mathbf{D} . \square

PLU Decomposition

So far, we have assumed that LU decomposition is applicable to systems of equations that do not require pivoting. However, in general, pivoting is necessary. To introduce LU decomposition with pivoting, we first define the permutation matrix.

Definition 2.3. A *permutation matrix* $\mathbf{P} = \{p_{ij}\}_{i,j=1}^n$ is an $n \times n$ matrix obtained by rearranging the rows of the identity matrix.

For example, the matrix

$$\mathbf{P} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

is a 3×3 permutation matrix where the second and third rows are interchanged. For any 3×3 matrix \mathbf{A} , multiplying by \mathbf{P} on the left swaps these two rows:

$$\mathbf{PA} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{31} & a_{32} & a_{33} \\ a_{21} & a_{22} & a_{23} \end{bmatrix}.$$

Let k_1, k_2, \dots, k_n be a permutation of $1, 2, \dots, n$. The permutation matrix \mathbf{P} is then defined as:

$$p_{ij} = \begin{cases} 1 & \text{if } j = k_i, \\ 0 & \text{otherwise.} \end{cases}$$

This satisfies the following properties:

1. \mathbf{PA} permutes the rows of \mathbf{A} :

$$\mathbf{PA} = \begin{bmatrix} a_{k_1 1} & a_{k_1 2} & \cdots & a_{k_1 n} \\ a_{k_2 1} & a_{k_2 2} & \cdots & a_{k_2 n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{k_n 1} & a_{k_n 2} & \cdots & a_{k_n n} \end{bmatrix}.$$

2. The inverse of a permutation matrix exists and is given by $\mathbf{P}^{-1} = \mathbf{P}^\top$.

In the previous section, we saw that for any nonsingular matrix \mathbf{A} , the linear system $\mathbf{Ax} = \mathbf{b}$ can be solved using Gaussian elimination with row interchanges. If the required row interchanges are known beforehand, we can apply them initially, allowing us to use LU decomposition without further row swaps. That is, for any nonsingular matrix \mathbf{A} , there exists a permutation matrix \mathbf{P} such that the system

$$\mathbf{PAx} = \mathbf{Pb}$$

can be solved without row interchanges. Consequently, we can factorize \mathbf{PA} as

$$\mathbf{PA} = \mathbf{LU}.$$

Since \mathbf{P} is a permutation matrix, we have $\mathbf{P}^{-1} = \mathbf{P}^\top$, which implies

$$\mathbf{A} = (\mathbf{P}^\top \mathbf{L}) \mathbf{U}.$$

While \mathbf{U} remains upper triangular, the matrix $\mathbf{P}^\top \mathbf{L}$ may not necessarily be lower triangular unless $\mathbf{P} = \mathbf{I}$.

Based on this, we establish the following lemma.

Lemma 2.4. *Let \mathbf{A} be an $n \times n$ matrix. Then, there exists a permutation matrix \mathbf{P} such that \mathbf{PA} has an LU decomposition, i.e., $\mathbf{PA} = \mathbf{LU}$.*

The next theorem addresses the uniqueness of the LU decomposition.

Theorem 2.5. *Let \mathbf{A} be an $n \times n$ matrix, and let \mathbf{P} be an $n \times n$ permutation matrix such that \mathbf{PA} has an LU decomposition. If \mathbf{A} is invertible, then there exists a unique $n \times n$ lower triangular matrix \mathbf{L} with all diagonal entries equal to 1, and a unique $n \times n$ upper triangular matrix \mathbf{U} such that*

$$\mathbf{PA} = \mathbf{LU}.$$

Proof. The existence of the LU decomposition follows from Lemma 2.4. We now prove the uniqueness.

Suppose \mathbf{L} is not unit lower triangular. Then, we can express the decomposition as

$$\mathbf{PA} = \mathbf{LU}.$$

Rewriting,

$$\mathbf{PA} = \mathbf{LD}^{-1}\mathbf{DU},$$

where \mathbf{D} is a diagonal matrix whose diagonal entries match those of \mathbf{L} . Since \mathbf{A} is invertible, \mathbf{L} is also invertible, ensuring that \mathbf{D}^{-1} exists. Defining

$$\mathbf{L}_1 = \mathbf{LD}^{-1}, \quad \mathbf{U}_1 = \mathbf{DU},$$

we obtain a new factorization with \mathbf{L}_1 as a unit lower triangular matrix and \mathbf{U}_1 as an upper triangular matrix:

$$\mathbf{PA} = \mathbf{L}_1\mathbf{U}_1.$$

Now, suppose there exists another decomposition:

$$\mathbf{PA} = \mathbf{L}_2\mathbf{U}_2,$$

where \mathbf{L}_2 is also unit lower triangular. Then, we equate the two decompositions:

$$\mathbf{L}_1\mathbf{U}_1 = \mathbf{L}_2\mathbf{U}_2.$$

Since \mathbf{A} is invertible, both \mathbf{L}_1 and \mathbf{P} are invertible, implying that $\mathbf{U}_1 = \mathbf{L}_1^{-1}\mathbf{PA}$ is also invertible.

Thus, we obtain

$$\mathbf{L}_2^{-1}\mathbf{L}_1 = \mathbf{U}_2\mathbf{U}_1^{-1}. \quad (2.5)$$

Since:

1. The inverse of a lower (upper) triangular matrix is lower (upper) triangular.
2. The product of lower (upper) triangular matrices remains lower (upper) triangular.

it follows that $\mathbf{L}_2^{-1}\mathbf{L}_1$ is lower triangular, and $\mathbf{U}_2\mathbf{U}_1^{-1}$ is upper triangular. Since $\mathbf{L}_2^{-1}\mathbf{L}_1$ is also unit diagonal, the only possibility is

$$\mathbf{L}_2^{-1}\mathbf{L}_1 = \mathbf{I} \Rightarrow \mathbf{L}_2 = \mathbf{L}_1.$$

Similarly, we obtain $\mathbf{U}_1 = \mathbf{U}_2$, proving uniqueness. □

2.2.2 LDL^T Decomposition

In linear algebra we have certain special matrices and they enjoy certain “good” properties. This is true with respect to their LU decomposition as well. We first mention certain matrices, followed by their properties, and then their special kind of factorisation.

Definition 2.6. A matrix \mathbf{A} is said to be *diagonally dominant* when

$$|a_{ii}| \geq \sum_{j=1, j \neq i}^n |a_{ij}| \quad \forall i = 1, 2, \dots, n.$$

If the inequality is strict then it is called as *strictly diagonally dominant*.

Algorithm 8 LU Decomposition with Partial Pivoting

Given: Matrix \mathbf{A} of size $n \times n$.

Find: Matrices \mathbf{L} , \mathbf{U} , and \mathbf{P} such that $\mathbf{PA} = \mathbf{LU}$.

Step 1: Initialize Matrices

Initialize \mathbf{L} as an $n \times n$ identity matrix.

Initialize \mathbf{P} as an $n \times n$ identity matrix.

Initialize \mathbf{U} as \mathbf{A} .

Step 2: Perform LU Decomposition

for $i = 1$ to $n - 1$ do

Step 2.1: Check Pivot

Initialize $p = -1$

for $q = i$ to n do

 if $U_{qi} \neq 0$ then

$p = q$

 break

 end if

end for

if $p = -1$ then

 Output("Matrix is singular but the LU decomposition still exists!")

 continue

end if

Step 2.2: Exchange Rows for \mathbf{P} and \mathbf{U} , $R_i \leftrightarrow R_p$

if $p \neq i$ then

 temp1 = 0; temp2 = 0.

 for $j = 1$ to $n + 1$ do

 temp1 = \mathbf{P}_{ij} ; temp2 = \mathbf{U}_{ij}

$\mathbf{P}_{ij} = \mathbf{P}_{pj}$; $\mathbf{A}_{ij} = \mathbf{A}_{pj}$.

$\mathbf{P}_{pj} = \text{temp1}$; $\mathbf{A}_{pj} = \text{temp2}$.

 end for

end if

Step 2.3: Matrix Reduction

for $k = i + 1$ to n do

$m_{ki} = \mathbf{U}_{ki} / \mathbf{U}_{ii}$

$\mathbf{L}_{ki} = m_{ki}$

 for $j = i$ to n do

$\mathbf{U}_{kj} = \mathbf{U}_{kj} - m_{ki} \mathbf{U}_{ij}$

 end for

end for

end for

return $\mathbf{L}, \mathbf{U}, \mathbf{P}$

Theorem 2.7. *A strictly diagonally dominant matrix \mathbf{A} is non-singular. Moreover, in this case, Gaussian elimination can be performed on any linear system of the form $\mathbf{Ax} = \mathbf{b}$ to obtain its unique solution without row or column interchanges, and the computations will be stable with respect to the growth of round-off errors.*

Proof. This theorem has three parts:

1. Non-Singularity of \mathbf{A} .
2. Unique solution using Gaussian Elimination and no row-interchange.
3. Stability of the solution.

We will prove the first two parts, as the proof of the third part is out of the scope of this lecture. For the first part we use the method of contradiction. Suppose \mathbf{A} is singular. Then the system $\mathbf{Ax} = \mathbf{0}$ has non-trivial solution, say $\mathbf{x} = \{x_i\}$. Let k be an index for which

$$0 < |x_k| = \max_{1 \leq j \leq n} |x_j|.$$

As $\mathbf{Ax} = \mathbf{0}$, we get $\sum_{j=1}^n a_{ij}x_j = 0$ for $i = 1, 2, \dots, n$. At $i = k$

$$\sum_{j=1}^n a_{kj}x_j = 0 \Rightarrow a_{kk}x_k = - \sum_{j=1, j \neq k}^n a_{kj}x_j.$$

From the triangular inequality we have

$$\begin{aligned} |a_{kk}||x_k| &= \left| \sum_{j=1, j \neq k}^n a_{kj}x_j \right| \\ &\leq \sum_{j=1, j \neq k}^n |a_{kj}||x_j| \\ &< \sum_{j=1, j \neq k}^n |a_{kj}||x_k| \end{aligned}$$

Hence, $|a_{kk}| < \sum_{j=1, j \neq k}^n |a_{kj}|$ which is a contradiction as \mathbf{A} is strictly diagonally dominant. Hence the matrix \mathbf{A} is singular.

For the second part we show that the matrices $\mathbf{A}^{(k)}$ for $k = 2, 3, \dots, n$ generated during the Gaussian elimination is strictly diagonally dominant. Hence it ensure that each pivot element is non-zero.

Since \mathbf{A} is strictly diagonally dominant, $a_{11} \neq 0$ and $\mathbf{A}^{(2)}$ can be formed. Thus for each $i = 2, 3, \dots, n$,

$$a_{ij}^{(2)} = a_{ij}^{(1)} - \frac{a_{1j}^{(1)}a_{i1}^{(1)}}{a_{11}^{(1)}}, \quad \text{for } 2 \leq j \leq n.$$

First, $a_{i1}^{(2)} = 0$. Now using the triangle inequality

$$\sum_{j=2, j \neq i}^n |a_{ij}^{(2)}| = \sum_{j=2, j \neq i}^n \left| a_{ij}^{(1)} - \frac{a_{1j}^{(1)} a_{i1}^{(1)}}{a_{11}^{(1)}} \right| \leq \sum_{j=2, j \neq i}^n |a_{ij}^{(1)}| + \sum_{j=2, j \neq i}^n \left| \frac{a_{1j}^{(1)} a_{i1}^{(1)}}{a_{11}^{(1)}} \right|.$$

But since \mathbf{A} is strictly diagonally dominant,

$$\begin{aligned} \sum_{j=1, j \neq i}^n |a_{ij}^{(1)}| &< |a_{ii}^{(1)}| \\ \sum_{j=2, j \neq i}^n |a_{ij}^{(1)}| &< |a_{ii}^{(1)}| - |a_{i1}^{(1)}|, \end{aligned}$$

and similarly

$$\begin{aligned} \sum_{j=1, j \neq i}^n |a_{1j}^{(1)}| &< |a_{11}^{(1)}| \\ \sum_{j=2, j \neq i}^n |a_{1j}^{(1)}| &< |a_{11}^{(1)}| - |a_{1i}^{(1)}|, \end{aligned}$$

so

$$\sum_{j=2, j \neq i}^n |a_{ij}^{(2)}| < |a_{ii}^{(1)}| - |a_{i1}^{(1)}| + \frac{|a_{i1}^{(1)}|}{|a_{11}^{(1)}|} \left(|a_{11}^{(1)}| - |a_{1i}^{(1)}| \right) = |a_{ii}^{(1)}| - \frac{|a_{i1}^{(1)}| |a_{1i}^{(1)}|}{|a_{11}^{(1)}|}.$$

The reverse triangle inequality implies

$$|a_{ii}^{(1)}| - \frac{|a_{i1}^{(1)}| |a_{1i}^{(1)}|}{|a_{11}^{(1)}|} \leq \left| a_{ii}^{(1)} - \frac{|a_{i1}^{(1)}| |a_{1i}^{(1)}|}{|a_{11}^{(1)}|} \right| = |a_{ii}^{(2)}|,$$

which gives

$$\sum_{j=2, j \neq i}^n |a_{ij}^{(2)}| < |a_{ii}^{(2)}|.$$

This establish the strict diagonal dominance for rows $2, \dots, n$. But the first row of $\mathbf{A}^{(2)}$ and \mathbf{A} are the same, so $\mathbf{A}^{(2)}$ is strictly diagonally dominant.

We can continue this process inductively and see that the result holds. \square

Definition 2.8. A matrix \mathbf{A} is said to be *positive definite* if $\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$ for all $\mathbf{x} \neq \mathbf{0}$. If the matrix is symmetric then it is referred to as *symmetric positive definite*.

For the next few theorems and corollaries, we will not be presenting the proofs can be found in Linear Algebra books. If you are interested, you can refer to [6].

Theorem 2.9. (Necessary Conditions for Symmetric Positive Definite) *If \mathbf{A} is a $n \times n$ symmetric positive definite matrix then*

1. \mathbf{A} has an inverse.
2. $a_{ii} > 0$ for each $i = 1, 2, \dots, n$.
3. $\max_{1 \leq k, j \leq n} |a_{kj}| \leq \max_{1 \leq i \leq n} |a_{ii}|$
4. $(a_{ij})^2 < a_{ii}a_{jj}$ for each $i \neq j$.

These conditions are only necessary conditions. For sufficient and necessary condition we introduce the notion of leading principal sub-matrix.

Definition 2.10. *A leading principal sub-matrix of a matrix \mathbf{A} is a matrix of the form*

$$\mathbf{A}_k = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1k} \\ a_{21} & a_{22} & \cdots & a_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ a_{k1} & a_{k2} & \cdots & a_{kk} \end{bmatrix},$$

for some $k = 1, 2, \dots, n$.

Theorem 2.11. (Necessary and Sufficient Condition for Symmetric Positive Definite) *A symmetric matrix \mathbf{A} is symmetric positive definite if and only if its leading principal sub-matrices have a positive determinant.*

Example: Consider the matrix

$$\mathbf{A} = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}.$$

It has three principal sub-matrix, \mathbf{A}_1 , \mathbf{A}_2 , and \mathbf{A}_3 each having determinants as 2, 3, and 4 respectively.

Theorem 2.12. *A symmetric matrix \mathbf{A} is symmetric positive definite if and only if Gaussian elimination without row interchanges can be performed on the linear system $\mathbf{Ax} = \mathbf{b}$ with all positive pivot element. Moreover in this case, the computations are stable with respect to the growth of round-off errors.*

Again we are not interested in the proof of the above theorem but rather certain corollaries that come while proving this theorem.

Corollary 2.13. (LDL^T Factorisation) *The matrix \mathbf{A} is symmetric positive definite if and only if \mathbf{A} can be factored in the form \mathbf{LDL}^T where \mathbf{L} is a unit lower triangular matrix and \mathbf{D} is a diagonal matrix with positive diagonal entries.*

Corollary 2.2.2 has a counterpart in case we have \mathbf{A} as a symmetric matrix.

Corollary 2.14. *Let \mathbf{A} be a symmetric matrix for which Gaussian elimination can be applied without row interchange. Then \mathbf{A} can be factored into \mathbf{LDL}^\top , where \mathbf{L} is lower unit triangular matrix with ones on the diagonal and \mathbf{D} is the diagonal matrix with $a_{11}^{(1)}, a_{22}^{(2)}, \dots, a_{nn}^{(n)}$.*

The algorithm to compute \mathbf{LDL}^\top is presented in Algorithm 9.

Algorithm 9 \mathbf{LDL}^\top Decomposition

Given: Symmetric matrix \mathbf{A} of size $n \times n$.

Find: Matrix \mathbf{L} (with unit diagonal) and \mathbf{D} such that $\mathbf{A} = \mathbf{LDL}^\top$.

Step 1: Initialize Matrices

Initialize \mathbf{L} as an identity matrix of size $n \times n$.

Initialize \mathbf{D} as a zero matrix of size $n \times n$.

Step 2: Compute \mathbf{D} and \mathbf{L}

$\mathbf{D}_{11} = \mathbf{A}_{11}$

for $i = 1$ **to** n **do**

Step 2.1: Compute \mathbf{D}

if $i \neq 1$ **then**

$\text{sum} = 0$

for $j = 1$ **to** $i - 1$ **do**

$\text{sum} = \text{sum} + \mathbf{D}_{jj} \mathbf{L}_{ij}^2$

end for

$\mathbf{D}_{ii} = \mathbf{A}_{ii} - \text{sum}$

end if

Step 2.2: Compute \mathbf{L}

for $j = i + 1$ **to** n **do**

$\text{sum} = 0$

if $i \neq 1$ **then**

for $k = 1$ **to** $i - 1$ **do**

$\text{sum} = \text{sum} + \mathbf{D}_{kk} \mathbf{L}_{ik} \mathbf{L}_{jk}$

end for

end if

$\mathbf{L}_{ji} = \frac{\mathbf{A}_{ji} - \text{sum}}{\mathbf{D}_{ii}}$

end for

end for

return \mathbf{L}, \mathbf{D}

Algorithm 9 is based on the computation of individual entries of \mathbf{L} and \mathbf{D} . Let us take

an example of how these entries actually look like or to be more precise how this algorithm is created.

Example: Let \mathbf{A} be a 3×3 symmetric matrix having LDL^\top decomposition. Then

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ \ell_{21} & 1 & 0 \\ \ell_{31} & \ell_{32} & 1 \end{bmatrix} \begin{bmatrix} d_{11} & 0 & 0 \\ 0 & d_{22} & 0 \\ 0 & 0 & d_{33} \end{bmatrix} \begin{bmatrix} 1 & \ell_{21} & \ell_{31} \\ 0 & 1 & \ell_{32} \\ 0 & 0 & 1 \end{bmatrix} \\ &= \begin{bmatrix} d_{11} & d_{11}\ell_{21} & d_{11}\ell_{31} \\ d_{11}\ell_{21} & d_{22} + d_{11}\ell_{21}^2 & d_{22}\ell_{32} + d_{11}\ell_{21}\ell_{31} \\ d_{11}\ell_{31} & d_{22}\ell_{32} + d_{11}\ell_{21}\ell_{31} & d_{11}\ell_{31}^2 + d_{22}\ell_{32}^2 + d_{33} \end{bmatrix}. \end{aligned}$$

We notice that $d_{11} = a_{11}$ and $\ell_{i1} = a_{i1}/d_{11}$ for $i = 2, 3$. After this we can compute d_{22} and then ℓ_{32} . Finally we compute d_{33} . Same process can be extended to a $n \times n$ matrix.

2.2.3 Cholesky Decomposition

From Theorem 2.12 we have another corollary related to symmetric positive definite matrix which gives another decomposition.

Corollary 2.15. (Cholesky Decomposition) *The matrix \mathbf{A} is symmetric positive definite if and only if \mathbf{A} can be factored in the form \mathbf{LL}^\top where \mathbf{L} is a lower triangular matrix with non-zero diagonal entries.*

The Cholesky decomposition was discovered by André-Louis Cholesky who was a French military officer (along with being a mathematician).



Figure 2.5: André-Louis Cholesky: 15 October 1875-31 August 1918.

The algorithm for the Cholesky decomposition can be found in Algorithm 10.

Algorithm 10 Cholesky Decomposition

Given: Symmetric positive definite matrix \mathbf{A} of size $n \times n$.

Find: Matrices \mathbf{L} such that $\mathbf{A} = \mathbf{L}\mathbf{L}^T$.

Step 1: Initialize Matrix

Initialize \mathbf{L} as a zero matrix of size $n \times n$.

Step 2: Compute L

$\mathbf{L}_{11} = \sqrt{\mathbf{A}_{11}}$

for $j = 2$ to n do

$\mathbf{L}_{j1} = \frac{\mathbf{A}_{j1}}{\mathbf{L}_{11}}$

end for

for $i = 2$ to n do

Step 2.1: Compute \mathbf{L}_{ii}

sum = 0

for $k = 1$ to $i - 1$ do

 sum = sum + \mathbf{L}_{ik}^2

end for

$\mathbf{L}_{ii} = \sqrt{\mathbf{A}_{ii} - \text{sum}}$

for $j = i + 1$ to n do

Step 2.2: Compute \mathbf{L}_{ji}

sum = 0

for $k = 1$ to $i - 1$ do

 sum = sum + $\mathbf{L}_{jk}\mathbf{L}_{ik}$

end for

$\mathbf{L}_{ji} = \frac{1}{\mathbf{L}_{ii}} (\mathbf{A}_{ji} - \text{sum})$

end for

end for

return \mathbf{L}

Algorithm 10 is based on the computation of individual entries of \mathbf{L} . Let us take an example of how these entries actually look like.

Example: Let \mathbf{A} be a 3×3 symmetric positive definite matrix having a $\mathbf{L}\mathbf{L}^T$ decomposition. Then

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} \ell_{11} & 0 & 0 \\ \ell_{21} & \ell_{22} & 0 \\ \ell_{31} & \ell_{32} & \ell_{33} \end{bmatrix} \begin{bmatrix} \ell_{11} & \ell_{21} & \ell_{31} \\ 0 & \ell_{22} & \ell_{32} \\ 0 & 0 & \ell_{33} \end{bmatrix} \\ &= \begin{bmatrix} \ell_{11}^2 & \ell_{11}\ell_{21} & \ell_{11}\ell_{31} \\ \ell_{11}\ell_{21} & \ell_{21}^2 + \ell_{22}^2 & \ell_{21}\ell_{31} + \ell_{22}\ell_{32} \\ \ell_{11}\ell_{31} & \ell_{21}\ell_{31} + \ell_{22}\ell_{32} & \ell_{31}^2 + \ell_{32}^2 + \ell_{33}^2 \end{bmatrix}. \end{aligned}$$

We notice that $l_{11} = \sqrt{a_{11}}$ and $l_{i1} = a_{i1}/l_{11}$ for $i = 2, 3$. After this we can compute l_{22} and then l_{32} . Finally we compute l_{33} . Same process can be extended to a $n \times n$ matrix.

Until now we have not discussed about the computational complexity of any of the three factorisation methods. Table 2.1 gives the applicability, computational complexity, and advantages of the three factorisation methods. For brevity, we would not derive the computational complexity for the methods but interested students can try it for their own.

Method	Applicability	Advantages	Computational Complexity	
			M/D	A/S
LU	General square matrices	Works for any matrix but requires pivoting	$\frac{n^3}{3} - \frac{n}{3}$	$\frac{n^3}{3} - \frac{n^2}{2} + \frac{n}{6}$
LDL ^T	Symmetric matrices	More stable than LU, reduces storage, avoids pivoting	$\frac{n^3}{6} + n^2 - \frac{7n}{6}$	$\frac{n^3}{6} - \frac{n}{6}$
Cholesky	Symmetric positive definite matrices	Fastest and most efficient, lowest computation cost	$\frac{n^3}{6} + \frac{n^2}{2} - \frac{2n}{3}$	$\frac{n^3}{6} - \frac{n}{6}$

Table 2.1: Applicability, Computational Complexity, and Advantages for LU, LDL^T, and Cholesky Decomposition. M/D: Multiplication and Division, A/S: Addition and Subtraction.

We notice that the Cholesky decomposition requires the least number of operations while factorisation but it can be a little misleading as it requires extracting n square roots. However the computation of square root is a linear factor of n and will decrease significantly as n increases.

2.3 Iterative Methods

Root-finding methods are a class of iterative methods that we are aware of. In this part of the chapter, we will translate these ideas into a system of equations. Before delving into iterative methods for a system of equations, we need to find a way to measure the distance between n -dimensional column vectors. This will help us determine the sequence of vectors that converge to the solution of the system.

Definition 2.16. A *vector norm* on \mathbb{R}^n is a function, $\|\cdot\|$ from \mathbb{R}^n into \mathbb{R} with the following properties:

1. $\|\mathbf{x}\| \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$.
2. $\|\mathbf{x}\| = 0$ if and only if $\mathbf{x} = \mathbf{0}$.
3. $\|\alpha\mathbf{x}\| = |\alpha|\|\mathbf{x}\|$ for all $\alpha \in \mathbb{R}$ and $\mathbf{x} \in \mathbb{R}^n$.
4. $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$.

Definition 2.17. The ℓ_2 and the ℓ_∞ norm for the vector $\mathbf{x} = (x_1, x_2, \dots, x_n)$ are defined by

$$\|\mathbf{x}\|_2^2 = \sum_{i=1}^n x_i^2 \quad \text{and} \quad \|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|.$$

If we define a unit ball in \mathbb{R}^2 using these norms, then they are given by $\|\mathbf{x}\|_2 \leq 1$ which is an unit disc centred at $(0, 0)$ and $\|\mathbf{x}\|_\infty \leq 1$ which is a square.

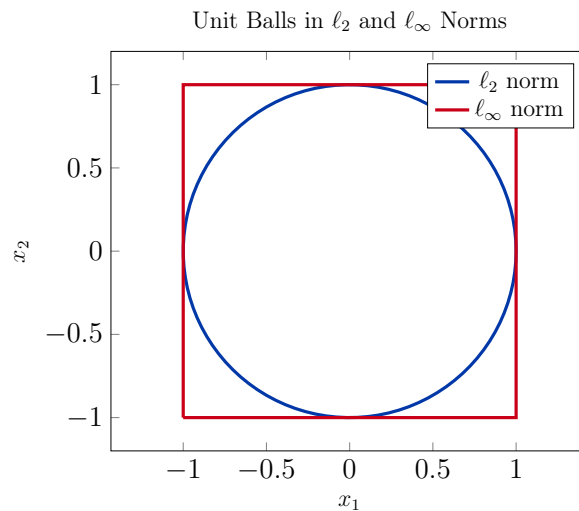


Figure 2.6: Unit balls in ℓ_2 and ℓ_∞ norm.

A fundamental property of these norms that is widely used is the Cauchy-Schwarz inequality.

Theorem 2.18. (Cauchy-Schwarz Inequality) For each $\mathbf{x} = (x_1, x_2, \dots, x_n)^\top$ and $\mathbf{y} = (y_1, y_2, \dots, y_n)^\top$ in \mathbb{R}^n

$$\mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i \leq \left\{ \sum_{i=1}^n x_i^2 \right\}^{1/2} \left\{ \sum_{i=1}^n y_i^2 \right\}^{1/2} = \|\mathbf{x}\|_2 \|\mathbf{y}\|_2.$$

Proof. The result is immediate if $\mathbf{x} = \mathbf{0}$ or $\mathbf{y} = \mathbf{0}$. Suppose $\mathbf{x} \neq \mathbf{0}$ and $\mathbf{y} \neq \mathbf{0}$. Now, note that for each $\lambda \in \mathbb{R}$ we have

$$0 \leq \|\mathbf{x} - \lambda \mathbf{y}\|_2^2 = \sum_{i=1}^n (x_i - \lambda y_i)^2 = \sum_{i=1}^n x_i^2 - 2\lambda \sum_{i=1}^n x_i y_i + \lambda^2 \sum_{i=1}^n y_i^2,$$

so that

$$2\lambda \sum_{i=1}^n x_i y_i \leq \|\mathbf{x}\|_2^2 + \lambda^2 \|\mathbf{y}\|_2^2.$$

As, $\|\mathbf{x}\|_2 > 0$ and $\|\mathbf{y}\|_2 > 0$ so we let $\lambda = \frac{\|\mathbf{x}\|_2}{\|\mathbf{y}\|_2}$, which gives

$$\frac{2\|\mathbf{x}\|_2}{\|\mathbf{y}\|_2} \left(\sum_{i=1}^n x_i y_i \right) \leq \|\mathbf{x}\|_2^2 + \frac{\|\mathbf{x}\|_2^2}{\|\mathbf{y}\|_2^2} \|\mathbf{y}\|_2^2 = 2\|\mathbf{x}\|_2^2,$$

which, after simplification, gives us the result. \square

The norm of a vector measures the distance between an arbitrary vector and the zero vector. We define the distance between two vectors as

$$\|\mathbf{x} - \mathbf{y}\|_2 = \left\{ \sum_{i=1}^n (x_i - y_i)^2 \right\}^{1/2} \quad \text{and} \quad \|\mathbf{x} - \mathbf{y}\|_\infty = \max_{1 \leq i \leq n} |x_i - y_i|.$$

Now, we define the convergence of a sequence of vectors in \mathbb{R}^n .

Definition 2.19. A sequence $\{\mathbf{x}^{(k)}\}_{k=1}^\infty$ of vectors in \mathbb{R}^n is said to converge to \mathbf{x} with respect to $\|\cdot\|$ if given for any $\varepsilon > 0$ there exist a $N(\varepsilon)$ such that

$$\|\mathbf{x}^{(k)} - \mathbf{x}\| < \varepsilon \quad \forall k \geq N(\varepsilon).$$

Next, we present the result regarding the equivalence of norms.

Theorem 2.20. For each $\mathbf{x} \in \mathbb{R}^n$

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2 \leq \sqrt{n} \|\mathbf{x}\|_\infty.$$

Proof. Let x_j be the coordinate such that $\|\mathbf{x}\|_\infty = |x_j| = \max_{1 \leq i \leq n} |x_i|$.

Now,

$$\|\mathbf{x}\|_\infty^2 = |x_j|^2 = x_j^2 \leq \sum_{i=1}^n x_i^2 = \|\mathbf{x}\|_2^2.$$

Similarly,

$$\|\mathbf{x}\|_2^2 = \sum_{i=1}^n x_i^2 \leq \sum_{i=1}^n x_j^2 \leq x_j^2 n \leq \|\mathbf{x}\|_\infty^2 n.$$

Hence, $\|\mathbf{x}\|_2 \leq \sqrt{n} \|\mathbf{x}\|_\infty$. \square

Similar to vector norms, we also have matrix norms. The measure given to a matrix under a natural norm describes how the matrix stretches unit vectors relative to that norm. The maximum stretch is the norm of the matrix. The definition of the matrix norm is similar to that of the vector norm.

Theorem 2.21. If $\|\cdot\|$ is a vector norm in \mathbb{R}^n , then

$$\|\mathbf{A}\| = \max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|$$

is a matrix norm.

Note that we have an abuse of notation here; we denote $\|\cdot\|$ to show both the vector and the matrix norm.

Matrix norms defined by vector norms are called *natural* or *induced* matrix norm. We can also write the natural matrix norms as

$$\|\mathbf{A}\| = \max_{\mathbf{y} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{y}\|}{\|\mathbf{y}\|},$$

as $\mathbf{y}/\|\mathbf{y}\|$ is a unit vector.

Corollary 2.22. For any vector $\mathbf{y} \neq \mathbf{0}$, matrix \mathbf{A} and any natural norm $\|\cdot\|$, we have

$$\|\mathbf{A}\mathbf{y}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{y}\|.$$

The matrix norm that we consider are the ∞ norm, i.e., $\|\mathbf{A}\|_\infty = \max_{\|\mathbf{x}\|_\infty=1} \|\mathbf{A}\mathbf{x}\|_\infty$ and the ℓ_2 norm, i.e., $\|\mathbf{A}\|_\infty = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2$.

Lastly, we define the $\|\cdot\|_\infty$ norm of a matrix.

Theorem 2.23. If $\mathbf{A} = \{a_{ij}\}_{i,j=1}^n$ is a $n \times n$ matrix then

$$\|\mathbf{A}\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|.$$

A square matrix \mathbf{A} takes the set of n -dimensional vectors into itself, which gives a linear function from \mathbb{R}^n to \mathbb{R}^n . After this transformation, certain vectors might be parallel to the original vector, i.e., \mathbf{x} is parallel to $\mathbf{A}\mathbf{x}$. It might be stretched, shrunk, or remains unchanged. The magnitude with which it stretches or shrinks is called the eigen or characteristic value. But why do we care about these eigenvalues? There is a close relation between these eigenvalues and the convergence of the iterative methods.

Definition 2.24. If \mathbf{A} is a square matrix, then the *characteristic polynomial* of \mathbf{A} is defined by

$$p(\lambda) = \det(\mathbf{A} - \lambda\mathbf{I}).$$

Definition 2.25. If $p(\lambda)$ is the characteristic polynomial of the matrix \mathbf{A} , the zeros of p are *eigenvalues* or *characteristic values* of the matrix \mathbf{A} . If λ is an eigenvalue of \mathbf{A} and $\mathbf{x} \neq \mathbf{0}$ satisfies $(\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = \mathbf{0}$ then \mathbf{x} is an *eigenvector* or characteristic vector of \mathbf{A} corresponding the λ .

Note that if \mathbf{x} is an eigenvector of \mathbf{A} associated with λ and $\alpha \in \mathbb{R} \setminus \{0\}$ then $\alpha\mathbf{x}$ is an eigenvector since

$$\mathbf{A}(\alpha\mathbf{x}) = \alpha(\mathbf{A}\mathbf{x}) = \alpha(\lambda\mathbf{x}) = \lambda(\alpha\mathbf{x}).$$

As an immediate consequence of this is \mathbf{x} is an eigenvector then we can choose $\alpha = \pm\|\mathbf{x}\|^{-1}$, which would make $\alpha\mathbf{x}$ an eigenvector with norm one. So far any eigenvalue and any vector norm we have eigenvectors with norm one.

Definition 2.26. The *spectral radius* $\rho(\mathbf{A})$ of a matrix is defined by

$$\rho(\mathbf{A}) = \max\{|\lambda|\},$$

where λ is an eigenvalue of \mathbf{A} . For $\lambda(:= \alpha + i\beta) \in \mathbb{C}$, $|\lambda| = (\alpha^2 + \beta^2)^{1/2}$.

Next, we have a relation between the spectral radius and the matrix norm.

Theorem 2.27. If \mathbf{A} is a $n \times n$ matrix then:

1. $\|\mathbf{A}\|_2 = [\rho(\mathbf{A}^\top \mathbf{A})]^{1/2}$.
2. $\rho(\mathbf{A}) \leq \|\mathbf{A}\|$, for any natural norm $\|\cdot\|$.

If \mathbf{A} is symmetric then $\|\mathbf{A}\|_2 = \rho(\mathbf{A})$. Apart from spectral radius, another important property to study is how a matrix's power behaves.

Definition 2.28. We call $n \times n$ matrix \mathbf{A} *convergent* if

$$\lim_{k \rightarrow \infty} (\mathbf{A}^k)_{ij} = 0 \quad \forall i = 1, 2, \dots, n \quad \text{and} \quad j = 1, 2, \dots, n.$$

Now, convergent matrices have a special connection with the spectral radius.

Theorem 2.29. The following statements are equivalent:

1. \mathbf{A} is a convergent matrix.
2. $\lim_{n \rightarrow \infty} \|\mathbf{A}^n\| = 0$ for some natural norm.
3. $\lim_{n \rightarrow \infty} \|\mathbf{A}^n\| = 0$ for all natural norms.
4. $\rho(\mathbf{A}) < 1$.
5. $\lim_{n \rightarrow \infty} \mathbf{A}^n \mathbf{x} = \mathbf{0}$ for all \mathbf{x} .

2.3.1 Jacobi Method

After having a quick glance at the basics of linear algebra, we move back toward the domain of numerical analysis.

An iterative technique to solve the $n \times n$ linear system $\mathbf{Ax} = \mathbf{b}$ starts with an initial approximation $\mathbf{x}^{(0)}$ to the solution \mathbf{x} and generates a sequence of vectors $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$ such that $\mathbf{x}^{(k)} \rightarrow \mathbf{x}$ as $k \rightarrow \infty$.

Carl Gustav Jacob Jacobi was a German mathematician who proposed the Jacobi eigenvalue algorithm, an iterative method for calculating the eigenvalues and eigenvectors of a real symmetric matrix. The *Jacobi method* that we study is the stripped-down version of this algorithm only.

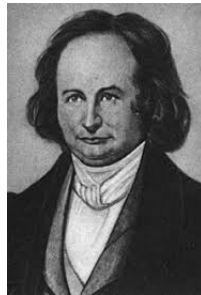


Figure 2.7: Carl Gustav Jacob Jacobi: 10 December 1804-18 February 1851

The Jacobi iterative method is obtained by solving the i^{th} equation in $\mathbf{Ax} = \mathbf{b}$ for x_i to obtain

$$x_i = \sum_{j=1, j \neq i}^n \left(-\frac{a_{ij}x_j}{a_{ii}} \right) + \frac{b_i}{a_{ii}} \quad \text{for } i = 1, 2, \dots, n. \quad (2.6)$$

For each $k \geq 1$ we generate the components $x_i^{(k)}$ of $\mathbf{x}^{(k)}$ from the components of $\mathbf{x}^{(k-1)}$ by

$$x_i^{(k)} = \frac{1}{a_{ii}} \left[\sum_{j=1, j \neq i}^n \left(-a_{ij}x_j^{(k-1)} \right) + b_i \right] \quad \text{for } i = 1, 2, \dots, n. \quad (2.7)$$

Example: Say we have the system of equations of the form

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= b_1 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n &= b_2 \\ &\vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n &= b_n. \end{aligned}$$

Suppose $\mathbf{x}^{(0)}$ is the initial iterate, then the first iterative solution is given by

$$\begin{aligned} x_1^{(1)} &= \frac{1}{a_{11}} \left(b_1 - \left(a_{12}x_2^{(0)} + \cdots + a_{1n}x_n^{(0)} \right) \right) \\ x_2^{(1)} &= \frac{1}{a_{22}} \left(b_2 - \left(a_{21}x_1^{(0)} + \cdots + a_{2n}x_n^{(0)} \right) \right) \\ &\vdots \\ x_n^{(1)} &= \frac{1}{a_{nn}} \left(b_n - \left(a_{n1}x_1^{(0)} + \cdots + a_{n,n-1}x_{n-1}^{(0)} \right) \right), \end{aligned}$$

and similarly, we compute for $k \geq 1$.

In general, iterative techniques for solving linear systems of equations involve a process that converts $\mathbf{Ax} = \mathbf{b}$ into an equivalent system $\mathbf{x} = \mathbf{T}\mathbf{x} + \mathbf{c}$ for some fixed-matrix \mathbf{T} and vector \mathbf{c} . Once the initial approximation is selected say $\mathbf{x}^{(0)}$ we get

$$\mathbf{x}^{(k)} = \mathbf{T}\mathbf{x}^{(k-1)} + \mathbf{c} \quad \text{for each } k = 1, 2, \dots$$

We can have an equivalent formulation for the Jacobi method by splitting \mathbf{A} as

$$\mathbf{A} = \mathbf{D} - \mathbf{L} - \mathbf{U},$$

where \mathbf{L} is strict lower triangular part of \mathbf{A} , \mathbf{U} is the strict upper triangular part of \mathbf{A} , and \mathbf{D} is the diagonal. Say,

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix}.$$

Then,

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

Then we can re-write $\mathbf{Ax} = \mathbf{b}$ as

$$\mathbf{D}\mathbf{x} = (\mathbf{L} + \mathbf{U})\mathbf{x} + \mathbf{b},$$

and if \mathbf{D}^{-1} exist, $\mathbf{x} = \mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})\mathbf{x} + \mathbf{D}^{-1}\mathbf{b}$.

Then, the Jacobi iterative is given by

$$\mathbf{x}^{(k)} = \mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})\mathbf{x}^{(k-1)} + \mathbf{D}^{-1}\mathbf{b} \quad \text{for } k \geq 1. \quad (2.8)$$

Denoting $\mathbf{T}_J = \mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})$ and $\mathbf{c}_J = \mathbf{D}^{-1}\mathbf{b}$ then we get the Jacobi iteration as

$$\mathbf{x}^{(k)} = \mathbf{T}_J\mathbf{x}^{(k-1)} + \mathbf{c}_J \quad \text{for } k \geq 1. \quad (2.9)$$

We need $a_{ii} \neq 0$ for each $i = 1, 2, \dots, n$. If one of the $a_{ii} = 0$ and the system is not singular, then the equations can be reordered so that no a_{ii} is zero.

The algorithm for the Jacobi method is provided in Algorithm 11.

Algorithm 11 Jacobi Iteration

Given: Matrix \mathbf{A} with non-zero pivots, right hand side \mathbf{b} , dimension n , `max_iterations`, and `tolerance`.

Find: Solution \mathbf{x} .

Step 1: Jacobi Iterations

Initialize $\mathbf{x}^{\text{old}} = \mathbf{0}$

for $k = 1$ **to** `max_iterations` **do**

for $i = 1$ **to** n **do**

$\text{sum} = \mathbf{b}_i$

for $j = 1$ **to** n **do**

if $j \neq i$ **then**

$\text{sum} = \text{sum} - \mathbf{A}_{ij}\mathbf{x}_j^{\text{old}}$

end if

end for

$\mathbf{x}_i = \frac{\text{sum}}{\mathbf{A}_{ii}}$

end for

$\text{Error} = \|\mathbf{x} - \mathbf{x}^{\text{old}}\|_{\infty}$

if $\text{Error} < \text{tolerance}$ **then**

Output(“Convergence reached”)

break

end if

$\mathbf{x}^{\text{old}} = \mathbf{x}$

end for

if $k == \text{max_iterations}$ **then**

Output(“Maximum Number of iterations reached”)

end if

return \mathbf{x}

2.3.2 Gauss-Seidel Method

In the Jacobi method we require all the components of $\mathbf{x}^{(k-1)}$ are used to compute the components $x_i^{(k)}$ of $\mathbf{x}^{(k)}$. But, for $i > 1$, the component $x_1^{(k)}, \dots, x_{i-1}^{(k)}$ of $\mathbf{x}^{(k)}$ have already being computed. If we use these values, then it is expected to give better approximations to the actual solutions than $x_1^{(k-1)}, \dots, x_{i-1}^{(k-1)}$.

Then, it is reasonable to compute $x_i^{(k)}$ using the most recently calculated values.

$$x_i^{(k)} = \frac{1}{a_{ii}} \left[- \sum_{j=1}^{i-1} a_{ij}x_j^{(k)} - \sum_{j=i+1}^n a_{ij}x_j^{(k-1)} + b_i \right], \quad (2.10)$$

for $i = 1, 2, \dots, n$. This is called the *Gauss Seidel method*. Gauss initially developed the

concept in the mid-1820s; it was only published and fully detailed by Seidel in 1874 through a private letter from Gauss to his student Gerling, making the method primarily attributed to both mathematicians.



Figure 2.8: Philipp Ludwig von Seidel: 24 October 1821-13 August 1896

Example: Say we have the system of equations of the form

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= b_1 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n &= b_2 \\ &\vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n &= b_n. \end{aligned}$$

Suppose $\mathbf{x}^{(0)}$ is the initial iterate, then the first iterative solution is given by

$$\begin{aligned} x_1^{(1)} &= \frac{1}{a_{11}} \left(b_1 - \left(a_{12}x_2^{(0)} + a_{13}x_3^{(0)} + \cdots + a_{1n}x_n^{(0)} \right) \right) \\ x_2^{(1)} &= \frac{1}{a_{22}} \left(b_2 - \left(a_{21}x_1^{(1)} + a_{23}x_3^{(0)} + \cdots + a_{2n}x_n^{(0)} \right) \right) \\ &\vdots \\ x_n^{(1)} &= \frac{1}{a_{nn}} \left(b_n - \left(a_{n1}x_1^{(1)} + a_{n2}x_2^{(1)} + \cdots + a_{n,n-1}x_{n-1}^{(1)} \right) \right), \end{aligned}$$

and similarly, we compute for $k \geq 1$.

To write the Gauss-Seidel method in matrix form, we multiply Eq. (2.10) with a_{ii} and collect the k^{th} iterate term to get

$$a_{i1}x_1^{(k)} + a_{i2}x_2^{(k)} + \cdots + a_{ii}x_i^{(k)} = -a_{i,i+1}x_{i+1}^{(k-1)} - \cdots - a_{in}x_n^{(k-1)} + b_i,$$

for $i = 1, 2, \dots, n$. Then

$$\begin{aligned} a_{11}x_1^{(k)} &= -a_{12}x_2^{(k-1)} - a_{13}x_3^{(k-1)} - \cdots - a_{1n}x_n^{(k-1)} + b_1, \\ a_{21}x_1^{(k)} + a_{22}x_2^{(k)} &= -a_{23}x_3^{(k-1)} - \cdots - a_{2n}x_n^{(k-1)} + b_2, \\ &\vdots \\ a_{n1}x_1^{(k)} + a_{n2}x_2^{(k)} + \cdots + a_{nn}x_n^{(k)} &= b_n. \end{aligned}$$

Then, we can write this system as

$$(\mathbf{D} - \mathbf{L})\mathbf{x}^{(k)} = \mathbf{U}\mathbf{x}^{(k-1)} + \mathbf{b},$$

and $\mathbf{x}^{(k)} = (\mathbf{D} - \mathbf{L})^{-1} \mathbf{U}\mathbf{x}^{(k-1)} + (\mathbf{D} - \mathbf{L})^{-1} \mathbf{b}$ for $k \geq 1$, where \mathbf{D} , \mathbf{L} , and \mathbf{U} are defined in the same way as Jacobi method. Then denoting $\mathbf{T}_{\text{GS}} = (\mathbf{D} - \mathbf{L})^{-1} \mathbf{U}$ and $\mathbf{c}_{\text{GS}} = (\mathbf{D} - \mathbf{L})^{-1} \mathbf{b}$ we get the Gauss-Seidel method as

$$\mathbf{x}^{(k)} = \mathbf{T}_{\text{GS}}\mathbf{x}^{(k-1)} + \mathbf{c}_{\text{GS}}.$$

Now, $\mathbf{D} - \mathbf{L}$ is non singular if and only if $a_{ii} \neq 0$ for all $i = 1, 2, \dots, n$.

It appears that the Gauss-Seidel method is always a better approximation to the Jacobi method, which is “mostly” true, but we have cases where this might not hold.

The algorithm for the Gauss-Seidel method is provided in Algorithm 12.

Algorithm 12 Gauss-Seidel Iteration

Given: Matrix \mathbf{A} with non-zero pivots, right hand side \mathbf{b} , dimension n , `max_iterations`, and `tolerance`.

Find: Solution \mathbf{x} .

Step 1: Gauss-Seidel Iterations

Initialize $\mathbf{x}^{\text{old}} = \mathbf{0}$

for $k = 1$ to `max_iterations` do

 for $i = 1$ to n do

 sum = \mathbf{b}_i

 for $j = 1$ to n do

 if $j < i$ then

 sum = sum - $\mathbf{A}_{ij}\mathbf{x}_j$

 else if $i < j$ then

 sum = sum - $\mathbf{A}_{ij}\mathbf{x}_j^{\text{old}}$

 end if

 end for

$\mathbf{x}_i = \frac{\text{sum}}{\mathbf{A}_{ii}}$

 end for

 Error = $\|\mathbf{x} - \mathbf{x}^{\text{old}}\|_{\infty}$

 if Error < `tolerance` then

Output(“Convergence reached”)

 break

 end if

$\mathbf{x}^{\text{old}} = \mathbf{x}$

end for

if $k == \text{max_iterations}$ then

Output(“Maximum Number of iterations reached”)

end if

return \mathbf{x}

General Iteration Matrices

We need to analyze the formula

$$\mathbf{x}^{(k)} = \mathbf{T}\mathbf{x}^{(k-1)} + \mathbf{c} \quad \text{for } k \geq 1,$$

to study the convergence of general iteration techniques where $\mathbf{x}^{(0)}$ is arbitrary.

Lemma 2.30. *If the spectral radius $\rho(\mathbf{T}) < 1$ then $(\mathbf{I} - \mathbf{T})^{-1}$ exists, and*

$$(\mathbf{I} - \mathbf{T})^{-1} = \mathbf{I} + \mathbf{T} + \mathbf{T}^2 + \dots = \sum_{j=0}^{\infty} \mathbf{T}^j.$$

Proof. Now, let λ be an eigenvalue of \mathbf{T} with eigenvector \mathbf{x} then

$$\mathbf{T}\mathbf{x} = \lambda\mathbf{x} \iff (\mathbf{I} - \mathbf{T})\mathbf{x} = (1 - \lambda)\mathbf{x}.$$

Hence, λ is an eigenvalue of \mathbf{T} if and only if $1 - \lambda$ is a eigenvalue of $\mathbf{I} - \mathbf{T}$.

However, by the definition of spectral radius $|\lambda| \leq \rho(\mathbf{T}) < 1$, so $\lambda = 1$ is not an eigenvalue of \mathbf{T} which implies 0 is not an eigenvalue of $\mathbf{I} - \mathbf{T}$.

Hence, $\mathbf{I} - \mathbf{T}$ is invertible. Let $\mathbf{S}_m = \mathbf{I} + \mathbf{T} + \mathbf{T}^2 + \dots + \mathbf{T}^m$, then

$$(\mathbf{I} - \mathbf{T})\mathbf{S}_m = \mathbf{I} + \mathbf{T} + \mathbf{T}^2 + \dots + \mathbf{T}^m - \mathbf{T} - \mathbf{T}^2 - \dots - \mathbf{T}^{m+1} = \mathbf{I} - \mathbf{T}^{m+1}.$$

As $\rho(\mathbf{T}) < 1$ then by Theorem 2.29 we have \mathbf{T} is convergent. Again using Theorem 2.29 we get

$$\lim_{m \rightarrow \infty} (\mathbf{I} - \mathbf{T})\mathbf{S}_m = \lim_{m \rightarrow \infty} (\mathbf{I} - \mathbf{T}^{m+1}) = \mathbf{I}.$$

Thus

$$(\mathbf{I} - \mathbf{T})^{-1} = \lim_{m \rightarrow \infty} \mathbf{S}_m = \sum_{j=0}^{\infty} \mathbf{T}^j.$$

□

Theorem 2.31. *For any $\mathbf{x}^{(0)} \in \mathbb{R}^n$ the sequence $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$ defined by*

$$\mathbf{x}^{(k)} = \mathbf{T}\mathbf{x}^{(k-1)} + \mathbf{c} \quad \text{for each } k \geq 1,$$

converges to the unique solution $\mathbf{x} = \mathbf{T}\mathbf{x} + \mathbf{c}$ if and only if $\rho(\mathbf{T}) < 1$.

Proof. Let $\rho(\mathbf{T}) < 1$. Then

$$\begin{aligned} \mathbf{x}^{(k)} &= \mathbf{T}\mathbf{x}^{(k-1)} + \mathbf{c} \\ &= \mathbf{T}(\mathbf{T}\mathbf{x}^{(k-2)} + \mathbf{c}) + \mathbf{c} \\ &= \mathbf{T}^2\mathbf{x}^{(k-2)} + \mathbf{T}\mathbf{c} + \mathbf{c} \\ &\vdots \\ &= \mathbf{T}^k\mathbf{x}^{(0)} + (\mathbf{T}^{k-1} + \dots + \mathbf{I})\mathbf{c}. \end{aligned} \tag{2.11}$$

As $\rho(\mathbf{T}) < 1$ from Theorem 2.29 we get that \mathbf{T} is convergent and $\lim_{k \rightarrow \infty} \mathbf{T}^k \mathbf{x}^{(0)} = \mathbf{0}$.

In Eq. (2.11) passing the limit of $k \rightarrow \infty$, and then using the previous lemma, we get,

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \lim_{k \rightarrow \infty} \mathbf{T}^{(k)} \mathbf{x}^{(0)} + \left(\sum_{j=0}^{\infty} \mathbf{T}^j \right) \mathbf{c} = \mathbf{0} + (\mathbf{I} - \mathbf{T})^{-1} \mathbf{c}.$$

Hence, $\{\mathbf{x}^{(k)}\} \rightarrow \mathbf{x}$ as $k \rightarrow \infty$ and $\mathbf{x} = \mathbf{T}\mathbf{x} + \mathbf{c}$ is the unique limit.

Conversely we will show that for any $\mathbf{y} \in \mathbb{R}^n$ we have $\lim_{k \rightarrow \infty} \mathbf{T}^k \mathbf{y} = \mathbf{0}$ which is equivalent to $\rho(\mathbf{T}) < 1$.

Let $\mathbf{y} \in \mathbb{R}^n$ be arbitrary and \mathbf{x} be the unique solution to $\mathbf{x} = \mathbf{T}\mathbf{x} + \mathbf{c}$. Define $\mathbf{x}^{(0)} = \mathbf{x} - \mathbf{y}$ and for $k \geq 1$

$$\mathbf{x}^{(k)} = \mathbf{T}\mathbf{x}^{(k-1)} + \mathbf{c}.$$

Now, by the hypothesis $\{\mathbf{x}^{(k)}\} \rightarrow \mathbf{x}$. Also,

$$\mathbf{x} - \mathbf{x}^{(k)} = (\mathbf{T}\mathbf{x} + \mathbf{c}) - (\mathbf{T}\mathbf{x}^{(k-1)} + \mathbf{c}) = \mathbf{T}(\mathbf{x} - \mathbf{x}^{(k-1)}).$$

Inductively $\mathbf{x} - \mathbf{x}^{(k)} = \mathbf{T}(\mathbf{x} - \mathbf{x}^{(k-1)}) = \mathbf{T}^2(\mathbf{x} - \mathbf{x}^{(k-2)}) = \dots = \mathbf{T}^k(\mathbf{x} - \mathbf{x}^{(0)}) = \mathbf{T}^k \mathbf{y}$. Hence, $\lim_{k \rightarrow \infty} \mathbf{T}^k \mathbf{y} = \lim_{k \rightarrow \infty} \mathbf{T}^k(\mathbf{x} - \mathbf{x}^{(0)}) = \lim_{k \rightarrow \infty} (\mathbf{x} - \mathbf{x}^{(k)}) = \mathbf{0}$.

As $\mathbf{y} \in \mathbb{R}^n$ was arbitrary. By Theorem 2.29 we get that $\rho(\mathbf{T}) < 1$. □

Based on this theorem, a nice corollary bounds the error.

Corollary 2.32. *If $\|\mathbf{T}\| < 1$ for any natural norm and \mathbf{c} is a given vector, then the sequence $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$ defined by $\mathbf{x}^{(k)} = \mathbf{T}\mathbf{x}^{(k-1)} + \mathbf{c}$ converges, for any $\mathbf{x}^{(0)} \in \mathbb{R}^n$, to a vector $\mathbf{x} \in \mathbb{R}^n$ such that $\mathbf{x} = \mathbf{T}\mathbf{x} + \mathbf{c}$. Furthermore, the following error bounds hold:*

1. $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|\mathbf{T}\|^k \|\mathbf{x}^{(0)} - \mathbf{x}\|.$
2. $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \frac{\|\mathbf{T}\|^k}{1 - \|\mathbf{T}\|} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|.$

Proof. Let $\|\mathbf{T}\| < 1$. By Theorem 2.27, we have $\rho(\mathbf{T}) < 1$. Therefore, the method converges to the solution $\mathbf{x} = \mathbf{T}\mathbf{x} + \mathbf{c}$ by the previous theorem. We now prove the error bounds.

Starting from the iterative formula $\mathbf{x}^{(k)} = \mathbf{T}\mathbf{x}^{(k-1)} + \mathbf{c}$, we subtract $\mathbf{x} = \mathbf{T}\mathbf{x} + \mathbf{c}$ from both sides and taking the norm:

$$\begin{aligned} \|\mathbf{x}^{(k)} - \mathbf{x}\| &= \|\mathbf{T}\mathbf{x}^{(k-1)} + \mathbf{c} - \mathbf{T}\mathbf{x} - \mathbf{c}\| \\ &= \|\mathbf{T}(\mathbf{x}^{(k-1)} - \mathbf{x})\| \\ &\leq \|\mathbf{T}\| \cdot \|\mathbf{x}^{(k-1)} - \mathbf{x}\|. \end{aligned}$$

Applying this inequality recursively:

$$\begin{aligned} \|\mathbf{x}^{(k)} - \mathbf{x}\| &\leq \|\mathbf{T}\| \cdot \|\mathbf{x}^{(k-1)} - \mathbf{x}\| \\ &\leq \|\mathbf{T}\|^2 \cdot \|\mathbf{x}^{(k-2)} - \mathbf{x}\| \\ &\vdots \\ &\leq \|\mathbf{T}\|^k \cdot \|\mathbf{x}^{(0)} - \mathbf{x}\|. \end{aligned}$$

This proves the first error bound.

For the second error bound, consider the norm of the difference between successive iterates:

$$\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| = \|\mathbf{T}(\mathbf{x}^{(k-1)} - \mathbf{x}^{(k-2)})\| \leq \|\mathbf{T}\| \cdot \|\mathbf{x}^{(k-1)} - \mathbf{x}^{(k-2)}\|.$$

Applying this inequality recursively:

$$\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \|\mathbf{T}\|^{k-1} \cdot \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|.$$

For $m > k \geq 1$, we have:

$$\begin{aligned} \|\mathbf{x}^{(m)} - \mathbf{x}^{(k)}\| &= \|\mathbf{x}^{(m)} - \mathbf{x}^{(m-1)} + \mathbf{x}^{(m-1)} - \dots + \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \\ &\leq \|\mathbf{x}^{(m)} - \mathbf{x}^{(m-1)}\| + \|\mathbf{x}^{(m-1)} - \mathbf{x}^{(m-2)}\| + \dots + \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \\ &\leq \|\mathbf{T}\|^{m-1} \cdot \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| + \|\mathbf{T}\|^{m-2} \cdot \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| + \dots + \|\mathbf{T}\|^k \cdot \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|. \end{aligned}$$

Factoring out $\|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|$:

$$\|\mathbf{x}^{(m)} - \mathbf{x}^{(k)}\| \leq \|\mathbf{T}\|^k \cdot \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| (1 + \|\mathbf{T}\| + \dots + \|\mathbf{T}\|^{m-k-1}).$$

Taking the limit as $m \rightarrow \infty$:

$$\lim_{m \rightarrow \infty} \|\mathbf{x}^{(m)} - \mathbf{x}^{(k)}\| \leq \|\mathbf{T}\|^k \cdot \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| \sum_{i=0}^{\infty} \|\mathbf{T}\|^i.$$

Since $\|\mathbf{T}\| < 1$, we have:

$$\sum_{i=0}^{\infty} \|\mathbf{T}\|^i = \frac{1}{1 - \|\mathbf{T}\|}.$$

Thus:

$$\lim_{m \rightarrow \infty} \|\mathbf{x}^{(m)} - \mathbf{x}^{(k)}\| \leq \frac{\|\mathbf{T}\|^k}{1 - \|\mathbf{T}\|} \cdot \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|.$$

Since $\lim_{m \rightarrow \infty} \mathbf{x}^{(m)} = \mathbf{x}$, the second error bound follows. \square

Hence, if $\rho(\mathbf{T}_J)$ and $\rho(\mathbf{T}_{GS}) < 1$ then we get the convergence of these schemes. We also note from this corollary that the convergence of the method is based on $\|\mathbf{T}\|^k$. Now according to Theorem 2.27 we have $\rho(\mathbf{T}) \leq \|\mathbf{T}\|$. We have a more general result, i.e., for given $\varepsilon > 0$, $\rho(\mathbf{T}) \leq \|\mathbf{T}\| \leq \rho(\mathbf{T}) + \varepsilon$. Hence, we can say that

$$\|\mathbf{x}^{(k)} - \mathbf{x}\| \approx \rho(\mathbf{T})^k \|\mathbf{x} - \mathbf{x}^{(0)}\|.$$

Thus, we would like to select methods with minimal $\rho(\mathbf{T})$.

We have a sufficient method to show the convergence of the Jacobi and the Gauss-Seidel method.

Theorem 2.33. If \mathbf{A} is strictly diagonally dominant then for any choice of $\mathbf{x}^{(0)}$ both the Jacobi and the Gauss-Seidel methods give sequences $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$ that converge to the unique solution of $\mathbf{Ax} = \mathbf{b}$.

There is no best method in general between the Gauss-Seidel and the Jacobi method. However, in exceptional cases, we have some results.

Theorem 2.34. If $a_{ij} \leq 0$ for each $i \neq j$ and $a_{ii} > 0$ for each $i = 1, 2, \dots, n$, then one and only one of the following statement holds:

1. $0 \leq \rho(\mathbf{T}_{\text{GS}}) < \rho(\mathbf{T}_{\text{J}}) < 1$,
2. $1 < \rho(\mathbf{T}_{\text{J}}) < \rho(\mathbf{T}_{\text{GS}})$,
3. $\rho(\mathbf{T}_{\text{J}}) = \rho(\mathbf{T}_{\text{GS}}) = 0$,
4. $\rho(\mathbf{T}_{\text{J}}) = \rho(\mathbf{T}_{\text{GS}}) = 1$.

In the above theorem, we note that if 1. holds, then both methods converge together with \mathbf{T}_{GS} being better, and if 2. holds, then both diverge and \mathbf{T}_{GS} has “better” divergence.

2.3.3 Successive Over Relaxation

Definition 2.35. Suppose $\tilde{\mathbf{x}} \in \mathbb{R}^n$ is an approximation to the solution of the linear system defined by $\mathbf{Ax} = \mathbf{b}$. The *residual vector* for $\tilde{\mathbf{x}}$ with respect to the system is $\mathbf{r} = \mathbf{b} - \mathbf{A}\tilde{\mathbf{x}}$.

In procedure such as Jacobi or Gauss-Seidel method, a residual vector is associated with each calculation of an approximate component to the solution. The true objective of an iterative method is to generate a sequence of approximation that allows the residual vector to converge rapidly to zero. Suppose

$$\mathbf{r}_i^{(k)} = \left(r_{1i}^{(k)}, r_{2i}^{(k)}, \dots, r_{ni}^{(k)} \right)^\top,$$

denote the residual vector for the Gauss-Seidel method corresponding to the approximate solution vector $\mathbf{x}_i^{(k)}$ which is defined by

$$\mathbf{x}_i^{(k)} = \left(x_1^{(k)}, x_2^{(k)}, \dots, x_{i-1}^{(k)}, x_i^{(k-1)}, \dots, x_n^{(k-1)} \right)^\top.$$

Now, the m^{th} component of $\mathbf{r}_i^{(k)}$ is

$$r_{mi}^{(k)} = b_m - \sum_{j=1}^{i-1} a_{mj} x_j^{(k)} - \sum_{j=i}^n a_{mj} x_j^{(k-1)}, \quad (2.12)$$

or equivalently

$$r_{mi}^{(k)} = b_m - \sum_{j=1}^{i-1} a_{mj} x_j^{(k)} - \sum_{j=i+1}^n a_{mj} x_j^{(k-1)} - a_{mi} x_i^{(k-1)},$$

for each $m = 1, 2, \dots, n$.

In particular for $m = i$

$$r_{ii}^{(k)} = b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k)} - \sum_{j=i+1}^n a_{ij}x_j^{(k-1)} - a_{ii}x_i^{(k-1)},$$

which is equivalent to

$$a_{ii}x_i^{(k-1)} + r_{ii}^{(k)} = b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k)} - \sum_{j=i+1}^n a_{ij}x_j^{(k-1)}. \quad (2.13)$$

However, we know from the Gauss Seidel method that

$$x_i^{(k)} = \frac{1}{a_{ii}} \left[b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k)} - \sum_{j=i+1}^n a_{ij}x_j^{(k-1)} \right]. \quad (2.14)$$

So we can write Eq. (2.13) as

$$a_{ii}x_i^{(k-1)} + r_{ii}^{(k)} = a_{ii}x_i^{(k)}.$$

Consequently the Gauss Seidel method can be characterised as choosing $\mathbf{x}_i^{(k)}$ to satisfy

$$x_i^{(k)} = x_i^{(k-1)} + \frac{r_{ii}^{(k)}}{a_{ii}}. \quad (2.15)$$

Hence the update is determined by the residual at the current step.

Now let us look at the residual vector $\mathbf{r}_{i+1}^{(k)}$ associated with $\mathbf{x}_{(i+1)}^{(k)} = (x_1^{(k)}, x_2^{(k)}, \dots, x_{i+1}^{(k-1)}, \dots, x_n^{(k-1)})$

By Eq. (2.12) the i^{th} component of $\mathbf{r}_{i+1}^{(k)}$ is

$$\begin{aligned} r_{i,i+1}^{(k)} &= b_i - \sum_{j=1}^i a_{ij}x_j^{(k)} - \sum_{j=i+1}^n a_{ij}x_j^{(k-1)} \\ &= b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k)} - \sum_{j=i+1}^n a_{ij}x_j^{(k-1)} - a_{ii}x_i^{(k)}. \end{aligned}$$

By Eq. (2.14) we notice that the RHS is zero. In a way $x_{i+1}^{(k)}$ is chosen such that the i^{th} component of $\mathbf{r}_{i+1}^{(k)}$ is zero. But here only one component is zero which may not be the most efficient way to reduce the norm of the vector $\mathbf{r}_{i+1}^{(k)}$. Hence if we modify Eq. (2.15) to

$$x_i^{(k)} = x_i^{(k-1)} + \omega \frac{r_{ii}^{(k)}}{a_{ii}}, \quad (2.16)$$

then for certain values of ω we can reduce the norm of the residual. Eq. (2.16) refers to *relaxation methods*. If $\omega \in (0, 1)$ we get *under relaxation method* and if $\omega > 1$ we get *over relaxation methods*. Generally we refer to them as *successive over relaxation (SOR) methods*.

The idea of the SOR methods were devised simultaneously by Stan Frankel and David M. Young Jr. in the 1950s but the idea of the relaxation methods can be traced back way



Figure 2.9: Stan Frankel (1919 – May 1978, left) and David M. Young Jr. (20 October 1923 – 21 December 2008, right).

earlier. Interestingly, Frankel was also a part of the Manhattan project and was a PostDoc under Oppenheimer.

We first reformulate the SOR method. By Eq. (2.16)

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

which is $a_{ii} x_i^{(k)} + \omega \sum_{j=1}^{i-1} a_{ij} x_j^{(k)} = (1 - \omega) a_{ii} x_i^{(k-1)} - \omega \sum_{j=i+1}^n a_{ij} x_j^{(k-1)} + \omega b_i$. So in vectorise form this is

$$(\mathbf{D} - \omega \mathbf{L}) \mathbf{x}^{(k)} = [(1 - \omega) \mathbf{D} + \omega \mathbf{U}] \mathbf{x}^{(k-1)} + \omega \mathbf{b},$$

i.e.,

$$\mathbf{x}^{(k)} = (\mathbf{D} - \omega \mathbf{L})^{-1} [(1 - \omega) \mathbf{D} + \omega \mathbf{U}] \mathbf{x}^{(k-1)} + \omega (\mathbf{D} - \omega \mathbf{L})^{-1} \mathbf{b}.$$

Letting $\mathbf{T}_\omega = (\mathbf{D} - \omega \mathbf{L})^{-1} [(1 - \omega) \mathbf{D} + \omega \mathbf{U}]$ and $\mathbf{c}_\omega = \omega (\mathbf{D} - \omega \mathbf{L})^{-1} \mathbf{b}$ we get the SOR method as

$$\mathbf{x}^{(k)} = \mathbf{T}_\omega \mathbf{x}^{(k-1)} + \mathbf{c}_\omega. \quad (2.17)$$

Now, the next big question is what should be the appropriate value of ω . In general for $n \times n$ system we cannot say it but for particular cases we have the answer.

Theorem 2.36. *If $a_{ii} \neq 0$ for all $i = 1, 2, \dots, n$ then $\rho(\mathbf{T}_\omega) \geq |\omega - 1|$. This means SOR can only converge if $0 < \omega < 2$.*

Proof. Let $\{\lambda_i\}_{i=1}^n$ be the eigenvalues of \mathbf{T}_ω . Then

$$\begin{aligned} \rho(\mathbf{T}_\omega)^n &\geq \prod_{i=1}^n \lambda_i = \det(\mathbf{T}_\omega) \\ &= \det((\mathbf{D} - \omega \mathbf{L})^{-1} [(1 - \omega) \mathbf{D} + \omega \mathbf{U}]) \\ &= \det(\mathbf{D} - \omega \mathbf{L})^{-1} \det((1 - \omega) \mathbf{D} + \omega \mathbf{U}) \\ &= \det(\mathbf{D}^{-1}) \det((1 - \omega) \mathbf{D}) \\ &= \frac{1}{\prod_{i=1}^n a_{ii}} (1 - \omega)^n \prod_{i=1}^n a_{ii} = (1 - \omega)^n. \end{aligned}$$

Now, $\rho(\mathbf{T}_\omega) = \max_{1 \leq i \leq n} |\lambda_i| \geq |1 - \omega|$. Now the method will converge if $\rho(\mathbf{T}_\omega) < 1$. Hence $\omega \in (0, 2)$. \square

Next, we present a theorem regarding the convergence of the SOR method for symmetric positive definite matrices.

Theorem 2.37. *If \mathbf{A} is symmetric positive definite matrix and $\omega \in (0, 2)$ then the SOR method converges for any choice of initial approximation.*

Theorem 2.38. *If \mathbf{A} is symmetric positive definite and tridiagonal, then $\rho(\mathbf{T}_{\text{GS}}) = [\rho(\mathbf{T}_J)]^2 < 1$ and the optimal choice of ω for the SOR method is*

$$\omega = \frac{2}{1 + \sqrt{1 - [\rho(\mathbf{T}_J)]^2}},$$

with this choice of ω we have $\rho(\mathbf{T}_\omega) = \omega - 1$.

The SOR algorithm is presented in Algorithm 13.

2.3.4 Condition Number

When solving systems of linear equations, either using *iterative solvers* or *direct solvers*, numerical errors are inevitable. In iterative solvers, errors can accumulate due to finite precision and approximation, while in direct solvers, round-off errors from finite precision arithmetic can also affect the solution.

Consider the system of linear equations

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

where \mathbf{A} is an invertible matrix. Suppose we introduce a small perturbation $\delta\mathbf{b}$ to the right-hand side \mathbf{b} , resulting in a perturbed solution $\mathbf{x} + \delta\mathbf{x}$. The perturbed system is given by:

$$\mathbf{A}(\mathbf{x} + \delta\mathbf{x}) = \mathbf{b} + \delta\mathbf{b}.$$

Expanding this and subtracting the original system $\mathbf{A}\mathbf{x} = \mathbf{b}$ yields:

$$\mathbf{A}\delta\mathbf{x} = \delta\mathbf{b}.$$

Since \mathbf{A} is invertible, we can solve for $\delta\mathbf{x}$ as:

$$\delta\mathbf{x} = \mathbf{A}^{-1}\delta\mathbf{b}.$$

Taking the norm of both sides, we have:

$$\|\delta\mathbf{x}\| = \|\mathbf{A}^{-1}\delta\mathbf{b}\| \leq \|\mathbf{A}^{-1}\| \cdot \|\delta\mathbf{b}\|.$$

Dividing both sides by $\|\mathbf{x}\|$, we obtain the bound on the relative error in the solution:

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\|\mathbf{A}^{-1}\|}{\|\mathbf{x}\|} \cdot \|\delta\mathbf{b}\|.$$

Algorithm 13 SOR Iteration

Given: Matrix \mathbf{A} with non-zero pivots, right hand side \mathbf{b} , dimension n , ω , max_iterations , and tolerance .

Find: Solution \mathbf{x} .

Step 1: SOR Iterations

Initialize $\mathbf{x}^{\text{old}} = \mathbf{0}$

for $k = 1$ **to** max_iterations **do**

for $i = 1$ **to** n **do**

$\text{sum} = \mathbf{b}_i$

for $j = 1$ **to** n **do**

if $j < i$ **then**

$\text{sum} = \text{sum} - \mathbf{A}_{ij}\mathbf{x}_j$

else if $i < j$ **then**

$\text{sum} = \text{sum} - \mathbf{A}_{ij}\mathbf{x}_j^{\text{old}}$

end if

end for

$\mathbf{x}_i = (1 - \omega)\mathbf{x}_i^{\text{old}} + \omega \frac{\text{sum}}{\mathbf{A}_{ii}}$

end for

$\text{Error} = \|\mathbf{x} - \mathbf{x}^{\text{old}}\|_{\infty}$

if $\text{Error} < \text{tolerance}$ **then**

Output("Convergence reached")

break

end if

$\mathbf{x}^{\text{old}} = \mathbf{x}$

end for

if $k == \text{max_iterations}$ **then**

Output("Maximum Number of iterations reached")

end if

return \mathbf{x}

Since $\mathbf{Ax} = \mathbf{b}$, we know that:

$$\|\mathbf{b}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{x}\|.$$

Using this relationship, we can rewrite the relative error as:

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\| \cdot \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|}.$$

This inequality shows that the relative error in the solution is bounded by the relative error in the right-hand side, scaled by the factor $\|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\|$. We define the *condition number* of the matrix \mathbf{A} as:

$$\kappa(\mathbf{A}) = \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\|.$$

The condition number $\kappa(\mathbf{A})$ provides a measure of how sensitive the solution \mathbf{x} is to perturbations in \mathbf{b} . Specifically: - If $\kappa(\mathbf{A}) \approx 1$, the system is said to be *well-conditioned*, meaning small perturbations in \mathbf{b} lead to small errors in \mathbf{x} . - If $\kappa(\mathbf{A})$ is large, the system is *ill-conditioned*, and even small perturbations in \mathbf{b} may cause large errors in \mathbf{x} .

Preconditioners

To reduce the condition number and improve numerical stability, we can use a technique called *preconditioning*. Preconditioning involves transforming the system $\mathbf{Ax} = \mathbf{b}$ by multiplying both sides with a matrix \mathbf{P} to obtain an equivalent system with a lower condition number. There are two common types of preconditioning:

1. **Left Preconditioning:** Multiply both sides of the system by \mathbf{P}^{-1} :

$$\mathbf{P}^{-1}\mathbf{Ax} = \mathbf{P}^{-1}\mathbf{b}.$$

2. **Right Preconditioning:** Solve the system:

$$\mathbf{AP}^{-1}\mathbf{y} = \mathbf{b}, \quad \text{where } \mathbf{x} = \mathbf{P}^{-1}\mathbf{y}.$$

A good preconditioner \mathbf{P} should satisfy two key properties:

1. The convergence of the iterative method applied to the preconditioned system $\mathbf{P}^{-1}\mathbf{A}$ or \mathbf{AP}^{-1} should be faster than for the original system.
2. Solving the system involving \mathbf{P} should be computationally inexpensive.

In practice, a balance must be struck between these two requirements.

Some commonly used preconditioners are:

1. **Jacobi (Diagonal) Preconditioner:** $\mathbf{P} = \mathbf{D}$, where \mathbf{D} is the diagonal part of \mathbf{A} .
2. **Forward Gauss-Seidel Preconditioner:** $\mathbf{P} = \mathbf{D} + \mathbf{L}$, where \mathbf{L} is the strict lower triangular part of \mathbf{A} .
3. **Backward Gauss-Seidel Preconditioner:** $\mathbf{P} = \mathbf{D} + \mathbf{U}$, where \mathbf{U} is the strict upper triangular part of \mathbf{A} .
4. **Symmetric Gauss-Seidel Preconditioner:** $\mathbf{P} = (\mathbf{D} + \mathbf{L})\mathbf{D}^{-1}(\mathbf{D} + \mathbf{U})$.

For convenience, we often denote the preconditioner by \mathbf{P}^{-1} rather than \mathbf{P} . Preconditioning is a powerful tool in improving the stability and performance of numerical solvers, and iterative methods are often used in conjunction with preconditioners rather than as standalone solvers.

2.4 Least Square Methods

Least Square Problems has been quite helpful in different science areas, from physics to data science. In simple language, i.e., the mathematical language, we are trying to solve an over-determined system, i.e., $\mathbf{Ax} = \mathbf{b}$, by minimizing the ℓ_2 norm of the residual.

Consider a linear system of equation $\mathbf{Ax} = \mathbf{b}$ with n unknowns and m equations with $m > n$, i.e.,

$$\mathbf{A} = \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ a_{21} & \cdots & a_{2n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{bmatrix}, \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}.$$

Hence we need to compute $\mathbf{x} \in \mathbb{R}^n$ such that $\mathbf{Ax} = \mathbf{b}$ with $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. A direct solution to such a problem does not exist as the invertibility of the matrix is in question. Hence, instead, we try to reduce the residual \mathbf{r} given by

$$\mathbf{r} = \mathbf{b} - \mathbf{Ax} \in \mathbb{R}^m.$$

What do we mean by reduction? If we choose the ℓ_2 norm then the problem is: Given $\mathbf{A} \in \mathbb{R}^{m \times n}$, $m \geq n$ and $\mathbf{b} \in \mathbb{R}^m$, find $\mathbf{x} \in \mathbb{R}^n$ such that $\|\mathbf{b} - \mathbf{Ax}\|_2$ is minimized.

The choice of ℓ_2 norm can be justified geometrically. We seek a vector $\mathbf{x} \in \mathbb{R}^n$ such that $\mathbf{Ax} \in \mathbb{R}^m$ is closest to point \mathbf{b} in range of \mathbf{A} (see Fig. 2.4).

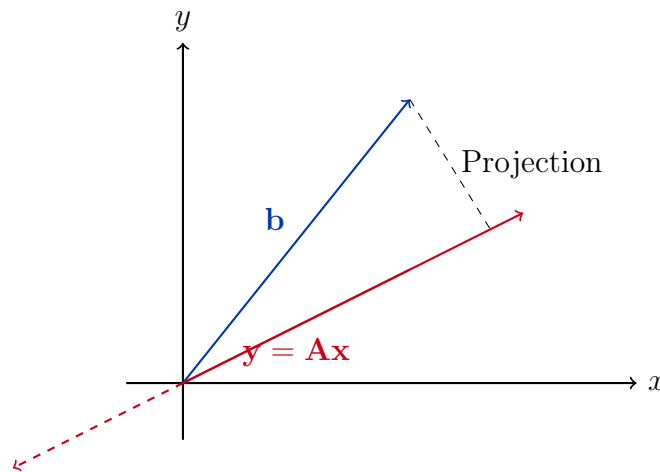


Figure 2.10: Orthogonal projection.

We want to find \mathbf{Ax} in the range of \mathbf{A} such that $\mathbf{r} = \mathbf{Ax} - \mathbf{b}$ is minimum. It is clear from the geometry that $\mathbf{Ax} = \mathbf{Pb}$ is the solution where $\mathbf{P} \in \mathbb{R}^{m \times m}$ is the orthogonal projection operator that maps \mathbb{R}^m to range of \mathbf{A} . In other words, the residual must be orthogonal to $\text{range}(\mathbf{A})$.

Theorem 2.39. Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ ($m \geq n$) and $\mathbf{b} \in \mathbb{R}^m$ be given. A vector $\mathbf{x} \in \mathbb{R}^n$ minimizes the residual norm $\|\mathbf{r}\|_2 = \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2$ thereby solving the least square problem if and only if $\mathbf{r} \perp \text{range}(\mathbf{A})$, i.e.,

$$\mathbf{A}^\top \mathbf{x} = \mathbf{0}, \quad (2.18)$$

or

$$\mathbf{A}^\top \mathbf{A}\mathbf{x} = \mathbf{A}^\top \mathbf{b}, \quad (2.19)$$

or

$$\mathbf{P}\mathbf{b} = \mathbf{A}\mathbf{x}, \quad (2.20)$$

where $\mathbf{P} \in \mathbb{R}^{m \times m}$ is the orthogonal projection onto $\text{range}(\mathbf{A})$. The $n \times n$ system Eq. (2.19) known as normal equation is non-singular if and only if \mathbf{A} has full rank.

We have not talked about the orthogonal projection \mathbf{P} , but we will use its certain formulations, the major one being $\mathbf{P} = \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top$. Interested readers can read any standard Linear Algebra book to read more about it; see [5, Section 6.6].

Now how do we actually solve Eq. (2.18), Eq. (2.19) or Eq. (2.20)? If \mathbf{A} has full rank, then the solution to the least square problem is unique and given by

$$\mathbf{x} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}.$$

The matrix $(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top$ is called the *pseudoinverse* of \mathbf{A} and is denoted by \mathbf{A}^+ . This is a matrix of size $n \times m$. The problem is to compute one or both vectors

$$\mathbf{x} = \mathbf{A}^+ \mathbf{b} \quad \mathbf{y} = \mathbf{P}\mathbf{b},$$

where \mathbf{A}^+ is the pseudoinverse of \mathbf{A} .

First, look at Eq. (2.19) and try to solve it. Now, we have that $\mathbf{A}^\top \mathbf{A}$ is a symmetric and positive definite matrix. Hence, we can apply Cholesky Decomposition (see 10) to write

$$\mathbf{A}^\top \mathbf{A} = \mathbf{L}\mathbf{L}^\top$$

and then solve $(\mathbf{L}\mathbf{L}^\top)\mathbf{x} = \mathbf{A}^\top \mathbf{b}$ to get \mathbf{x} . Here, it is important to note that we need to solve two systems of equations.

2.4.1 QR Decomposition

We have seen that matrix factorization has certain advantages. There is one factorization that is useful for least square methods.

We recall from Linear Algebra that give linearly independent vectors $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\}$ in \mathbb{R}^n we can compute an orthogonal linearly independent set of vectors $\{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n\}$ using Gram Schmidt orthogonalization [5, Theorem 6.4]. Further, we can also compute an orthonormal set of linearly independent vectors $\{\hat{\mathbf{q}}_1, \hat{\mathbf{q}}_2, \dots, \hat{\mathbf{q}}_n\}$. In matrix notation it means if we denote $\mathbf{A}|_{m \times n}$ by columns of $\{\mathbf{A}_i\}_{i=1}^n$ and $\mathbf{Q}|_{m \times n}$ by $\{\hat{\mathbf{q}}_i\}_{i=1}^n$ then

$$\mathbf{A} = \mathbf{Q}\mathbf{R},$$

where \mathbf{R} is a $n \times n$ matrix. To compute \mathbf{R} , we re-write the Gram-Schmidt orthogonalization as

$$\begin{aligned}\mathbf{A}_1 &= \hat{\mathbf{q}}_1 \|\mathbf{q}_1\| \\ \mathbf{A}_2 &= \hat{\mathbf{q}}_2 \|\mathbf{q}_2\| + \langle \mathbf{A}_2, \hat{\mathbf{q}}_1 \rangle \|\hat{\mathbf{q}}_1\| \\ \mathbf{A}_3 &= \hat{\mathbf{q}}_3 \|\mathbf{q}_3\| + \langle \mathbf{A}_3, \hat{\mathbf{q}}_1 \rangle \|\hat{\mathbf{q}}_1\| + \langle \mathbf{A}_3, \hat{\mathbf{q}}_2 \rangle \|\hat{\mathbf{q}}_2\|,\end{aligned}$$

where $\langle \cdot, \cdot \rangle$ denotes the ℓ_2 inner product. Looking at the pattern, we notice that \mathbf{R} is an upper triangular matrix such that

$$r_{ij} = \begin{cases} \langle \mathbf{A}_j, \hat{\mathbf{q}}_i \rangle & \text{if } i < j, \\ \|\hat{\mathbf{q}}_i\| & \text{if } i = j, \\ 0 & \text{else.} \end{cases}$$

The existence of the QR factorization comes from Gram-Schmidt orthogonalization, and the uniqueness follows the same logic as in Theorem 2.5.

Now, to solve Eq. (2.20) using $\mathbf{A} = \mathbf{QR}$ and the definition of \mathbf{P} we note that

$$\begin{aligned}\mathbf{P} &= \mathbf{A} (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \\ &= (\mathbf{QR}) \left[(\mathbf{QR})^\top \mathbf{QR} \right]^{-1} (\mathbf{QR})^\top \\ &= \mathbf{QR} \left[\mathbf{R}^\top \mathbf{Q}^\top \mathbf{QR} \right]^{-1} \mathbf{R}^\top \mathbf{Q}^\top \\ &= \mathbf{QR} \left[\mathbf{R}^\top \mathbf{R} \right]^{-1} \mathbf{R}^\top \mathbf{Q}^\top \\ &= \mathbf{Q} (\mathbf{RR}^{-1}) (\mathbf{R}^\top)^{-1} \mathbf{R}^\top \mathbf{Q}^\top \\ &= \mathbf{QQ}^\top.\end{aligned}$$

The above result holds as \mathbf{Q}^\top is the left-inverse of \mathbf{Q} . Using this, we re-write

$$\begin{aligned}\mathbf{Pb} &= \mathbf{Ax} \\ \mathbf{QQ}^\top \mathbf{b} &= \mathbf{QRx} \\ \mathbf{Q}^\top \mathbf{b} &= \mathbf{Rx}.\end{aligned}$$

Hence, we get $\mathbf{x} = \mathbf{R}^{-1} \mathbf{Q}^\top \mathbf{b}$. Here, we note that we only need to solve one system of equations.

The algorithm for QR decomposition using Gram-Schmidt orthogonalization is present in Algorithm 14.

Algorithm 14 QR Decomposition

Given: Matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m \geq n$.

Find: Orthonormal matrix \mathbf{Q} and upper triangular matrix \mathbf{R} such that $\mathbf{A} = \mathbf{QR}$

Step 1: Initialize Matrices

Initialize \mathbf{Q} as a zero matrix of size $m \times n$.

Initialize \mathbf{R} as a zero matrix of size $n \times n$.

for $j = 1$ **to** n **do**

 Set $\mathbf{q}_j = \mathbf{A}_j$

for $i = 1$ **to** $j - 1$ **do**

$r_{ij} = \langle \mathbf{q}_i, \mathbf{A}_j \rangle$

$\mathbf{q}_j = \mathbf{q}_j - r_{i,j} \mathbf{q}_i$

end for

$r_{jj} = \|\mathbf{q}_j\|_2$.

$\hat{\mathbf{q}}_j = \mathbf{q}_j / r_{jj}$.

end for

return $\mathbf{Q} = [\hat{\mathbf{q}}_1, \hat{\mathbf{q}}_2, \dots, \hat{\mathbf{q}}_n], \mathbf{R}$

Chapter 3

Computing

The word **computing** has different meanings based on context and definition. According to Wikipedia:

Computing is any goal-oriented activity requiring, benefitting from, or creating computing machinery.

This definition creates a recursive loop, as it uses “computing” to define itself. To break this loop, let us explore what a **computer** is. Wikipedia defines it as:

A computer is a machine that can be programmed to automatically carry out sequences of arithmetic or logical operations (computation).

Here, two terms stand out: *arithmetic* and *logical*. These are fundamental concepts that mathematicians are familiar with. Thus, we have some basic understanding of computing.

In this course, we will not delve deeply into the workings of a computer. It is assumed that students are familiar with components like the keyboard, mouse (or trackpad), CPU, and monitor. For a refresher, see Figure 3.1.

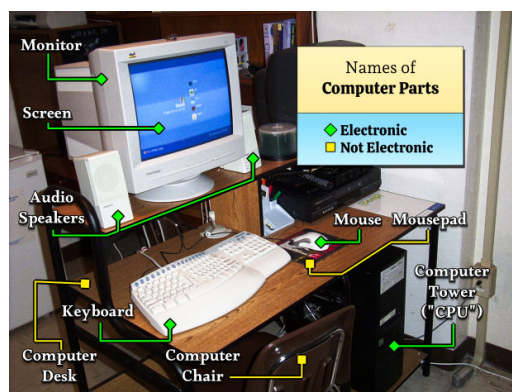


Figure 3.1: Parts of a computer from the early 2000s.

The primary aim of this course is to teach the fundamentals of programming using the Python language. This is not an *Introduction to Python* course. Instead, the focus is on how

to think about coding problems, understand common paradigms across languages, adopt good coding practices, maintain code, and debug effectively.

3.1 Good Practices in Coding

Coding is an art form, and like any art, its true audience is those who interact with it (in this case, the users of the code). A good codebase should be **well-documented**. Comments should explain the purpose of each function or variable. The beginning of the code should clearly state its objective.

3.1.1 Variable Initialization and Naming

Variables and functions should have **descriptive names**. For example, if a variable represents the number of oranges, naming it `n_oranges` is much clearer than simply using `n`. Additionally, variables should be initialized to prevent the use of garbage values.

Indentation of the conditional and iterative statements are important as it helps to differentiate different loops (or if-else statements).

In Python, indentation is mandatory, making this practice less error-prone. However, for languages like C++, proper formatting and indentation are crucial. Below we give an example of bad coding vs good coding in C++.

Bad Example:

```
int n_oranges;
for (int i = 0; i < 10; i++)
{
std::cout << i;
for(int j=0;j<2;j++
{
std::cout<<i+j;
}
std::cout<<n_oranges;
}
```

Good Example:

```
int n_oranges = 0;
for (int i = 0; i < 10; i++)
{
    std::cout << i;
    for (int j = 0; j < 2; j++)
    {
        std::cout << i + j;
    }
    std::cout<<n_oranges;
}
```


3.1.2 Reusability and Modularity

Code should be **reusable** and **modular**. For instance, consider a program that computes the Taylor series of a function. Instead of hardcoding the factorial computation in the main function, create a separate function for it. This is called *modularity*. This approach makes the code reusable. If another project requires the computation of 2C_k , the factorial function can be reused without rewriting it.

3.2 Testing and Continuous Integration

Testing is a critical aspect of programming to ensure correctness and reliability. Continuous integration ensures that code changes do not break existing functionality.

After writing code, how do we know if it is correct? One effective approach is to verify the solution produced by the code against a pre-existing known solution. For example, if we write code to find the roots of a function, we can test its accuracy by using values with known solutions, such as $x^2 = 2$.

It is always advisable to run the code on multiple test cases to validate its correctness. Once the code is verified, we can create specific *test routines* to ensure its reliability in various scenarios.

3.3 Introduction to Computing Using Python

This course covers various aspects of computing, but we begin with the basics to build a strong foundation.

3.3.1 Variables

In Python, there are three commonly used variable types:

- **int**: Represents integers.
- **str**: Represents strings (text).
- **float**: Represents floating-point numbers (decimals).

Python does not require explicit variable declaration; you can assign values directly. For example:

```
n_oranges = 10 # An integer
price_oranges = 10.4 # A floating-point number
```

To define strings, use double quotes ("):

```
fruit = "oranges" # A string
```

To check the type of a variable, use the `type()` function:

```
print(type(n_oranges)) # Output: <class 'int'>
```

Another important variable type is the **list**, which can contain multiple values of different types:

```
list_fruits = [n_oranges, price_oranges, fruit] # A list with mixed
            types
```

While there are more variable types in Python, these four are essential for now.

Note: The `print()` function is used to display information. We will explore more advanced printing techniques later.

3.3.2 Arithmetic Operations

Arithmetic operations are fundamental in any programming language. Python provides the following operations:

Operation	Description	Example
+	Addition	$2 + 2 = 4$
-	Subtraction	$6 - 2 = 4$
*	Multiplication	$2 * 2 = 4$
/	Division	$2/2 = 1$
**	Exponentiation	$2 ** 2 = 4$
==	Equality comparison	$2 == 2$
%	Modulus (remainder)	$3\%2 = 1$

Table 3.1: Arithmetic Operations in Python

Additionally, the `!=` operator means "not equal to," as in $3 \neq 2$.

These operations allow us to build more advanced functions and logic.

3.3.3 Compound Assignment

Python supports shorthand operations for self-assignment:

```
A = 10
A = A + 10 # Equivalent to:
A += 10
```

This shorthand applies to all arithmetic operations (`+=`, `-=`, `*=`, `/=`, etc.).

Note: When performing operations on variables of different types, such as `int` and `float`, Python automatically converts the result to `float`. For example:

```
result = 10 + 10.5 # result is a float (20.5)
```

However, operations combining `str` with `int` or `float` will result in errors:

```
"10" + 10 # This will raise a TypeError
```

Experiment with different cases to understand how Python handles these scenarios. Also, remember that Python follows the BODMAS convention.

3.3.4 Logical Operations

Another important class of operations is **logical operations**. These operations are used when you need to run specific parts of the code based on multiple conditions. For example:

- To check if a number is greater than 5 **and** divisible by 3.
- To check if a number is greater than 5 **or** divisible by 3.

In both cases, you have a logical expression to evaluate. Python provides three logical operators to handle such cases:

- (i) **and**: Evaluates to **True** if both conditions A and B are true, otherwise **False**.
- (ii) **or**: Evaluates to **True** if at least one of the conditions A or B is true, otherwise **False**.
- (iii) **not**: Returns the negation of condition A.

The truth table for these logical operators is shown below:

A	B	A and B	A or B	not A
T	T	T	T	F
T	F	F	T	F
F	T	F	T	T
F	F	F	F	T

Table 3.2: Truth table for logical operators.

Logical operations can also involve more than two conditions. For example, suppose you have three conditions: **A**, **B**, and **C**. In such cases, you can group conditions using parentheses to control the order of evaluation. For instance:

- Check (**A and B**) first, and then combine the result with **C**.
- Evaluate **A or (B and C)** to prioritize **B and C**.

This flexibility allows for constructing complex logical expressions tailored to your requirements.

3.4 Conditional Statements

Conditional statements allow executing specific code blocks depending on conditions. In Python, the syntax is as follows:

```
if condition_1:
    execute_1
else:
    execute_2
```

Example: Checking if a number is even or odd:

```
eval_point = 5
if eval_point % 2 == 0:
    print(f"The number {eval_point} is even.")
else:
    print(f"The number {eval_point} is odd.")
```

For multiple conditions, we use if-elif-else:

```
eval_point = 5
if eval_point % 2 == 0:
    print(f"The number {eval_point} is divisible by 2.")
elif eval_point % 3 == 0:
    print(f"The number {eval_point} is divisible by 3.")
else:
    print(f"The number {eval_point} is not divisible by 2 or 3.")
```

Note: An else statement is usually not necessary for an if statement. Suppose we want to check if a number is even we can use

```
eval_point = 5
if eval_point % 2 == 0:
    print(f"The number {eval_point} is divisible by 2.")
```

Here we want to check if the number is even without checking if it odd or not.

Note: Here we have introduced a new way to print. The `print(f"...")` command is printing a formatted string. It prints the characters as well as the variable values defined in the curly braces `{·}`.

3.5 Recursive Statements

Recursive Statements allow repetitive execution of code blocks.

3.5.1 For Loop

The syntax for a for loop in Python is:

```
for i in range(a, b):
    execute_1
```

Here:

- `i`: The loop iterator.
- `range(a, b)`: Specifies the range of values, starting at `a` and stopping before `b`, i.e., it goes over `a, a + 1, ..., b - 1`.

Example: Summing numbers from 1 to 9:

```
total_sum = 0
for i in range(1, 10):
    print(i)
    total_sum += i
print(f"The summation of 9 points: {total_sum}")
```

3.5.2 Custom Step Size

The default step size of a for loop is one. If we want to use a custom step-size then we can use the following modification:

```
for i in range(a, b, step):
    execute_1
```

Example: Summing odd numbers from 1 to 9:

```
total_sum = 0
for i in range(1, 10, 2):
    print(i)
    total_sum += i
print(f"The summation of 10 points with step 2: {total_sum}")
```

3.5.3 Break and Continue

While using loops there can be cases when we want to exit the loop due to some condition. Also we can have cases when we want to skip some iteration. In this case we use `break` and `continue`, respectively.

- `break` : Exits the loop entirely.
- `continue` : Skips the current iteration and moves to the next.

Example: Adding even numbers up to 10 but stopping at 7:

```
total_sum = 0
for i in range(10):
    if i == 7:
        break
    if i % 2 == 1:
        continue
    print(i)
    total_sum += i
print(f"The summation of even numbers: {total_sum}")
```

3.5.4 Nested Loops

A for loop can be nested within another for loop. For example, to generate multiplication tables:

```
for i in range(1, 5):
    print(f"The table of {i}")
    for j in range(1, 11):
        print(f"{i} x {j} = {i * j}")
```

3.6 Functions

Until now, we have focused on sequential coding. However, to enhance reusability and maintainability, modular coding is essential. Functions enable modular programming by allowing code reuse. The syntax for creating a function is:

```
def function_name(input_1, input_2):
    # Function body
    result = some_operation(input_1, input_2)
    return result
```

Example: Function to check if a number is even:

```
def is_even(number):
    if number % 2 == 0:
        return True
    else:
        return False

value = 20
result = is_even(value)
print(f"The number {value} is even: {result}")
```

Note:

- A function can accept multiple inputs, a single input, or no input at all.
- A function may include multiple `return` statements or omit a `return` entirely, in which case it returns `None` by default.

Example: Function with multiple return statements:

```
def analyze_number(number):
    if number > 0:
        return "positive", number
    elif number < 0:
        return "negative", number
    else:
        return "zero", number

result_type, result_value = analyze_number(-5)
print(f"The number {result_value} is {result_type}.")
```

3.7 NumPy Library

In this section we provide an introduction to the NumPy library, a fundamental Python library for mathematical computations.

Importing the Library

To use NumPy in Python, the library must be imported. The standard convention is to import it with the alias `np`:

Import the NumPy library as follows: `import numpy as np.`

3.7.1 Arrays and Matrices

NumPy arrays are versatile tools used as vectors (one-dimensional) or matrices (two-dimensional). For instance:

- A one-dimensional array can be thought of as a row vector, e.g., $[1, 2, 3, 4]$. This kind of array can be created using `temp_array = np.array([1, 2, 3, 4])`.
- A two-dimensional matrix is defined by nesting arrays, e.g., $\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{bmatrix}$. This kind of array can be created using `temp_array = np.array([[1, 2, 3, 4], [5, 6, 7, 8]])`.

Key commands include:

- `np.size(array)`: Returns the total number of elements in an array.
- `np.shape(array)`: Provides the dimensions of an array.

Special Arrays: Zeros and Ones

It is common to initialize arrays with default values such as zeros or ones. This helps avoid uninitialized or garbage values in computations. For example:

- A zero matrix of size 10×10 can be created using `np.zeros((10, 10))`.
- Similarly, a ones matrix of the same size is created with `np.ones((10, 10))`.

Indexing in Arrays and Matrices

In NumPy:

- Indexing starts at 0.
- Negative indexing allows access to elements from the end, e.g., -1 refers to the last element.

For matrices, indexing uses row and column coordinates. For example, the element at row 0, column 0 in a matrix is accessed as `A[0][0]`. Suppose we have the following $n \times n$ matrix

$$\mathbf{A} = \begin{bmatrix} a_{00} & a_{01} & \dots & a_{0 \ n-1} \\ a_{10} & a_{11} & \dots & a_{1 \ n-1} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n-1 \ 0} & a_{n-1 \ 1} & \dots & a_{n-1 \ n-1} \end{bmatrix},$$

then the $a_{n-1 \ n-1}$ entry can be accessed using `A[-1][-1]` as well as `A[n - 1][n - 1]`. Similarly, $a_{0 \ n-1}$ can be accessed using `A[0][-1]` and `A[0][n - 1]`.

3.7.2 Linspace

To generate arrays with evenly spaced points, the `np.linspace` function is used:

For an array between a and b with n elements, use the syntax:
`np.linspace(a, b, n)`.

This is particularly useful for numerical methods. It is important to note that the end points a, b are included and the spacing between the points is $(b - a)/(n - 1)$.

3.7.3 Mathematical Functions

NumPy provides a wide range of mathematical functions, including:

- Trigonometric functions: `np.sin`, `np.cos`, etc.
- Hyperbolic functions: `np.sinh`, `np.cosh`, etc.
- Absolute value: `np.abs`.

For example, the sine and absolute value of $-\pi$ can be computed using `np.sin(-np.pi)` and `np.abs(-np.pi)`. For a comprehensive list of available mathematical routines, refer to the official documentation: <https://numpy.org/doc/stable/reference/routines.math.html>.

List of Algorithms

1	Vandermonde Interpolation	11
2	Lagrange Interpolation	15
3	Newton Interpolation	20
4	Hermite Interpolation	26
5	Cubic Natural Spline Interpolation	32
6	Gauss Elimination	40
7	Gauss Jordan	43
8	LU Decomposition with Partial Pivoting	49
9	LDL ^T Decomposition	53
10	Cholesky Decomposition	55
11	Jacobi Iteration	63
12	Gauss-Seidel Iteration	65
13	SOR Iteration	73
14	QR Decomposition	78

Index

- break, 85
- continue, 85
- if and else, 84

- Arithmetic Operations, 82
- Augmented Matrix, 36

- B-Splines, 31
- Band Matrix, 35

- Characteristic Polynomial, 59
- Cholesky Decomposition, 54
- Clamped Boundary, 28
- Compact Support, 33
- Condition Number, 73
- Conditional Statements, 83
- Convergent Matrix, 60
- Cubic Spline, 27

- Diagonally Dominant Matrix, 48

- Eigenvalues, 60
- Eigenvector, 60

- For Loop, 84

- Gauss Jordan Algorithm, 42
- Gauss Seidel Method, 63
- Gaussian Elimination, 36
- Generalized Rolle's Theorem, 14

- Hermite Interpolation, 22

- Jacobi Method, 61

- Knots, 28

- Lagrange Interpolation, 11
- LDL^T Decomposition, 48
- Leading Principal Sub-Matrix, 52
- Least Square Problems, 74
- Linspace, 88
- Logical Operations, 83

- LU Decomposition, 42

- Mean Value Theorem, 18
- Modular Coding, 81

- Natural Boundary, 28
- Natural Norm, 59
- Natural Spline, 28
- Newton Divided Differences, 17
- Not-A-Knot, 31
- NumPy, 87

- Orthogonal Projection, 76
- Over Relaxation Method, 70

- Permutation Matrix, 46
- Pivot, 38
- PLU Decomposition, 46
- Polynomial Interpolation, 9
- Positive Definite Matrix, 51
- Preconditioning, 74
- Pseudoinverse, 76

- QR Decomposition, 76

- Recursive Statements, 84
- Relaxation Methods, 70
- Residual Vector, 69
- Reusable Coding, 81
- Runge Function, 15
- Runge Phenomena, 15

- Spectral Radius, 60
- Splines, 27
- Strictly Diagonal Dominant Matrix, 31
- Successive Over Relaxation Methods, 70
- Symmetric Positive Definite, 51

- Taylor's Theorem, 8
- Testing, 81

- Under Relaxation Method, 70

Vandermonde Matrix, 10

Variables, 81

Vector Norm, 56

Weierstrass Approximation Theorem, 8

Bibliography

- [1] Robert G. Bartle and Donald R. Sherbert. *Introduction to real analysis*. Second. John Wiley & Sons, Inc., New York, 1992, pp. xii+404. ISBN: 0-471-51000-9.
- [2] Jean-Paul Berrut and Lloyd N. Trefethen. “Barycentric Lagrange interpolation”. In: *SIAM Rev.* 46.3 (2004), pp. 501–517. ISSN: 0036-1445,1095-7200. DOI: 10.1137/S0036144502417715. URL: <https://doi.org/10.1137/S0036144502417715>.
- [3] J. Douglas Faires and Richard Burden. *Numerical methods*. Second. With 1 IBM-PC floppy disk (3.5 inch; HD). Brooks/Cole Publishing Co., Pacific Grove, CA, 1998, pp. xii+594. ISBN: 0-534-35187-5.
- [4] L. R. Ford Jr. and D. R. Fulkerson. “Maximal flow through a network”. In: *Canadian J. Math.* 8 (1956), pp. 399–404. ISSN: 0008-414X,1496-4279. DOI: 10.4153/CJM-1956-045-5. URL: <https://doi.org/10.4153/CJM-1956-045-5>.
- [5] S.H. Friedberg, A.J. Insel, and L.E. Spence. *Linear Algebra*. Pearson Education, 2014. ISBN: 9780321998897. URL: <https://books.google.co.in/books?id=KyB0DAAAQBAJ>.
- [6] Gilbert Strang. *Linear algebra and its applications*. Academic Press [Harcourt Brace Jovanovich, Publishers], New York-London, 1976, pp. xi+374.