

# Numerical Methods

Xuemei Chen

New Mexico State University

(Last updated on April 30, 2020)

# Contents

<b>Notations</b>	<b>1</b>
<b>1 Floating Point Arithmetic</b>	<b>2</b>
1.1 Floating Point Representation . . . . .	2
1.2 Floating Point Operations . . . . .	4
1.3 IEEE 754 Standard . . . . .	5
1.4 Errors in Scientific Computing . . . . .	6
Exercises . . . . .	7
<b>2 Solving nonlinear equations</b>	<b>10</b>
2.1 The Bisection Method . . . . .	10
2.2 The Newton's Method . . . . .	11
2.3 The Secant Method . . . . .	14
Exercises . . . . .	14
<b>3 Solving system of linear equations: Direct methods</b>	<b>16</b>
3.1 Gaussian Elimination and $LU$ factorization . . . . .	16
3.1.1 Forward elimination = $LU$ factorization . . . . .	18
3.1.2 A 'second way' to solve $Ax = b$ . . . . .	19
3.2 Pivoting (row exchange) . . . . .	20
3.2.1 Partial pivoting . . . . .	21
3.2.2 PLU factorization . . . . .	22
3.2.3 Methods and software . . . . .	24
Exercises . . . . .	24
<b>4 Polynomial Interpolation</b>	<b>29</b>
4.1 Linear Interpolation . . . . .	29
4.2 Interpolation by One Polynomial . . . . .	30
4.2.1 Vandermonde matrix . . . . .	30
4.2.2 Lagrange interpolation . . . . .	31
4.2.3 Newton's interpolation . . . . .	32
4.2.4 Flop count . . . . .	33
4.2.5 Limitation of polynomial interpolant . . . . .	33
4.3 Chebyshev nodes . . . . .	34
4.4 Spline . . . . .	36
4.4.1 Cubic Spline . . . . .	38
Exercises . . . . .	40

<b>5</b>	<b>Numerical Integration</b>	<b>44</b>
5.1	Basic Quadrature Rules . . . . .	44
5.2	Composite Quadrature Rules . . . . .	46
5.3	Gaussian Quadrature . . . . .	48
	Exercises . . . . .	50
<b>6</b>	<b>Least Squares and Singular Value Decomposition</b>	<b>52</b>
6.1	Positive definite matrices . . . . .	52
6.2	Cholesky Factorization . . . . .	53
6.3	Least Squares . . . . .	53
6.4	Orthogonal matrices . . . . .	55
6.5	Singular Value Decomposition . . . . .	56
	Exercises . . . . .	57
<b>7</b>	<b>Preliminaries</b>	<b>60</b>
7.1	Calculus . . . . .	60
7.2	Linear Algebra . . . . .	61
7.2.1	Vector/Matrix algebra . . . . .	61
7.2.2	Matrix Inversion . . . . .	62
7.2.3	Determinant . . . . .	63
7.2.4	Linear Independence and span . . . . .	64
7.2.5	Eigenvalues . . . . .	64
7.2.6	Norm . . . . .	64
7.3	Complexity . . . . .	65

# Notations

Before this course, you need to be familiar with

- Calculus I
- Calculus II: Integrations. Sequences and Series.
- It is not necessary, but helpful to know some multivariate Calculus
- Some basic Linear Algebra.

## Notations

- $\{x | \text{descriptions of } x\}$  is the set of all  $x$  that meets the description after the “.” sign. For example,  $\{x | x^2 - 1 < 0\}$  is the set of all number  $x$  such that  $x^2 - 1 < 0$ . We can solve this inequality further and see that  $\{x | x^2 - 1 < 0\} = (-1, 1)$

Remark: Some books use colon instead of verticle bar as  $\{x : \text{descriptions of } x\}$

- $\in$ : belongs to/is contained in. For example,  $a \in \{x | x^2 - 1 < 0\}$  means that  $a$  is a number satisfying  $a^2 - 1 < 0$ .
- A vector in  $\mathbb{R}^n$  is a column vector in default.
- In all exercises at the end of each chapter, \* means extra credit problems.

# Chapter 1

## Floating Point Arithmetic

The arithmetic performed by a calculator or computer is different from the arithmetic that we use in our algebra or calculus class. In the ideal world, we permit numbers with an infinite number of digits. For example,  $\sqrt{2}$  is the unique positive number that, when multiplied by itself, produces the integer 2. It is an irrational number, which means that there are infinitely many digits associated with  $\sqrt{2}$  when represented in decimals. In the computational world, we are only able to assign a fixed and finite number of digits to each number.

Try add up 0.1 and 0.2 in Python, you will get

```
>>> 0.1 + 0.2
0.30000000000000004
```

This can cause serious issues and create bugs. See Section 5.1 of [2]. The above error is caused by rounding in floating point numbers, with which numerical analysis is traditionally concerned.

### 1.1 Floating Point Representation

Floating point numbers are represented in terms of a base  $\beta$ , a precision  $p$ , and an exponent  $e$ . For example, in base 10, 0.34 can be represented as  $3.4 \times 10^{-1}$  or  $34 \times 10^{-2}$ . In order to guarantee uniqueness, we take a floating point number to have the specific form.

$$\pm (d_0 + d_1\beta^{-1} + \cdots + d_{p-1}\beta^{-(p-1)}) \times \beta^e, \quad (1.1)$$

where each  $d_i$  is an integer in  $[0, \beta)$  and  $d_0 \neq 0$ . Such a representation is said to be a *normalized floating point number*. In this representation,  $d_0 + d_1\beta^{-1} + \cdots + d_{p-1}\beta^{-(p-1)}$  is called *mantissa* (or *significand*). Here are a few examples where  $\beta = 10$ ,  $p = 5$  and  $e$  is in the range  $-10 \leq e \leq 10$ :

$$1.2345 \times 10^8, \quad 3.4000 \times 10^{-2}, \quad -5.6780 \times 10^5$$

In the above examples, we are using base  $\beta = 10$  to ease into the material. Computers work more naturally in binary (base 2). When  $\beta = 2$ , each digit  $d_i$  is 0 or 1 in equation (1.1). I will leave it to the readers to review binary representation, but some examples are listed below:

- Conversion from binary to decimal:

a.  $11.0101_2 = 1 \cdot 2^1 + 1 \cdot 2^0 + 0 \cdot 2^{-1} + 1 \cdot 2^{-2} + 0 \cdot 2^{-3} + 1 \cdot 2^{-4} = 2 + 1 + 0.25 + 0.0625 = 3.3125$

b.  $1100_2 = 1.1_2 \times 2^3 = (1 \cdot 2^0 + 1 \cdot 2^{-1}) \cdot 2^3 = 12$

c.  $0.01_2 = 1.00_2 \times 2^{-2} = 0.25$

- Conversion from decimal to binary:

a.  $41 = 2^5 + 9 = 1 * 2^5 + 1 * 2^3 + 1 * 2^0 = 101001_2$

- b. To figure out the binary representation for 0.1, we first observe that  $2^{-4} = 0.0625$  is the biggest component of 0.1. The leftover  $0.1 - 0.0625 = 0.0375 > 2^{-5}$ . So  $0.1 = 2^{-4} + 2^{-5} + 0.00625 = 2^{-4} + 2^{-5} + 2^{-4} * 0.1 = 2^{-4} + 2^{-5} + 2^{-4}(2^{-4} + 2^{-5} + 2^{-4} * 0.1)$ . This can keep going, so

$$0.1 = 2^{-4} + 2^{-5} + 2^{-8} + 2^{-9} + \dots = 0.0001\overline{1}_2 = 1.\overline{1001}_2 \times 2^{-4}. \quad (1.2)$$

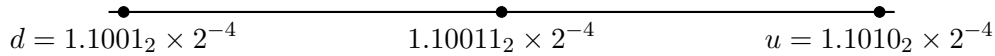
As seen, if a base (subscript) is not specified, then it is understood that it is a decimal number.

The representation in (1.2) has to be rounded at some point, and this is causing the computation inaccuracy at the beginning of this chapter.

There are usually 4 rounding models: rounding down, rounding up, rounding towards 0, and rounding to nearest. rounding to the nearest is, either round down or round up, whichever is closer. In case of a tie, we choose the one whose least significant (rightmost) bit is smaller. Given a real number  $x$ , we denote its floating point representation by  $\text{fl}(x)$ . For convenience of illustration, we use examples where we retain 4 digits after the binary point (this means the representation has precision  $p = 5$ ).

**Example 1.1.** Round  $0.1 = 1.1001100\dots_2 \times 2^{-4}$  to the nearest so that there are only 4 digits after the binary point. This means that in this machine,  $p = 5$ .

Rounding up will be  $u = 1.1010_2 \times 2^{-4}$ . Rounding down will be  $d = 1.1001_2 \times 2^{-4}$ . To see which one is nearest, we can compute the midpoint of  $u$  and  $d$ .



0.1 is clearly bigger than the middle point, so we pick  $1.1010_2 \times 2^{-4}$  with rounding to the nearest. We can write  $\text{fl}(0.1) = 1.1010_2 \times 2^{-4}$  in this particular machine.

**Question:** In this same machine, meaning we still have  $p = 5$  and we still use rounding to the nearest, what is  $\text{fl}(1.10011_2 \times 2^{-4})$ ? (Hint: in binary, we break the tie by choosing the one whose least significant (rightmost) bit is 0.)

In the model where  $\beta = 10, p = 5$ , 1 is represented as  $1.0000 \times 10^0$ , the number that is closest to 1 (and bigger than 1) is  $1.0001 \times 10^0$ . There is a gap  $10^{-4}$ . This gap is relative. The closest number to  $1.0000 \times 10^4$  is  $1.0001 \times 10^4$  (taking the bigger one again). This means we cannot represent any number between 10000 and 10001 with this machine. The gap ( $=1$ ) is a lot bigger now, but the relative gap is still  $1/10000 = 10^{-4}$ .

In general, with base  $\beta$  and precision level  $p$ , this gap (between 1 and its closest number) is defined to be the *machine precision* or *machine epsilon*  $\epsilon_m$ .<sup>1</sup> We have  $\epsilon_m = \beta^{1-p}$ . Machine epsilon is related to the *round-off error*, which is the error that results from replacing a number with its floating point form. The approximation  $x^*$  to  $x$  has *absolute error*  $|x - x^*|$  and *relative error*  $\frac{|x - x^*|}{|x|}$ , provided that  $x \neq 0$ .

**Example 1.2.** Let us use base 10 with  $p = 3$ , in which case the machine precision is  $\epsilon_m = 10^{-2}$ . If we use the rounding to nearest method, 1.345 becomes 1.34 (both 1.34 and 1.35 are equally

<sup>1</sup>Some sources define the machine precision to be half of this gap.

close, we break the tie by choosing the number whose right most bit is smaller). Relative error  $= \left| \frac{1.345 - 1.34}{1.345} \right| < \frac{0.005}{1} = 0.5 * \epsilon_m$  (The relative error is about 0.003.)

The number  $1.345 \times 10^4$  will be rounded to  $1.34 \times 10^4$ . The absolute error is very big, but the relative error is  $\left| \frac{1.345 \times 10^4 - 1.34 \times 10^4}{1.345 \times 10^4} \right| = \left| \frac{1.345 - 1.34}{1.345} \right|$ , the same as before.

**Example 1.3.** Following Example 1.1, the absolute error in representing 0.1 is  $|0.1 - \text{fl}(0.1)|$ , which is less than half of the difference between  $u$  and  $d$ . The relative error is

$$\frac{|0.1 - \text{fl}(0.1)|}{0.1} < \frac{2^{-1}(u - d)}{d} = \frac{2^{-5} \cdot 2^{-4}}{1.1001_2 \times 2^{-4}} < 2^{-5}.$$

If we use rounding to the nearest, one has

$$\frac{|\text{fl}(x) - x|}{|x|} \leq 0.5\epsilon_m. \quad (1.3)$$

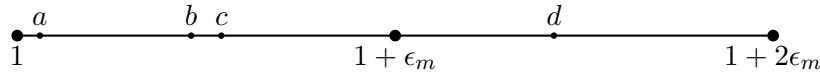
In other words, we may write

$$\text{fl}(x) = x(1 + \epsilon), \quad \text{where } |\epsilon| < 0.5\epsilon_m.$$

The idea of machine epsilon is that

- Two numbers may be indistinguishable in a machine if their relative error is less than  $0.5\epsilon_m$ ;
- Two numbers are definitely distinguishable in a machine if their relative error is more than  $\epsilon_m$ .

Let's consider four numbers  $a, b, c, d$  as indicated below:



Using rounding to the nearest,  $\text{fl}(a) = \text{fl}(b) = 1$  and  $a, b$  are within  $0.5\epsilon_m$  relative error.  $\text{fl}(b) \neq \text{fl}(c)$  even though they are within  $0.5\epsilon_m$  relative error.  $\text{fl}(c) = \text{fl}(d)$  even though they are almost  $\epsilon_m$  apart. The only way to guarantee that two numbers have different representations is make sure they are relatively  $\epsilon_m$  apart, like  $b$  and  $d$ .

## 1.2 Floating Point Operations

On a computer, all computations are reduced to  $+$ ,  $-$ ,  $\times$ , or  $/$ . Due to rounding, floating point arithmetic are different from the classical ones, and we will denote these *floating point operations* (*flop*) by  $\oplus, \ominus, \otimes, \oslash$ .

**Example 1.4** (Addition). Suppose we retain 4 digits after the binary point. To add up two numbers with different powers, we first need to adjust one of the numbers to align the mantissa (use bigger power), as:

$$1.1000_2 \times 2^1 \oplus 1.1001_2 \times 2^{-1} = 1.1000_2 \times 2^1 \oplus \text{fl}(0.011001_2) \times 2^1 = 1.1000_2 \times 2^1 \oplus 0.0110_2 \times 2^1 = \text{fl}(1.1110_2) \times 2^1 = 1.1110_2 \times 2^1.$$

**Fundamental axiom of floating point arithmetic:** Let  $*$  be one of the operations  $(+, -, \times, /)$ , and let  $\oplus$  be its floating point analogue, then for all **floating point numbers**  $x, y$ , we have that

$$x \oplus y = (x * y)(1 + \epsilon), \text{ for some } \epsilon \text{ such that } |\epsilon| \leq \epsilon_m.$$

or equivalently

$$\left| \frac{x \oplus y - x * y}{x * y} \right| \leq \epsilon_m.$$

We will not prove this, but we will quickly check this axiom in one example.

**Example 1.5.** We assume a machine, base 2, can only have 4 digits precision after the binary point, rounding to the nearest. Let  $x = 1.1010_2 \times 2^{-4}$  and  $y = 1.1010_2 \times 2^{-3}$  be two floating point numbers. ( $x$  is the floating point representation of 0.1 by Example 1.1 and  $y$  is the floating point representation of 0.2.)

Similar to Example 1.4, we first rewrite  $x$  as  $0.11010_2 \times 2^{-3}$ , which rounds to  $0.1110_2 \times 2^{-3}$  (break the tie by picking the one whose right most bit is 0).

$$x \oplus y = \text{fl}(0.1110_2 \times 2^{-3} + 1.1010_2 \times 2^{-3}) = \text{fl}(10.1000_2 \times 2^{-3}) = 1.0100_2 \times 2^{-2} = 0.3125.$$

The relative error for this addition is

$$\left| \frac{x \oplus y - (x + y)}{x + y} \right| = \left| \frac{1.0100_2 \times 2^{-2} - 1.001110_2 \times 2^{-2}}{1.001110_2 \times 2^{-2}} \right| = \frac{0.000010_2}{1.001110_2} < 2^{-5} < 2^{-4} = \epsilon_m.$$

Finally, we present the example of adding 0.1 and 0.2, but in a “very outdated” machine.

**Example 1.6.** If we assume a machine, base 2, precision 5. This is how it adds up 0.1 and 0.2.

Step 1: 0.1 is rounded to  $\text{fl}(0.1) = 1.1010_2 \times 2^{-4}$  as shown in Example 1.1.

Step 2:  $\text{fl}(0.2) = 1.1010_2 \times 2^{-3}$  since 0.2 is twice of 0.1.

Step 3:  $\text{fl}(0.1) \oplus \text{fl}(0.2) = 0.3125$  as shown in Example 1.5.

The overall relative error is  $|0.3 - 0.3125|/0.3 \approx 0.042$ .

To sum up, if we use  $m(\cdot)$  to indicate the machine output using floating point numbers, then

$$m(x + y) = \text{fl}(x) \oplus \text{fl}(y) = (x(1 + \epsilon_1) + y(1 + \epsilon_2))(1 + \epsilon_3)$$

### 1.3 IEEE 754 Standard

In 1985, the Institute for Electrical and Electronic Engineers (IEEE) published the *Binary Floating Point Arithmetic Standard 754-1985* which is followed by all microcomputer manufacturers. The double precision real numbers require a 64-bit representation. The first bit is a sign indicator, denoted  $s$ . This is followed by an 11-bit exponent  $c$ , and a 52-bit binary fraction  $f$  (also named mantissa or significand). See the following example:

$$0.1 = +1. \underbrace{10011001100110011001100110011001100110011010_2}_{52 \text{ digits mantissa}} \times 2^{-4},$$

This is 0.100000000000000055511151231257827021181583404541015625 in decimal. In the IEEE 754 standard,  $p = 53, \epsilon_m = 2^{-52} \approx 2.22 \times 10^{-16}$ . To retrieve this machine epsilon, type `eps = numpy.finfo(float).eps` in Python. Moreover, execute `1 + 0.4*eps`, what do you get?

The default IEEE standard is rounding to nearest. The 11-bit exponent can represent  $2^{11} = 2048$  numbers. The numbers 0 and 2047 are reserved for other use (like 0,  $\infty$ , NaN), so  $c$  can only take on integers ranging from 1 to 2046. In order to include both positive and negative exponents,





Due to rounding, subtracting two nearly equal numbers will cause a big relative error. Such an error is called cancellation error and we should try to avoid it if possible. As shown in Example 1.7, rationalizing a formula could be used when appropriate.

The following example is from Example 6 of Section 1.2 of [1].

**Example 1.8.** Evaluate  $f(x) = x^3 - 6.1x^2 + 3.2x + 1.5$  at  $x = 4.71$  using  $\beta = 10, p = 3$ .

The exact value is  $f(4.71) = -14.263899$ . With rounding, we will have

$$\text{fl}(4.71^3) - 6.1\text{fl}(4.71^2) + 3.2\text{fl}(4.71) + 1.5 = 104 - 135 + 15.1 + 1.50 = -14.4 \quad (1.7)$$

which produces a relative error  $|-14.4 + 14.263899|/14.263899 \approx 0.00954$ .

As an alternative approach,  $f(x)$  can be written in a nested manner as

$$f(x) = ((x - 6.1)x + 3.2)x + 1.5. \quad (1.8)$$

Evaluating (1.8) at  $x = 4.71$  using  $p = 3$  gives

$$\text{fl}\{\text{fl}\{[(4.71 - 6.1)4.71] + 3.2\}4.71\} + 1.5 = \text{fl}\{(-6.55 + 3.2)4.71\} + 1.5 = -15.8 + 1.5 = -14.3$$

whose relative error is reduced to 0.0025.

Why does the nested form provide a more accurate evaluation? This is because the nested form requires fewer flops to evaluate, and consequently fewer roundings. Indeed, the original evaluation (1.7) requires  $2+2+1$  multiplications and 3 additions, so a total of 8 flops. The nested form (1.8) requires only 5 flops.

Let  $p(x) = \sum_{i=0}^n a_i x^i$  be its original form, then it requires  $n + (n - 1) + \cdots + 2 + 1 = \frac{n(n+1)}{2}$

multiplications and  $n$  additions, so a total of  $\frac{n(n+1)}{2} + n = \frac{n(n+3)}{2}$  flops are needed.

It is easy to see that the nested form of  $p(x) = \sum_{i=0}^n a_i x^i$  is

$$p(x) = a_0 + x(a_1 + x(a_2 + x(a_3 + \cdots + x(a_{n-1} + xa_n) \cdots)))$$

In this form, the flop count is

$$\underbrace{2 + 2 + \cdots + 2}_n = 2n,$$

which is much less than the original form.

Nested multiplication should be performed whenever a polynomial is evaluated. This is because the nested form minimizes the number of flops needed, hence minimizing the rounding off error. The nested form is faster (fewer flops) and more accurate.

## Chapter 1 Exercises

1. Convert binary to decimal.

(a)  $110_2$

(b)  $11.0101_2$

2. Convert decimal to binary

- (a) 11.5
  - (b) 1.3125
3. Find the **normalized** floating point representation
- (a) 10.345, base 10,  $p = 4$ , rounding to the nearest.
  - (b) 10.395, base 10,  $p = 4$ , rounding to the nearest.
  - (c) 295, base 10,  $p = 2$ , rounding to the nearest.
  - (d)  $1.0101_2$ , base 2,  $p = 3$ , rounding up.
  - (e)  $1.0101_2$ , base 2,  $p = 3$ , rounding down.
  - (f)  $1.0101_2$ , base 2,  $p = 3$ , rounding to the nearest.
  - (g)  $1.0101_2$ , base 2,  $p = 4$ , rounding to the nearest.
4. Let us use base 10 with  $p = 2$ . If we use the rounding to the nearest method, what is the relative error of 1.23? What is the relative error of 10.61?
5. Following Exercise 3(g), what is the machine epsilon? what is the relative error in rounding this number? Does the rule (1.3) check out?
6. Directly add in binary (No rounding involved).
- (a)  $0.111_2 + 11.101_2$
  - (b)  $1.1010_2 \times 2^{-4} + 1.1010_2 \times 2^{-3}$ .
7. What is the biggest number that can be represented in IEEE 754 standard using double precision? Express this number in the form of  $2^a - 2^b$ .
8. In a machine of  $\beta = 2, p = 3$ , compute the machine result of  $0.1+0.2$  (rounding to the nearest).
9. In a machine of  $\beta = 10, p = 4$ , compute the machine result of  $6400.1 + 4.12$  (rounding to the nearest).
10. This problem is from Chapter 5 problem 15 of [2].

In the 1991 Gulf War, the Patriot missile defense system failed due to roundoff error. The troubles stemmed from a computer that performed the tracking calculations with an internal clock whose integer values in tenths of a second were converted to seconds by multiplying by a 24-bit binary approximation to 0.1:

$$0.1 \approx 0.00011001100110011001100_2$$

- (a) Convert the binary number above to a fraction. Call it  $x$ .
- (b) Compute the absolute difference between 0.1 and  $x$ .
- (c) What is the time error in seconds after 100 hours of operation (i.e., the value of  $|360,000 - 3,600,000x|$ )?

On February 25, 1991, a Patriot battery system, which was to protect the Dhahran Air Base, had been operating for over 100 consecutive hours. The roundoff error caused the system not to track an incoming Scud missile, which slipped through the defense system and detonated on US Army barracks, killing 28 American soldiers.

11. Find the most accurate approximation of the roots of  $x^2 - 75.7x + 1$  with  $p = 4$ . Compute the relative errors for both roots.
12. Given  $p(x) = 1 + 2x + 3x^2 + 3x^3 + 4x^4$ ,
  - (a) Find the nested form of  $p(x)$ .
  - (b) How many flops are needed to evaluate this polynomial at one point, in the original form, and in the nested form?

## Chapter 1 Python Exercises

13. (a) Define a vector  $v$  whose entries are  $[1, 0.01, 0.02, \dots, 5]$ . Do not print.  
 (b) Define a  $10 \times 10$  matrix whose entries are the integers 1 to 100 ordered column-wise. Print the matrix. [Hint: look up `np.reshape`]  
 (c) Plot  $y = \cos x$  and  $y = x$  on the interval  $[0, 1]$  in the same graph.
14. Given the equation  $x^2 - 100000x + 1 = 0$ , use formula (1.5) and formula (1.6) to compute two different approximations of  $x_1$ . Print both values in the exponential notation with 15 digits after the decimal point. Which formula produces a better approximation of  $x_1$ ?
15. Given  $p(x) = x^9 - 18x^8 + 144x^7 - 672x^6 + 2016x^5 - 4032x^4 + 5376x^3 - 4608x^2 + 2304x - 512$ .
  - (a) Let  $n(x)$  be the nested form of  $p(x)$ . Define  $p(x)$ ,  $n(x)$ , and  $s(x) = (x - 2)^9$ . All 3 functions need to be able to take in arrays.
  - (b)  $p(x) = n(x) = s(x)$  algebraically, but makes a difference in computers. Plot all three functions on the interval  $[1.92, 2.08]$  (with 0.001 increment). You need to place all 3 plots in a  $1 \times 3$  grid. Be sure to give a title to each subplot to indicate which is which.
  - (c) Explain the difference in three plots in (b).

## Chapter 2

# Solving nonlinear equations

If we recall the equations that we are able to solve, probably we would realize that there are very few. We can surely solve  $x^2 - x - 1 = 0$  using the famous quadratic formula, but have you ever wondered formulas for finding roots of polynomials whose degree is higher than 2? Unfortunately, even as simple as a polynomial equation  $x^5 - x - 1 = 0$ , a math Ph.D feels hopeless to find an exact answer by analysis. In fact, there are no formula for finding roots of polynomials of degree bigger than 4. Equations like  $\cos x = x$  do not have a closed form solution either.

A linear equation has the form  $ax = b$ . It is trivial to solve if  $x$  is a scalar. When  $x$  is a vector, we often write  $Ax = b$ , and is a major topic itself.

In this chapter, we discuss several iterative methods for solving nonlinear equations when  $x$  is a scalar. Direct methods are formula or procedures that will produce the exact solution. Iterative methods produce a sequence  $\{x_1, x_2, \dots, x_n, \dots\}$  that (hopefully) converges to the true solution  $x$ .

### 2.1 The Bisection Method

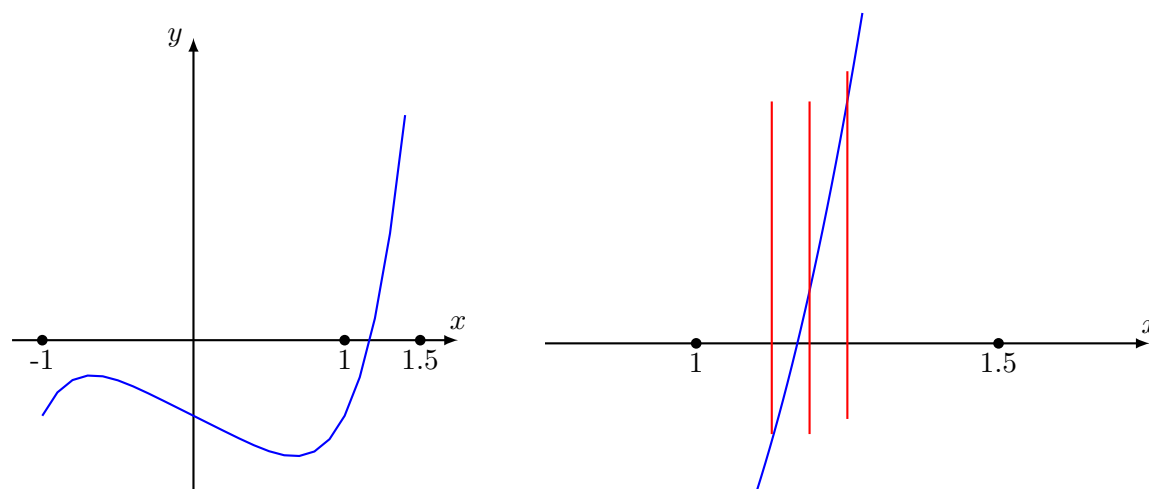


Figure 2.1: Bisection Method. The right is an enlargement of the left.

This is a simple, yet efficient method that is often covered in Calculus I as an application of the Intermediate Value Theorem (IVT). See Theorem 7.1.

Take  $f(x) = x^5 - x - 1$  as an example, whose graph is shown on the left of Figure 2.1. Since  $f(1) = -1 < 0$  and  $f(1.5) \approx 5.09 > 0$ , then by the IVT, there is  $c \in [1, 1.5]$  such that  $f(c) = 0$ , or

equivalently,  $f$  has a root  $c$  in  $[1, 1.5]$ .

To have a more precise estimate of  $c$ , we evaluate  $f$  at the middle point  $\frac{1+1.5}{2} = 1.25$ .  $f(1.25) \approx 0.8 > 0$ . By the IVT again, we conclude that the root  $c \in [1, 1.25]$ . We can keep bisecting the interval and eventually find the root to a sufficient precision. See the right of Figure 2.1.

Note that the key of the Bisection method is:

1. The function needs to be continuous.
2. One needs to be given initial points  $a, b$  such that  $f(a)$  and  $f(b)$  have different signs.

---

### The Bisection Algorithm

---

Input:  $a_0, b_0$  such that  $f(a_0), f(b_0)$  have different signs.

Output: an approximation of a root of  $f$

Repeat: for  $n \geq 0$

1.  $p_n = \frac{a_n + b_n}{2}$
  2. if  $f(a_n)f(p_n) < 0$ , set  $a_{n+1} = a_n, b_{n+1} = p_n$   
     otherwise, set  $a_{n+1} = p_n, b_{n+1} = b_n$
- until a stopping criteria has been reached
- 

**Example 2.1.** Find a solution of  $2^{-x} = x$ .

This is equivalent to finding a root of  $f(x) = (\frac{1}{2})^x - x$ . We have  $f(0) = 1 - 0 > 0, f(1) = 0.5 - 1 < 0$ , so we will let  $a_0 = 0, b_0 = 1$ . Therefore  $p_0 = 1/2$ .

$$f(1/2) = \frac{1}{\sqrt{2}} - \frac{1}{2} > 0, \text{ so } a_1 = 1/2, b_1 = 1, p_1 = 3/4.$$

$$f(3/4) = (0.5)^{3/4} - 3/4 < 0, \text{ so } a_2 = 1/2, b_2 = 3/4, p_2 = 5/8$$

After only two iterations, we get a coarse approximate of the root:  $p_2 = 5/8 = 0.625$ . For the reference, the actual root is 0.6411857445.

If we predetermine a tolerance level  $\text{tol}$  so that the algorithm will terminate if  $|p_n - p_{n+1}| < \text{Tol}$ , then this can be achieved when

$$|p_n - p_{n+1}| = \frac{b_n - a_n}{4} = \frac{b_0 - a_0}{2^{n+2}} < \text{Tol} \iff n \geq \log_2\left(\frac{b_0 - a_0}{4\text{Tol}}\right) \quad (2.1)$$

If  $p$  is the actual root, then we have

$$|p_n - p| \leq \frac{b_n - a_n}{2} = 2|p_n - p_{n+1}| < 2\text{Tol}.$$

Therefore the number of iterations needed can be computed via (2.1).

## 2.2 The Newton's Method

Newton's method (also called the Newton-Raphson method) is a little more sophisticated: it involves derivatives, but it is still Calculus material. This is usually how a calculator or a computer finds a root.

If we are to find the root of  $y = f(x)$ , we start with a random initial guess  $x_0$ . Consider the tangent line to the curve  $y = f(x)$ , at  $P_0 = (x_0, f(x_0))$ . The idea behind Newton's method is linearization. Since the tangent line (linearization of  $f$ ) is close to  $f$ , then its  $x$ -intercept,  $x_2$ , is close to the  $x$ -intercept of the curve  $y = f(x)$ . We can easily find  $x_2$ .

The tangent line at the point  $(x_0, f(x_0))$  is

$$y - f(x_0) = f'(x_0)(x - x_0).$$

Plugging in  $y = 0$  to solve for the  $x$ -intercept:

$$0 - f(x_0) = f'(x_0)(x - x_0) \implies x = x_0 - \frac{f(x_0)}{f'(x_0)}.$$

So  $x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$ . We use  $x_1$  as the next approximation and keep repeating this process, as shown in Figure 2.2, where we gain this iteration formula

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}, \quad k \geq 0. \quad (2.2)$$

Note that we require  $f'(x_k) \neq 0$  which is not a trivial condition.

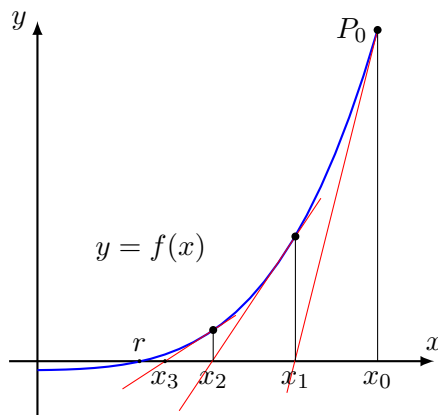


Figure 2.2: Newton's method

**Example 2.2.** Use the Newton's method to approximate  $\sqrt{2}$ .

Let  $f(x) = x^2 - 2$ , which makes  $\sqrt{2}$  a root of  $f(x)$ . We can compute that  $f'(x) = 2x$ . Let  $x_0 = 4$ . We can simplify (2.2) to

$$x_{k+1} = x_k - \frac{x_k^2 - 2}{2x_k} = \frac{x_k}{2} + \frac{1}{x_k}.$$

$$x_1 = 2 + 1/4 = 9/4 = 2.25$$

$$x_2 = 9/8 + 4/9 = 113/72 = 1.569\dot{4}$$

$$x_3 = 1.569\dot{4}/2 + 1/1.569\dot{4} \approx 1.42189$$

If we use the Bisection method for this problem starting at  $a_0 = 0, b_0 = 4$ , we then have  $p_0 = 2, p_1 = 1, p_2 = 1.5, p_3 = 1.25$ . This is a much worse approximation after the same number of iterations. See Table 2.1 for error comparison. This table also includes the secant method that will be introduced later.

The two drawbacks of the Newton's method are: (1) It requires the knowledge of the derivative function; (2) It may not converge if you have an unlucky initial guess  $x_0$ . For example, if one picks  $x_0 = 0$  for finding root of  $y = x^5 - x - 1$ , the sequence will not converge to the root (try to draw tangent lines on your own). However, if the initial guess is close enough, the Newton's method does enjoy fast convergence (faster than the Bisection method as shown in Example 2.2) .

$n$	Bisection $ p_n - \sqrt{2} $	Newton $ x_n - \sqrt{2} $	Secant
0	0.585786438	2.585786438	
1	0.414213562	0.835786438	
2	0.085786438	0.155230882	0.414213562
3	0.164213562	0.007676801	0.080880229
4	0.039213562	0.000020723	0.014357866
5	0.023286438	0.000000000	0.000420459

Table 2.1: Accuracy Comparison of Example 2.2

**Definition 2.3.** A method that produces a sequence  $\{x_n\}$  that converges to a number  $x$  *linearly* if there exists  $L \in (0, 1)$  such that for large values of  $n$

$$|x_{n+1} - x| \leq L|x_n - x|.$$

The sequence converges to  $x$  *quadratically* if there exists  $Q > 0$  such that for large values of  $n$

$$|x_{n+1} - x| \leq Q|x_n - x|^2.$$

For example, the sequence  $\{3^{-n}\}_{n=1}^{\infty}$  converges to 0 linearly. The sequence  $\{b_n = 3^{-2^n}\}_{n=1}^{\infty}$  converges quadratically to 0.

**Theorem 2.4.** Let  $f(x)$  be a function whose second derivative is continuous, and  $r$  be a root of  $f(x)$ . Let  $I$  be an interval contains  $r$  such that  $|f''(x)| \leq M$  for any  $x \in I$ . If  $x_k \in I$ , then the Newton's method converges to  $r$  quadratically as

$$|x_{k+1} - r| \leq \frac{M}{2|f'(x_k)|} |x_k - r|^2. \quad (2.3)$$

*Proof.* Expand  $f$  in the first Taylor polynomial ( $n = 1$  in (7.1)) at  $x_k$  gives

$$0 = f(r) = f(x_k) + f'(x_k)(r - x_k) + \frac{f''(\xi)}{2}(r - x_k)^2,$$

where  $\xi$  lies between  $x_k$  and  $r$ . Consequently, if  $f'(x_k) \neq 0$ , we have

$$r - x_k + \frac{f(x_k)}{f'(x_k)} = -\frac{f''(\xi)}{2f'(x_k)}(r - x_k)^2.$$

Combined with (2.2), we have

$$r - x_{k+1} = -\frac{f''(\xi)}{2f'(x_k)}(r - x_k)^2 \implies |r - x_{k+1}| \leq \frac{M}{2|f'(x_k)|} |r - x_k|^2$$

since  $\xi$  is in between  $x_k$  and  $r$ , hence  $\xi \in I$ . □

This implies that the Newton's method has the tendency to approximately double the number of digits of accuracy with each successive approximation, as demonstrated in Table 2.1.

The order of convergence for linear methods (like Bisection) is 1 and the order of convergence for quadratic methods (like Newton) is 2.



## 2.3 The Secant Method

As mentioned, one of the drawbacks of the Newton's method is that we have to know its derivative. For the simple example presented above, this is not difficult. In many problems, however,  $f$  can get complicated, or worse, cannot be written down analytically. In such scenarios, differentiating  $f$  analytically is difficult or impossible. For such problems, we approximate the derivatives. Iterations of the form

$$x_{k+1} = x_k - \frac{f(x_k)}{g(x_k)}, \quad \text{where } g_k \approx f'(x_k)$$

are called *quasi-Newton* methods.

Among many quasi-Newton methods, the *secant method* is defined by taking  $f'(x_k)$  to be

$$f'(x_k) \approx \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}.$$

So the iteration formula for the secant method is:

$$x_{k+1} = x_k - \frac{f(x_k)(x_k - x_{k-1})}{f(x_k) - f(x_{k-1})}, \quad k \geq 1. \quad (2.4)$$

To begin the secant method, one needs two starting points  $x_0, x_1$ .

It turns out that the order of convergence of the secant method is in between Bisection and Newton. In fact, the order is  $\frac{1+\sqrt{5}}{2} \approx 1.62$  the Golden ratio as for the secant method the iterates satisfy

$$|x_{n+1} - x| \leq M|x_n - x|^{\frac{1+\sqrt{5}}{2}}.$$

Interested readers can find a proof in [2, Section 4.3.3].

We will use the secant method on the same example.

**Example 2.5.** Use the secant method to find the positive root of  $f(x) = x^2 - 2$  with initial  $x_0 = 0, x_1 = 2$ .

$$x_{k+1} = x_k - \frac{f(x_k)(x_k - x_{k-1})}{f(x_k) - f(x_{k-1})} = x_k - \frac{(x_k^2 - 2)(x_k - x_{k-1})}{x_k^2 - x_{k-1}^2} = x_k - \frac{x_k^2 - 2}{x_k + x_{k-1}} = \frac{x_k x_{k-1} + 2}{x_k + x_{k-1}}.$$

$$x_2 = \frac{0 + 2}{2} = 1.$$

$$x_3 = \frac{2 + 2}{1 + 2} = 4/3.$$

$$x_4 = \frac{4/3 + 2}{7/3} = 10/7.$$

$$x_5 = \frac{\frac{10}{7} \cdot \frac{4}{3} + 2}{\frac{10}{7} + \frac{4}{3}} = \frac{40 + 42}{30 + 28} = \frac{41}{29}.$$

The errors are listed in the third column of Table 2.1.

## Chapter 2 Exercises

### Python Exercises

1. Bisection and Newton comparison.

- (a) Use the Bisection method to do Example 2.1, with the stopping criteria  $|p_{n+1} - p_n| \leq 1e-5$ , and initial interval  $a_0 = 0, b_0 = 1$ . How many iterations was run? (meaning what is the  $n$  value when  $|p_{n+1} - p_n| \leq 1e-5$  is reached?) How is this compared to the theoretical value in (2.1)?
- (b) Use the Newton's method to do Example 2.1, with the stopping criteria  $|x_{n+1} - x_n| \leq 1e-5$ , and initial value  $x_0 = 1$ . How many iterations was run?
2. Newton's method fails for finding a root of  $y = x^5 - x - 1$  with  $x_0 = 0$ . Print the first 6 approximations to illustrate this.

### Other Exercises

3. Use the Bisection method to find a solution of  $\cos x = x$  with initial interval  $[0, \pi]$ . Find  $p_i, i = 0, 1, 2$ . Use (2.1) to compute how many iterations are needed if  $\text{tol} = 0.01$ .
4. Let  $f(x) = 3(x+1)(x-1/2)(x-1)$ . Use the Bisection method to find  $p_3$  on the interval  $[-2, 1.5]$ . [Since  $f$  is factored, we know its roots right away. This problem is more of an exercise such that the function's sign can be determined easily.]
5. Write down the first 3 iterations  $(x_1, x_2, x_3)$  of the Newton's method for solving  $x^2 - 3 = 0$ , starting with  $x_0 = 1$ . Write your answers in terms of fractions.
6. Prove that Newton's method will converge to 0 given any initial value  $x_0$  if we are solving  $x^2 = 0$ .
7. Write down the first three iterates (2, 3, 4) of the secant method for solving  $x^2 - 3 = 0$ , starting with  $x_0 = 0$  and  $x_1 = 1$ . Write your answer in terms of fractions.
8. Show that  $\{3^{-n}\}_{n=1}^{\infty}$  converges to 0 linearly, and the sequence  $\{b_n = 3^{-2^n}\}_{n=1}^{\infty}$  converges to 0 quadratically.
9. We can compute  $1/3$  by solving  $f(x) = 0$  with  $f(x) = x^{-1} - 3$ .
  - (a) Write down the Newton iteration for this problem, and compute by hand the first 2 Newton iterates for approximating  $1/3$ , starting with  $x_0 = 0.5$ .
  - (b) What happens if you start with  $x_0 = 1$ ?
  - (c) \*In the case of (b), show that the iterates  $x_k \rightarrow -\infty$  as  $k \rightarrow \infty$ .
10. Problems involving the amount of money required to pay off a mortgage over a fixed period of time involve the formula

$$A = \frac{P}{i}[1 - (1+i)^{-n}],$$

known as an *ordinary annuity equation*. In this equation  $A$  is the amount of the mortgage,  $P$  is the amount of each payment, and  $i$  is the interest rate per period for the  $n$  payment periods. Suppose that a 30-year home mortgage in the amount of \$135,000 is needed and that the borrower can afford house payments of at most \$1000 per month. What is the maximal interest rate the borrower can afford to pay? (You can use `scipy.optimize.newton` to solve the equation in the end.)

## Chapter 3

# Solving system of linear equations: Direct methods

Solving a linear system  $Ax = b$  is everywhere in many applications. There are generally two approaches: Direct methods and iterative (indirect) methods. Direct methods give a direct answer at the end of the algorithm, whereas iterative methods follow an iterative formula where the iterates are approaching the real solution.

We also care about how efficient an algorithm is. This is to count how many flops are needed in a computer. Each flop is one of the following operations:  $+$ ,  $-$ ,  $\times$ ,  $\div$ . For example, the inner product  $x^T y$  requires  $n$  multiplications and  $n - 1$  summations, given both vectors have  $n$  coordinates. So  $x^T y$  needs  $2n - 1$  flops.

- Matrix-Vector multiplication: Given  $A \in \mathbb{R}^{m \times n}$ ,  $x \in \mathbb{R}^n$ ,  $Ax$  is computing  $m$  inner products, so needs  $m(2n - 1)$  flops.
- Matrix-Matrix multiplication: Given  $A \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{n \times p}$ ,  $AB$  is computing  $p$  matrix-vector multiplications. So the calculation of  $AB$  requires  $pm(2n - 1)$  flops. In particular, if both  $A$  and  $B$  are square matrices, then  $AB$  needs  $2n^3 - n^2 \approx 2n^3$  flops.

### 3.1 Gaussian Elimination and $LU$ factorization

Gaussian elimination is a core material in linear algebra and is the procedure one performs when finding key features like linear dependence, rank, null space, basis, dimension, etc. Gauss elimination is the simplest way to solve linear systems of equations by hand, and also the standard method for solving them on computers.

The following notations are standard:

$$Ax = b \iff \begin{array}{ccccccc} a_{11}x_1 & +a_{12}x_2 & +\cdots & +a_{1n}x_n & = & b_1 \\ a_{21}x_1 & +a_{22}x_2 & +\cdots & +a_{2n}x_n & = & b_2 \\ \vdots & \vdots & & \vdots & & \vdots \\ a_{m1}x_1 & +a_{m2}x_2 & +\cdots & +a_{mn}x_n & = & b_m \end{array} \iff [A|b].$$

The matrix  $A$  is called the coefficient matrix and the matrix  $[A|b]$  is called the augmented matrix.

Gaussian elimination is suitable for systems of all sizes, but **in this section our examples and analysis will focus on  $Ax = b$  where  $A$  is invertible**. When solving  $Ax = b$ , the Gaussian elimination consists of two procedures: forward elimination and backward substitution.

**Example 3.1.** Solve  $\left[ \begin{array}{ccc|c} 1 & 2 & 2 & 1 \\ 2 & 7 & 7 & 5 \\ 2 & 7 & 9 & 5 \end{array} \right]$ .

Forward Elimination:

$$\left[ \begin{array}{ccc|c} 1 & 2 & 2 & 1 \\ 2 & 7 & 7 & 5 \\ 2 & 7 & 9 & 5 \end{array} \right] \xrightarrow{\rho 1(-2)+\rho 2} \left[ \begin{array}{ccc|c} 1 & 2 & 2 & 1 \\ 0 & 3 & 3 & 3 \\ 2 & 7 & 9 & 5 \end{array} \right] \xrightarrow{\rho 1(-2)+\rho 3} \left[ \begin{array}{ccc|c} 1 & 2 & 2 & 1 \\ 0 & 3 & 3 & 3 \\ 0 & 3 & 5 & 3 \end{array} \right] \xrightarrow{\rho 2(-1)+\rho 3} \left[ \begin{array}{ccc|c} 1 & 2 & 2 & 1 \\ 0 & 3 & 3 & 3 \\ 0 & 0 & 2 & 0 \end{array} \right].$$

Backward substitution:

(1) Equation form:

$$2x_3 = 0 \Rightarrow x_3 = 0$$

$$3x_2 + 3 \cdot 0 = 3 \Rightarrow x_2 = 1$$

$$x_1 + 2 \cdot 1 + 2 \cdot 0 = 1 \Rightarrow x_1 = -1$$

(2) Matrix form:

$$\left[ \begin{array}{ccc|c} 1 & 2 & 2 & 1 \\ 0 & 3 & 3 & 3 \\ 0 & 0 & 2 & 0 \end{array} \right] \xrightarrow{\rho 3/2} \left[ \begin{array}{ccc|c} 1 & 2 & 2 & 1 \\ 0 & 3 & 3 & 3 \\ 0 & 0 & 1 & 0 \end{array} \right] \xrightarrow[\rho 3(-2)+\rho 1]{\rho 3(-3)+\rho 2} \left[ \begin{array}{ccc|c} 1 & 2 & 0 & 1 \\ 0 & 3 & 0 & 3 \\ 0 & 0 & 1 & 0 \end{array} \right] \xrightarrow{\rho 2/(3)} \left[ \begin{array}{ccc|c} 1 & 2 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{array} \right] \xrightarrow{\rho 2(-2)+\rho 1} \left[ \begin{array}{ccc|c} 1 & 0 & 0 & -1 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{array} \right].$$

As seen in Example 3.1, the Gaussian elimination consists of two procedures: Forward elimination and backward substitution. Next, we will describe them generally and also count the flops used in each procedure.

The following two formula will be useful:

$$\sum_{i=1}^n i = \frac{n(n+1)}{2}, \quad \sum_{i=1}^n i^2 = \frac{n(n+1)(2n+1)}{6}. \quad (3.1)$$

Forward elimination of Gaussian elimination: This is to get to an upper triangular matrix with row operations.

**Column 1**

We eliminate all entries below  $a_{11}$ .

$$\left[ \begin{array}{cccc|c} a_{11} & a_{12} & \cdots & a_{1n} & b_1 \\ a_{21} & a_{22} & \cdots & a_{2n} & b_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} & b_n \end{array} \right] \longrightarrow \left[ \begin{array}{cccc|c} a_{11} & a_{12} & \cdots & a_{1n} & b_1 \\ 0 & a'_{22} & \cdots & a'_{2n} & b'_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & a'_{n2} & \cdots & a'_{nn} & b'_n \end{array} \right] \quad (3.2)$$

At each row, we need  $n+1$  multiplications and  $n$  additions. [for example, to obtain the new row 2, we first compute  $-\frac{a_{21}}{a_{11}}$ , which takes one multiplication/division. Then we add  $-\frac{a_{21}}{a_{11}}$  copy of row 1 to row 2.] There are  $n-1$  rows to perform, so a total of  $(n-1) \cdot (2n+1) = 2n^2 - n - 1$  flops are needed to get to the right hand side of (3.2).

**Column 2:** Now we use  $a'_{22}$  to eliminate all entries below. This is the same work except we have an  $(n-1) \times (n-1)$  system now, so we need  $2(n-1)^2 - (n-1) - 1$  flops.

**Total:** We keep doing this until the coefficient matrix becomes upper triangular, so there are a total of  $\sum_{i=1}^n (2i^2 - i - 1) = \frac{n(n+1)(2n+1)}{3} - \frac{n(n+1)}{2} - n = \frac{2}{3}n^3 + \frac{1}{2}n^2 - \frac{7}{6}n \approx \frac{2}{3}n^3$  flops.

**Backward substitution of Gaussian Elimination:** After forward elimination, we obtain an upper triangular system

$$\left[ \begin{array}{cccc|c} c_{11} & c_{12} & \cdots & c_{1n} & y_1 \\ 0 & c_{22} & \cdots & c_{2n} & y_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & c_{nn} & y_n \end{array} \right]. \quad (3.3)$$

The backward substitution can be done in either the equation form or the matrix form, both of which are illustrated in Example 3.1. We will use the equation form in our operation counts here. We will start with last row and work our way up.

**Row  $n$ :**  $c_{nn}x_n = y_n$ . Takes 1 division to solve  $x_n$

**Row  $n-1$ :**  $c_{n-1,n-1}x_{n-1} + c_{n-1,n}x_n = y_{n-1}$ . Takes 1 multiplication, 1 addition, 1 division to solve  $x_{n-1}$ .

**Row  $n-2$ :**  $c_{n-2,n-2}x_{n-2} + c_{n-2,n-1}x_{n-1} + c_{n-2,n}x_n = y_{n-2}$ . Takes 2 multiplication, 2 addition, 1 division to solve  $x_{n-2}$ .

...

**Row 1:**  $c_{11}x_1 + c_{12}x_2 + \cdots + c_{1n}x_n = y_1$ . Takes  $n-1$  multiplications,  $n-1$  additions, 1 division to solve  $x_1$ .

Add up all counts, we have

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

**Gaussian elimination operation counts:**

$$\underbrace{\frac{2}{3}n^3 + \frac{1}{2}n^2 - \frac{7}{6}n}_{\text{forward elim.}} + \underbrace{n^2}_{\text{backward sub.}} = \frac{2}{3}n^3 + \frac{3}{2}n^2 - \frac{7}{6}n \approx \frac{2}{3}n^3.$$

### 3.1.1 Forward elimination = LU factorization

Each row operation corresponds to a matrix (often called an elementary row operation matrix). In

the forward elimination process of Example 3.1, let  $A = \begin{bmatrix} 1 & 2 & 2 \\ 2 & 7 & 7 \\ 2 & 7 & 9 \end{bmatrix}$ , then  $L_1A = \begin{bmatrix} 1 & 2 & 2 \\ 0 & 3 & 3 \\ 2 & 7 & 9 \end{bmatrix}$ ,

where  $L_1$  represents the first row operation. To figure out  $L_1$ , we perform the very row operation on

the identity matrix:  $\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \xrightarrow{\rho 1(-2)+\rho_2} \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = L_1$ . For the rest two row operations,

we can compute in the same fashion that  $L_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -2 & 0 & 1 \end{bmatrix}$ ,  $L_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -1 & 1 \end{bmatrix}$ .

The forward elimination process of Example 3.1 can be written as

$$L_3 L_2 L_1 A = U = \begin{bmatrix} 1 & 2 & 2 \\ 0 & 3 & 3 \\ 0 & 0 & 2 \end{bmatrix}.$$

So

$$A = (L_3 L_2 L_1)^{-1} U = L_1^{-1} L_2^{-1} L_3^{-1} U := LU.$$

$L_1, L_2, L_3$  are very special matrices, and their inverses and  $L$  can be figured out easily:

$$L_1^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, L_2^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 2 & 0 & 1 \end{bmatrix}, L_3^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}, L = L_1^{-1} L_2^{-1} L_3^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 2 & 1 & 1 \end{bmatrix}.$$

So we can decompose/factor  $A$  as

$$A = \begin{bmatrix} 1 & 2 & 2 \\ 2 & 7 & 7 \\ 2 & 7 & 9 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 2 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 2 \\ 0 & 3 & 3 \\ 0 & 0 & 2 \end{bmatrix} = LU.$$

In general, if there is no row exchange involved in forward elimination, then  $A$  can always be written as the product of a lower triangular matrix and an upper triangular matrix. This is called the  $LU$  factorization of  $A$ . It requires  $\frac{2}{3}n^3$  flops as well because it is done through forward elimination. In fact, the operation counts is slightly less because we don't need to consider the right hand side  $b$ , but it is still on the order of  $\frac{2}{3}n^3$ .

### 3.1.2 A 'second way' to solve $Ax = b$

If the  $LU$  factorization of  $A$  is already done, then it is very easy to solve  $LUx = b$ . We set  $Ux = y$ . We first solve  $y$  from  $Ly = b$ , which is very straightforward as we will see in Example 3.2. Then we solve  $x$  from  $Ux = y$ . We will use the same example to illustrate.

**Example 3.2.** Suppose we are given  $\begin{bmatrix} 1 & 2 & 2 \\ 2 & 7 & 7 \\ 2 & 7 & 9 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 2 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 2 \\ 0 & 3 & 3 \\ 0 & 0 & 2 \end{bmatrix}$ . In order to solve

$$\begin{bmatrix} 1 & 2 & 2 \\ 2 & 7 & 7 \\ 2 & 7 & 9 \end{bmatrix} x = \begin{bmatrix} 1 \\ 5 \\ 5 \end{bmatrix},$$

$$(1) \text{ We first solve } \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 2 & 1 & 1 \end{bmatrix} y = \begin{bmatrix} 1 \\ 5 \\ 5 \end{bmatrix} \text{ (} Ly = b \text{)}.$$

$$y_1 = 1$$

$$2 * 1 + y_2 = 5 \Rightarrow y_2 = 3$$

$$2 * 1 + 1 * 3 + y_3 = 5 \Rightarrow y_3 = 0$$

(2) Then we solve  $Ux = y$ , which is already done in the backward substitution of Example 3.1. It is not a coincidence that the  $y$  we solved in (1) is the same as the right hand side of the system after forward elimination.

It is obvious that solving  $Ly = b$  and  $Ux = y$  cost the same number of operations because the only difference is that  $Ly = b$  is solved forwardly, and  $Ux = y$  is solved backwardly.

#### **Solve $Ax = b$ via LU decomposition**

Step 1: $A = LU$ . (Now equation becomes $LUx = b$ .)	$\frac{2}{3}n^3$ flops
Step 2: Solve $Ly = b$ . ( $y$ will equal to $Ux$ .)	$n^2$ flops
Step 3: Solve $Ux = y$ .	$n^2$ flops

$LU$  decomposition approach is better when we need to solve multiple systems with the same coefficient matrix  $A$ .

We compare the flop counts of these two approaches for solving  $m$  systems:

$$Ax = b_1, Ax = b_2, \dots, Ax = b_m.$$

Regular Gauss Elimination:	$\frac{2}{3}n^3m$ flops
$LU$ approach:	$\frac{2}{3}n^3 + 2mn^2$ flops

If  $m = n$ , then regular GE approach takes around  $\frac{2}{3}n^4$  operations, whereas  $LU$  approach operation count is  $\frac{8}{3}n^3$ , still on the order of  $n^3$ . This makes a huge difference when the system is large (This summary is also included in Table 3.2).

We list in Table 3.1 the flop count of various matrix algebra.

Table 3.1: Flop count for various matrix operations:  $x, y \in \mathbb{R}^n, A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{n \times p}$ .

Operation	flop count	complexity	complexity if $m = n$
$x^T y$ ( $x$ dot product $y$ )	$= 2n - 1$	$O(n)$	$O(n)$
$Ax$	$= m(2n - 1)$	$O(mn)$	$O(n^2)$
$AB$	$= mp(2n - 1)$	$O(mpn)$	$O(n^3)$

We summarize in Table 3.2 the flop count related to solving a linear system.

### **3.2 Pivoting (row exchange)**

$LU$  decomposition is not the full story because it does not consider row exchanges, which is necessary in many scenarios.

Consider computing the  $LU$  factorization of  $A = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 2 & 1 \\ 2 & 7 & 9 \end{bmatrix}$ . We run into problem right away because  $a_{11} = 0$  is not able to eliminate entries below. In this case, a permutation of rows are necessary, and we can, for instance, switch row 1 and row 2. Switching row 1 and row 2 is the same as being left multiplied by the permutation matrix  $P = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ , i.e.  $PA = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 1 & 1 \\ 2 & 7 & 9 \end{bmatrix}$ . Then we can do the same old  $LU$  factorization to  $PA$ .

Table 3.2: Flop count related to solving a linear system:  $A \in \mathbb{R}^{n \times n}, B = [b_1, \dots, b_m] \in \mathbb{R}^{n \times m}$ .

Operation	flop count	complexity
Forward Elim. on $A$	$\approx \frac{2}{3}n^3$	$O(n^3)$
Backward Sub.	$= n^2$	$O(n^2)$
LU decomp of $A$	$\approx \frac{2}{3}n^3$	$O(n^3)$
Solve $Ax = b$ via Gauss Elim./LU	$\approx \frac{2}{3}n^3$	$O(n^3)$
Solve $AX = B$ by LU	$\approx \frac{2}{3}n^3 + 2mn^2$	$O(n^3 + mn^2)$
Finding $A^{-1}$ by reducing $[A I]$ to $[I A^{-1}]$	$\approx 2n^3$	$O(n^3)$

In the following example, it is not necessary to switch rows because the entry in the first row/first column is nonzero. But we will find out that switching rows will largely increase the accuracy of the solution.

**Example 3.3.** Let us look at another system  $\left[ \begin{array}{cc|c} 10^{-20} & 1 & 1 \\ 1 & 1 & 2 \end{array} \right]$ . The solution should be very close to  $x_1 = x_2 = 1$ , and we are able to get this solution if the computer had infinite precision. But with the forward elimination procedure and with double-precision arithmetic, we get

$$\left[ \begin{array}{cc|c} 10^{-20} & 1 & 1 \\ 1 & 1 & 2 \end{array} \right] \xrightarrow{\rho_1(-10^{20})+\rho_2} \left[ \begin{array}{cc|c} 10^{-20} & 1 & 1 \\ 0 & -10^{20} + 1 & -10^{20} + 2 \end{array} \right] \xrightarrow{\text{round}} \left[ \begin{array}{cc|c} 10^{-20} & 1 & 1 \\ 0 & -10^{20} & -10^{20} \end{array} \right]$$

The relative difference between 1 and  $-10^{20}$  is  $1/10^{20}$ , which is smaller than the machine epsilon. In other words, 1 is negligible compared to  $10^{20}$  (or  $-10^{20}$ ). You are encouraged to type `-1e20 + 1` in Python and see what happens.

From there, we solve  $x_2 = 1$ ,  $10^{-20}x_1 + 1 = 1 \implies x_1 = 0$ , which is far from the actual solution.

A remedy of this is to exchange these two rows:

$$\left[ \begin{array}{cc|c} 1 & 1 & 2 \\ 10^{-20} & 1 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 1 & 2 \\ 0 & -10^{-20} + 1 & -2 \times 10^{-20} + 1 \end{array} \right] \xrightarrow{\text{round}} \left[ \begin{array}{cc|c} 1 & 1 & 2 \\ 0 & 1 & 1 \end{array} \right]$$

And we can get a much more accurate solution  $x_1 = x_2 = 1$ .

This row exchange strategy is the content of the next subsection.

### 3.2.1 Partial pivoting

Partial pivoting is a straightforward method of row exchanging in the forward process of Gaussian elimination. At stage  $k$  (column  $k$ ), we will look for the biggest entry in absolute value of that column, and use that entry as a pivot. In Example 3.3, at stage 1, we compare all entries in column 1, and 1 is bigger than  $10^{-20}$ , so we exchange row 1 and row 2 to use 1 as the pivot. The next example is  $3 \times 3$  so this algorithm can be explained more clearly

**Example 3.4.** Use partial pivoting to solve the system  $\left[ \begin{array}{ccc|c} 2 & 1 & 1 & 1 \\ 4 & 3 & 3 & 1 \\ 8 & 7 & 9 & -1 \end{array} \right]$ .

First we focus on column 1, 8 is the biggest number in absolute value, so we switch row 1 and row 3, and then do two row operations.



$$\left[ \begin{array}{ccc|c} 2 & 1 & 1 & 1 \\ 4 & 3 & 3 & 1 \\ 8 & 7 & 9 & -1 \end{array} \right] \xrightarrow{\rho_1 \leftrightarrow \rho_3} \left[ \begin{array}{ccc|c} 8 & 7 & 9 & -1 \\ 4 & 3 & 3 & 1 \\ 2 & 1 & 1 & 1 \end{array} \right] \xrightarrow[\rho_1(-1/4)+\rho_3]{\rho_1(-1/2)+\rho_2} \left[ \begin{array}{ccc|c} 8 & 7 & 9 & -1 \\ 0 & -1/2 & -3/2 & 3/2 \\ 0 & -3/4 & -5/4 & 5/4 \end{array} \right]$$

Now we compare the entries of column 2. Keep in mind that we only focus on the smaller system so the first row is out of the picture. We compare the magnitude of  $-1/2$  and  $-3/4$ .  $-3/4$  is bigger so we have to switch rows again.

$$\left[ \begin{array}{ccc|c} 8 & 7 & 9 & -1 \\ 0 & -1/2 & -3/2 & 3/2 \\ 0 & -3/4 & -5/4 & 5/4 \end{array} \right] \xrightarrow{\rho_2 \leftrightarrow \rho_3} \left[ \begin{array}{ccc|c} 8 & 7 & 9 & -1 \\ 0 & -3/4 & -5/4 & 5/4 \\ 0 & -1/2 & -3/2 & 3/2 \end{array} \right] \xrightarrow{\rho_2(2/3)+\rho_3} \left[ \begin{array}{ccc|c} 8 & 7 & 9 & -1 \\ 0 & -3/4 & -5/4 & 5/4 \\ 0 & 0 & -2/3 & 2/3 \end{array} \right]$$

The rest backward substitution is the same as before:

$$\begin{aligned} x_3 &= -1 \\ -\frac{3}{4}x_2 + \frac{5}{4} &= \frac{5}{4} \implies x_2 = 0 \\ 8x_1 + 0 - 9 &= -1 \implies x_1 = 1 \end{aligned}$$

Note that when solving this problem by hand, partial pivoting does make the arithmetic harder due to all the fractions, but it is the same to a computer.

### 3.2.2 PLU factorization

With row exchanges, we have PLU factorization instead of LU factorization. For every matrix  $A$ , we are able to write

$$PA = LU,$$

where  $P$  is a permutation matrix,  $L$  is lower triangular and  $U$  is upper triangular.

First of all, we need to talk more about permutation matrices. A matrix  $P$  is a permutation matrix if it is a row permutation of the identity matrix.  $\begin{bmatrix} & & 1 \\ & 1 & \\ 1 & & \end{bmatrix}$  is a  $3 \times 3$  permutation matrix

because it is switching row 1 and row 3 of the identity matrix.  $\begin{bmatrix} & & & 1 \\ & 1 & & \\ & & & 1 \\ 1 & & & \end{bmatrix}$  is a  $4 \times 4$  permutation

matrix and it represents “place row 1 to row 3; place row 2 to row 1; place row 3 to row 4; place row 4 to row 2”, since this is the result of the identity matrix after doing these row operations.

Another useful fact is that for any permutation matrix  $P$ , we have  $P^{-1} = P^T$ .

Before figuring out the PLU factorization, it is also very useful to discuss the difference between left and right multiplications. In short, **left multiplication by an elementary matrix is to operate on the rows; right multiplication by an elementary matrix is to operate on the columns**. We have been left multiplying a matrix since we have been doing row operations. Take a permutation matrix for example,

$$\begin{bmatrix} & & 1 \\ & 1 & \\ 1 & & \end{bmatrix} \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} = \begin{bmatrix} g & h & i \\ d & e & f \\ a & b & c \end{bmatrix}, \quad \text{and} \quad \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \begin{bmatrix} & & 1 \\ & 1 & \\ 1 & & \end{bmatrix} = \begin{bmatrix} c & b & a \\ f & e & d \\ i & h & g \end{bmatrix}.$$

Same thing happens with scaling matrices (diagonal matrix):

$$\begin{bmatrix} 2 & & \\ & 3 & \\ & & 4 \end{bmatrix} \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} = \begin{bmatrix} 2a & 2b & 2c \\ 3d & 3e & 3f \\ 4g & 4h & 4i \end{bmatrix}, \quad \text{and} \quad \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \begin{bmatrix} 2 & & \\ & 3 & \\ & & 4 \end{bmatrix} = \begin{bmatrix} 2a & 3b & 4c \\ 2d & 3e & 4f \\ 2g & 3h & 4i \end{bmatrix}.$$

**Example 3.5.** Find the PLU decomposition of  $A = \begin{bmatrix} 2 & 1 & 1 \\ 4 & 3 & 3 \\ 8 & 7 & 9 \end{bmatrix}$ .

This is the same matrix as the coefficient matrix in Example 3.4, so we will use the work there, but specifying all the row operation matrices.

First we exchanged row 1 and row 3 since 8 is the biggest entry in first column:

$$P_1 A = \underbrace{\begin{bmatrix} & & 1 \\ & 1 & \\ 1 & & \end{bmatrix}}_{P_1} \begin{bmatrix} 2 & 1 & 1 \\ 4 & 3 & 3 \\ 8 & 7 & 9 \end{bmatrix} = \begin{bmatrix} 8 & 7 & 9 \\ 4 & 3 & 3 \\ 2 & 1 & 1 \end{bmatrix}.$$

Second we use the pivot 8 to eliminate the entries below:

$$L_1 P_1 A = \underbrace{\begin{bmatrix} 1 & & \\ -\frac{1}{2} & 1 & \\ -\frac{1}{4} & & 1 \end{bmatrix}}_{L_1} \begin{bmatrix} 8 & 7 & 9 \\ 4 & 3 & 3 \\ 2 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 8 & 7 & 9 \\ & -\frac{1}{2} & -\frac{3}{4} \\ & -\frac{3}{4} & -\frac{5}{4} \end{bmatrix} \quad \text{here } L_1 \text{ combines two row operations.}$$

Now we are at stage 2 (2nd column),  $-3/4$  is bigger than  $-1/2$  (remember we compare numbers in absolute value), so we exchange row 2 and row 3:

$$P_2 L_1 P_1 A = \underbrace{\begin{bmatrix} 1 & & \\ & & 1 \\ & 1 & \end{bmatrix}}_{P_2} \begin{bmatrix} 8 & 7 & 9 \\ & -\frac{1}{2} & -\frac{3}{4} \\ & -\frac{3}{4} & -\frac{5}{4} \end{bmatrix} = \begin{bmatrix} 8 & 7 & 9 \\ & -\frac{3}{4} & -\frac{5}{4} \\ & -\frac{1}{2} & -\frac{3}{4} \end{bmatrix}.$$

Finally we perform the last row operation to eliminate the entry below  $-3/4$ .

$$L_2 P_2 L_1 P_1 A = \underbrace{\begin{bmatrix} 1 & & \\ & 1 & \\ & \frac{2}{3} & 1 \end{bmatrix}}_{L_2} \begin{bmatrix} 8 & 7 & 9 \\ & -\frac{3}{4} & -\frac{5}{4} \\ & -\frac{1}{2} & -\frac{3}{4} \end{bmatrix} = \begin{bmatrix} 8 & 7 & 9 \\ & -\frac{3}{4} & -\frac{5}{4} \\ & & -\frac{2}{3} \end{bmatrix} = U.$$

Now we have  $L_2 P_2 L_1 P_1 A = U$ . We will rewrite the left hand side as  $L_2 L'_1 P_2 P_1 A$ , which means that  $P_2 L_1 = L'_1 P_2$ , so

$$L'_1 = P_2 L_1 P_2^{-1} = P_2 L_1 P_2^T = \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix} \begin{bmatrix} 1 & & \\ -\frac{1}{2} & 1 & \\ -\frac{1}{4} & & 1 \end{bmatrix} \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix} = \begin{bmatrix} 1 & & \\ -\frac{1}{4} & 1 & \\ -\frac{1}{2} & & 1 \end{bmatrix}.$$

It is crucial that  $L'_1$  is still lower triangular.

With  $\underbrace{L_2 L'_1}_{L^{-1}} \underbrace{P_2 P_1}_P A = U$ , we have

$$P \begin{bmatrix} 2 & 1 & 1 \\ 4 & 3 & 3 \\ 8 & 7 & 9 \end{bmatrix} = LU, \quad \text{where } P = \begin{bmatrix} & & 1 \\ 1 & & \\ & 1 & \end{bmatrix}, L = \begin{bmatrix} 1 & & \\ -\frac{1}{4} & 1 & \\ -\frac{1}{2} & \frac{2}{3} & 1 \end{bmatrix}, U = \begin{bmatrix} 8 & 7 & 9 \\ & -\frac{3}{4} & -\frac{5}{4} \\ & & -\frac{2}{3} \end{bmatrix}. \quad (3.4)$$

**Example 3.6.** Given the PLU decomposition in (3.4), solve  $\begin{bmatrix} 2 & 1 & 1 \\ 4 & 3 & 3 \\ 8 & 7 & 9 \end{bmatrix} x = \begin{bmatrix} -3 \\ -3 \\ -1 \end{bmatrix}$ .

Multiply both sides of the system by  $P$ , then we get  $LUx = \begin{bmatrix} -1 \\ -3 \\ -3 \end{bmatrix}$ . The rest is the same as

Example 3.2. We solve  $Ly = \begin{bmatrix} -1 \\ -3 \\ -3 \end{bmatrix}$  and get  $y = [-1, -11/4, -2/3]^T$ . We then solve  $Ux = y$  to get  $x = [-3, 2, 1]^T$ .

You can find more details about PLU factorization in [4, Lecture 21].

### 3.2.3 Methods and software

The direct methods are the methods of choice for most linear systems. To be more specific, Gaussian elimination with pivoting are recommended.

LAPACK (Linear Algebra Package) is a standard software library for numerical linear algebra. It provides routines for solving systems of linear equations and linear least squares, eigenvalue problems, and singular value decomposition. MATLAB started its life in the late 1970s as an interactive calculator built on top of LAPACK<sup>1</sup>. Python also uses the LAPACK routine. In Python, the command `P, L, U = scipy.linalg.lu(A)` provides the decomposition  $A = PLU$ . Many other functions of Python, like computing the determinant or finding the inverse, are based on the PLU factorization.

It can be very difficult to solve a large dense linear system. By difficult, we mean it takes too long using Gaussian elimination. But if a large system is sparse (most entries are 0), it can be solved efficiently using iterative methods. Systems of this type arise naturally, for example, when finite-difference techniques are used to solve boundary-value problems, a common application in the numerical solution of partial differential equations.

## Chapter 3 Exercises

- Write out the full matrix  $A = [a_{ij}]_{3 \times 4}$  where  $a_{ij} = i^{j-1}$ .
- Let  $A = \begin{bmatrix} 3 & 0 \\ -1 & 2 \\ 1 & 1 \end{bmatrix}$ ,  $B = \begin{bmatrix} 1 & 5 & 2 \\ -1 & 0 & 1 \\ 3 & 2 & 4 \end{bmatrix}$ ,  $C = \begin{bmatrix} 1 & 4 & 2 \\ 3 & 1 & 5 \end{bmatrix}$ . Compute the following (write undefined if operation is not legal).
  - $2A^T + C$
  - $CBB$
  - $BA$
- Compute  $x^T Ax$ , where  $x = (x_1, x_2, x_3)$ , and  $A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{12} & a_{22} & a_{23} \\ a_{13} & a_{23} & a_{33} \end{bmatrix}$  ( $A$  is symmetric.)
- Perform forward elimination.
  - $\left[ \begin{array}{ccc|c} 2 & 2 & 2 & 0 \\ -2 & 5 & 2 & 1 \\ 8 & 1 & 4 & -1 \end{array} \right]$
  - $\left[ \begin{array}{ccc|c} 0 & -2 & 3 & 1 \\ 3 & 6 & -3 & -2 \\ 6 & 6 & 3 & 5 \end{array} \right]$

<sup>1</sup>It was called LINPACK and EISPACK then, the predecessor of LAPACK.

5. Find scalars  $c_i$  for which the equation  $c_1(-1, 0, 2) + c_2(2, 2, -2) + c_3(1, -2, 1) = (-6, 12, 4)$  holds.
6. How many **exact** flops are needed for the following computations?
  - (a)  $x$  dot product  $y$ , where  $x, y$  are vectors in  $\mathbb{R}^9$ .
  - (b)  $A + B$ , where  $A, B$  are both  $m \times n$ .
  - (c)  $ABC$ , where  $A, B, C$  are  $5 \times 5$  matrices.
  - (d)  $Ux$ , where  $U$  is  $n \times n$  upper triangular matrix and  $x$  is a vector in  $\mathbb{R}^n$ .
  - (e)  $ABC$ , where  $A$  is  $10 \times 100$ ,  $B$  is  $100 \times 100$ ,  $C$  is  $100 \times 100$ .
  - \*(f)  $LU$ , where  $L$  is  $n \times n$  lower triangular,  $U$  is  $n \times n$  upper triangular.
7. Find the  $4 \times 4$  elementary row matrices that correspond to
  - (a)  $\rho 2(3) + \rho 3$
  - (b)  $\rho 1(-3) + \rho 4$
8. Find the inverse of the matrix in 7(a). What row operation does this inverse correspond to?
9. Find the inverse of  $\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$ .
10. Is  $\begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 1 \\ 2 & 3 & 2 \end{bmatrix}$  invertible?
11. Find the LU decomposition of the coefficient matrix in 4(a).
12. Does the system in 4(b) solvable? Is the matrix  $A = \begin{bmatrix} 0 & -2 & 3 \\ 3 & 6 & -3 \\ 6 & 6 & 3 \end{bmatrix}$  invertible?
13. Given  $\begin{bmatrix} 3 & -6 & -3 \\ 2 & 0 & 6 \\ -4 & 7 & 4 \end{bmatrix} = \begin{bmatrix} 3 & 0 & 0 \\ 2 & 4 & 0 \\ -4 & -1 & 2 \end{bmatrix} \begin{bmatrix} 1 & -2 & -1 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{bmatrix}$ , solve the system
 
$$\begin{array}{rrcr} 3x_1 & -6x_2 & -3x_3 & = -3 \\ 2x_1 & & +6x_3 & = -22 \\ -4x_1 & +7x_2 & +4x_3 & = 3 \end{array}$$
14. Solve  $\begin{bmatrix} 2 & 6 & 2 \\ -3 & -8 & 0 \\ 4 & 9 & 2 \end{bmatrix} x = \begin{bmatrix} 2 \\ 2 \\ 3 \end{bmatrix}$  ( $Ax = b$ ) by
  - (a) First find the  $LU$  factorization of  $A$ .
  - (b) Second solve  $Ly = b$ .
  - (c) Third solve  $Ux = y$ .
15. Consider the linear system  $\left[ \begin{array}{cc|c} 1 & 1 & 2 \\ 2 & 2 \times 10^{20} & 2 \times 10^{20} \end{array} \right]$ . Find its solution using partial pivoting, in a machine where numbers are in standard IEEE double-precision format. Do you think the solution you find is a fairly accurate approximation to the real solution?

16. The LINPACK Benchmarks are a measure of a system's floating point computing power. Introduced by Jack Dongarra, they measure how fast a computer solves a dense  $n \times n$  system of linear equations  $Ax = b$ . This is the standard measure for computer speed. Supercomputers (which usually have linux operating systems) can do about 10 petaflops in one second. One petaflop is  $10^{15}$  flops. Titan, a supercomputer in the Oak Ridge National lab, is the fastest computer in the US. It can perform 17.59 petaflops every second. We assume that it takes  $\frac{2}{3}n^3$  flops for Gauss Elimination.

- (a) Calculate how much time (in seconds) does Titan need to solve a 50000 by 50000 system  $Ax = b$ .
- (b) Suppose a laptop performs 20 gigaflops (a gigaflop is  $10^9$  flops) per second. Calculate how much time (in seconds) does this laptop need to solve a 50000 by 50000 system  $Ax = b$ .

17. Let  $V = \begin{bmatrix} 0 & -2 & 4 & 2 \\ 0 & 0 & -2 & 7 \\ 2 & 1 & 3 & 1 \\ 0 & 0 & 0 & -43/2 \end{bmatrix}$ . Find the permutation matrix  $P$  so that  $PV$  is upper triangular.

18. State the matrix that you need to multiply for the following transformations. Specify left or right multiplication.

(a)  $\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \rightarrow \begin{bmatrix} 100 & 200 \\ -30 & -40 \end{bmatrix}$  (b)  $\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \rightarrow \begin{bmatrix} 100 & -20 \\ 300 & -40 \end{bmatrix}$  (c)  $\begin{bmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \end{bmatrix} \rightarrow \begin{bmatrix} 2 & 1 & 3 \\ 2 & 1 & 3 \\ 2 & 1 & 3 \end{bmatrix}$

19. Solve the system  $\left[ \begin{array}{ccc|c} 1 & 2 & 3 & 1 \\ 4 & 5 & 6 & 0 \\ 7 & 8 & 0 & 2 \end{array} \right]$  using partial pivoting.

20. Given the following row operations

$$A = \begin{bmatrix} 1 & \frac{5}{2} & \frac{5}{2} \\ 2 & 3 & 1 \\ 2 & 5 & \frac{17}{6} \\ -3 & \frac{4}{3} & \frac{6}{5} \end{bmatrix} \xrightarrow{\rho_1 \leftrightarrow \rho_2} \begin{bmatrix} 2 & 3 & 1 \\ 1 & \frac{5}{2} & \frac{5}{2} \\ 2 & 5 & \frac{17}{6} \\ -3 & \frac{4}{3} & \frac{6}{5} \end{bmatrix} \xrightarrow[\rho_1(-1/3)+\rho_3]{\rho_1(-1/2)+\rho_2} \begin{bmatrix} 2 & 3 & 1 \\ 0 & 1 & 2 \\ 0 & \frac{1}{4} & \frac{5}{2} \\ 0 & \frac{1}{4} & 2 \end{bmatrix} \xrightarrow{\rho_2(-1/4)+\rho_3} \begin{bmatrix} 2 & 3 & 1 \\ 0 & 1 & 2 \\ 0 & 0 & 2 \\ 0 & 0 & 2 \end{bmatrix}$$

- (a) List the three matrices  $(P_1, L_1, L_2)$  that correspond to these three row operations.
- (b) Find the PLU factorization of  $A$ . You need to list  $P, L, U$  explicitly.

21. Given the following row operations

$$A = \begin{bmatrix} 2 & 1 & 1 \\ 4 & -6 & 0 \\ -2 & 8 & 2 \end{bmatrix} \xrightarrow{\rho_1 \leftrightarrow \rho_2} \begin{bmatrix} 4 & -6 & 0 \\ 2 & 1 & 1 \\ -2 & 8 & 2 \end{bmatrix} \xrightarrow[\rho_1 \frac{1}{2} + \rho_3]{\rho_1(-\frac{1}{2}) + \rho_2} \begin{bmatrix} 4 & -6 & 0 \\ 0 & 4 & 1 \\ 0 & 5 & 2 \end{bmatrix} \xrightarrow{\rho_2 \leftrightarrow \rho_3} \begin{bmatrix} 4 & -6 & 0 \\ 0 & 5 & 2 \\ 0 & 4 & 1 \end{bmatrix} \xrightarrow{\rho_2(-\frac{4}{5}) + \rho_3} \begin{bmatrix} 4 & -6 & 0 \\ 0 & 5 & 2 \\ 0 & 0 & -0.6 \end{bmatrix}$$

- (a) List the four matrices  $(P_1, L_1, P_2, L_2)$  that correspond to these four row operations.
- (b) Find the PLU factorization of  $A$ . You need to list  $P, L, U$  explicitly.

22. Given the PLU factorization of  $A$  as  $PA = LU$ , where

$$P = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, L = \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ \frac{1}{3} & \frac{1}{4} & 1 \end{bmatrix}, U = \begin{bmatrix} 2 & 3 & 1 \\ 0 & 1 & 2 \\ 0 & 0 & 2 \end{bmatrix}.$$

Using these factors, without forming the matrix  $A$  or taking any inverses, to solve the system  $Ax = b$ , where  $b = (2, 10, -12)$ .

### Python Exercises

23. Basic matrix operations:

- (a) Create a random  $4 \times 4$  matrix  $A$  (each entry follows standard normal distribution). Print  $A$ .
- (b) Display the 3rd column of  $A$ .
- (c) Display the upper right  $2 \times 2$  submatrix of  $A$ .
- (d) Compute the sum of all entries of  $A$ .
- (e) Define  $B$  to be the transpose of  $A$ .

24. How powerful is your computer?

- (a) Create a 5000 by 5000 random matrix  $R$ , where each entry of  $R$  is standard normal distribution. Define  $r$  to be the product of  $R$  and a  $5000 \times 1$  vector of all 1's.
- (b) Measure how much time does it take to solve  $Rx=r$  using Python's built-in function. You may want to run several times so that the time stabilizes. Use `time1` to store the elapsed time. Print it. Verify that the solution is the vector of all 1's.
- (c) Calculate how many gigaflops your computer can do in a second (see Problem 16.) Print your answer.

25. It is never necessary to compute the inverse of a matrix in practice. For example, some people may want to solve a linear system  $Ax = b$  by finding the inverse of  $A$  and then  $x = A^{-1}b$ . Let's try this with the system  $Rx=r$ , where  $R$  and  $r$  are from Problem 24. To be precise, you are going to measure how many seconds it takes to "first finding the inverse of  $R$ , and then performing the multiplication  $R^{-1}r$ ". Use `time2` to store the elapsed time. Both Problem 24(b) and this problem is solving the same linear system. Which method takes more time?

26. Write a function `x = backsub(U, b)` to do backward substitution from scratch. The input  $U$  must be an invertible upper triangular matrix (of any size) and  $b$  is a column vector of appropriate length. The function returns the solution of  $Ux=b$ . You only need one loop in your code. Test your function with the backward substitution step in Example 3.1.

27. More time comparison.

- (a) Define  $U$  to be the upper triangular portion of a  $5000 \times 5000$  matrix full of 1's. Define  $b$  to be the product of  $U$  and a  $5000 \times 1$  vector of all 1's.
- (b) Solve  $Ux=b$  using `scipy.linalg.solve_triangular`. Use `time3` to store the elapsed time.

- (c) Solve  $Ux=b$  using your own function `backsub`. Use `time4` to store the elapsed time. Check that the solution is the same as the one from (b) to make sure your code is correct.

- (d) Finally run the following command:

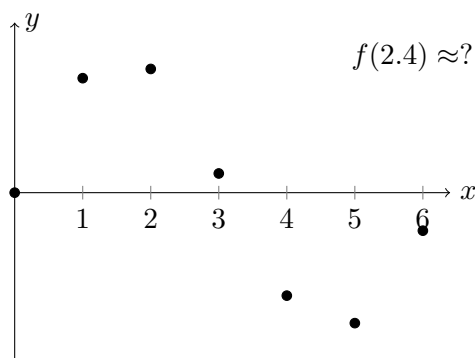
```
print ("For a 5000 by 5000 dense system, it took",time1,"seconds to solve via the  
Gaussian elimination. It took",time2,"seconds to solve via finding the inverse.")  
print ("For a 5000 by 5000 upper triangular system, it took",time3,"seconds to  
solve via using Python's built-in function. It took",time4,"seconds to solve via my code.")
```

Comment on the differences among these time span.

## Chapter 4

# Polynomial Interpolation

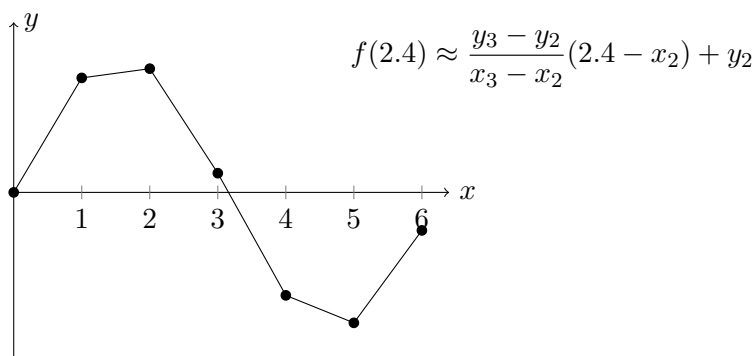
In many applications of engineering and science, one often has a number of data points, obtained by sampling or experimentation, say  $f(x_0), f(x_1), \dots, f(x_n)$ , and would like to estimate the value of this function at an intermediate value of the  $x_i$ 's. See the following figure for example.



This is the problem of interpolation. An estimation of the value at another point is often given by finding a function that at least agrees with  $f$  at these  $n+1$  points. In general, given  $n+1$  points in  $\mathbb{R}^2$ :  $\{(x_i, y_i)\}_{i=0}^n$ , any function that goes through these  $n+1$  points is called an *interpolant* of these points. This is a classical topic from approximation theory, and we will only present some important theories here.

### 4.1 Linear Interpolation

The simplest way is linear interpolation, that is, connect nearby points by lines which makes the interpolant a piecewise linear function. See below:





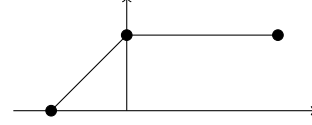
This is quick and easy, but not precise and the whole function is not differentiable at the nodes  $x_i$ . But linear interpolation is the simplest example of spline, about which we will talk later.

**Example 4.1.** Find the linear interpolant of the three points  $(-1, 0)$ ,  $(0, 1)$ ,  $(2, 1)$ , and then estimate  $f(-0.8)$ .

Solution: Let  $L(x)$  be the linear interpolant and it is a 2-piece line, that is

$$L(x) = \begin{cases} x + 1, & x \in [-1, 0] \\ 1, & x \in [0, 2] \end{cases}.$$

So  $f(-0.8) \approx L(-0.8) = 0.2$ .



## 4.2 Interpolation by One Polynomial

Computers can do additions and multiplications really well, therefore making polynomials a friendly function to a computer. For this reason, we would like to interpolate a given set of data points with polynomials.

We can find a line (degree 1 polynomial) going through two points, and we can find a quadratic function (degree 2 polynomial) going through 3 given points. In general, it is believable that we are able to find a polynomial of degree  $n$  to fit  $n + 1$  points because such polynomial has  $n + 1$  coefficients ( $n + 1$  degrees of freedom) and we have exactly  $n + 1$  pieces of information. We will often use  $p_n$  for a polynomial of degree  $n$ .

**Theorem 4.2.** Given  $\{(x_i, y_i)\}_{i=0}^n$  where the  $n + 1$  values of  $x_i$  are all different, then there is a **unique** polynomial  $p$  of degree at most  $n$  that goes through these  $n + 1$  points.

The proof is easy. Subsection 4.2.1 can be viewed as a proof.

We will introduce three well-known methods for finding this unique  $p$ .

### 4.2.1 Vandermonde matrix

Let  $p(x) = \sum_{i=0}^n c_i x^i$ . In order to figure out the coefficients  $c_i$ , we plug in the known points  $(x_j, y_j)$ ,

which generates  $n + 1$  equations  $\sum_{i=0}^n c_i x_j^i = y_j, j = 0, 1, \dots, n$ . These equations are all linear, so we get an  $(n + 1) \times (n + 1)$  system

$$\begin{bmatrix} 1 & x_0 & x_0^2 & \cdots & x_0^n \\ 1 & x_1 & x_1^2 & \cdots & x_1^n \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^n \end{bmatrix} \begin{bmatrix} c_0 \\ c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix} = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{bmatrix}. \quad (4.1)$$

This is a straightforward method. The coefficient matrix in the system (4.1) is in the form of a Vandermonde matrix.

In general, a  $k \times l$  *Vandermonde matrix* is a matrix generated by a vector  $v = (v_1, v_2, \dots, v_k)$  in the following fashion:

For  $1 \leq i \leq l$ , the  $i$ th column is  $[v_1^{i-1}, v_2^{i-1}, \dots, v_k^{i-1}]^T$ .

For example, the matrix in Chapter 3 Exercise 1 is a  $3 \times 4$  Vandermonde matrix generated by the vector  $(1, 2, 3)$ .

It can be proven that as long as  $x_i \neq x_j$  for any  $i \neq j$ , then the Vandermonde matrix in (4.1) is invertible, thus (4.1) has a unique solution. System (4.1) will then be solved via Gaussian elimination as described in Chapter 3. In the solution,  $c_n$  (or any other coefficient) could be 0, therefore the polynomial we get is of degree at most  $n$ .

**Example 4.3.** Find the quadratic that goes through the three points  $(-1, 0)$ ,  $(0, 1)$ ,  $(2, 1)$  (with partial pivoting). Then estimate  $f(-0.8)$ .

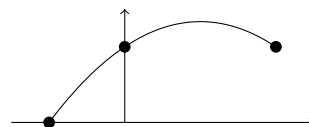
$$\left[ \begin{array}{ccc|c} 1 & -1 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 2 & 4 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & -1 & 1 & 0 \\ 0 & 1 & -1 & 1 \\ 0 & 3 & 3 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & -1 & 1 & 0 \\ 0 & 3 & 3 & 1 \\ 0 & 1 & -1 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & -1 & 1 & 0 \\ 0 & 3 & 3 & 1 \\ 0 & 0 & -2 & 2/3 \end{array} \right]$$

$$c_2 = -1/3; \quad 3c_1 - 1 = 1 \Rightarrow c_1 = 2/3; \quad c_0 - 2/3 - 1/3 = 0 \Rightarrow c_0 = 1.$$

$$p(x) = 1 + \frac{2}{3}x - \frac{1}{3}x^2.$$

$$\text{So } f(-0.8) \approx p(-0.8) = 0.2533.$$

This can be compared to Example 4.1.



## 4.2.2 Lagrange interpolation

Lagrange solved this problem in a different approach, where no linear system is involved. First we generate  $n + 1$  polynomials  $\varphi_i(x)$  such that

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

$\delta_{ij}$  is a useful notation and is heavily used in mathematics and engineering.

The construction of  $\varphi_i$  is immediate. For example, if  $x_0 = 2, x_1 = 2.5, x_2 = 4$ , then  $\varphi_0$  has roots 2.5 and 4, so  $\varphi_0(x) = c(x - 2.5)(x - 4)$ .  $c$  is determined by  $\varphi_0(x_0) = 1$ , so we can write  $\varphi_0(x) = \frac{(x - 2.5)(x - 4)}{(2 - 2.5)(2 - 4)}$ . Similarly,  $\varphi_1(x) = \frac{(x - 2)(x - 4)}{(2.5 - 2)(2.5 - 4)}$  and  $\varphi_2(x) = \frac{(x - 2)(x - 2.5)}{(4 - 2)(4 - 2.5)}$ .

Note that  $\varphi_i$ 's only depend on the nodes  $\{x_0, x_1, \dots, x_n\}$ . In general, they can be written as

$$\varphi_i(x) = \frac{\prod_{j \neq i} (x - x_j)}{\prod_{j \neq i} (x_i - x_j)}. \quad (4.2)$$

The second step of the Lagrange method is even easier, the answer that we are looking for is then

$$p(x) = \sum_{i=0}^n y_i \varphi_i(x) \quad (4.3)$$

**Example 4.4.** Redo Example 4.3 using Lagrange's method.

$$x_0 = -1, x_1 = 0, x_2 = 2.$$

Form  $\varphi_0(x) = \frac{(x - 0)(x - 2)}{(-1 - 0)(-1 - 2)}$ ,  $\varphi_1(x) = \frac{(x + 1)(x - 2)}{(0 + 1)(0 - 2)}$ ,  $\varphi_2(x) = \frac{(x + 1)(x - 0)}{(2 + 1)(2 - 0)}$ . After simplification, we have  $\varphi_0(x) = \frac{1}{3}x(x - 2)$ ,  $\varphi_1(x) = \frac{-1}{2}(x + 1)(x - 2)$ ,  $\varphi_2(x) = \frac{1}{6}x(x + 1)$ .

So  $p(x) = 0 \cdot \varphi_0(x) + 1 \cdot \varphi_1(x) + 1 \cdot \varphi_2(x) = \frac{-1}{2}(x + 1)(x - 2) + \frac{1}{6}x(x + 1)$ . One can verify that this is the same polynomial as the one from Example 4.3.

### 4.2.3 Newton's interpolation

The Newton's interpolation is similar to the first method. We only formulate  $p(x)$  differently as

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

Plugging in  $y_0 = p(x_0)$  implies  $y_0 = a_0$ . Plugging in  $y_1 = p(x_1)$  implies  $y_1 = a_0 + a_1(x_1 - x_0)$ . In the end,  $p(x_i) = y_i$  gives us a lower triangular system

$$\begin{bmatrix} 1 & 0 & 0 & 0 & \cdots & 0 \\ 1 & x_1 - x_0 & 0 & 0 & \cdots & 0 \\ 1 & x_2 - x_0 & (x_2 - x_0)(x_2 - x_1) & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & & \vdots \\ 1 & x_n - x_0 & (x_n - x_0)(x_n - x_1) & \cdots & \prod_{i=0}^{n-1} (x_n - x_i) & \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}. \quad (4.4)$$

It is very efficient to solve (4.4) because it is a lower triangular system.

**Example 4.5.** Redo Example 4.3 using Newton's interpolation.

We let  $p(x) = a_0 + a_1(x + 1) + a_2(x + 1)(x - 0)$ . Once again, in this method, The formulation of  $p$  depends on all, but the last nodes. We still have  $p(x_i) = y_i$  for  $i = 0, 1, 2$ , which produces three equations:

$$\begin{aligned} a_0 &= 0 \\ a_0 + a_1(0 + 1) &= 1 \\ a_0 + a_1(2 + 1) + a_2(2 + 1)(2 - 0) &= 1 \end{aligned} \quad (4.5)$$

This is the same as  $\begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 + 1 & 0 \\ 1 & 2 + 1 & (2 + 1)(2 - 0) \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$ , which matches (4.4). Solving this lower triangular system, we get  $a_0 = 0, a_1 = 1, a_2 = -1/3$ .  
So  $p(x) = (x + 1) - \frac{1}{3}(x + 1)(x - 0)$ .

One benefit of using Newton's formulation is that the computation is minimal for adding a data point. Suppose now we are trying to find the degree 3 polynomial going through  $(-1, 0), (0, 1), (1.5, 3), (2, 1)$ . We observe that  $(1.5, 3)$  is the extra point, so we will reorder the four points as  $(-1, 0), (0, 1), (2, 1), (1.5, 3)$ . [The ordering of the points has no effect on the polynomial interpolant. These points are considered collectively.]

Letting  $q(x) = a_0 + a_1(x + 1) + a_2(x + 1)(x - 0) + a_3(x + 1)(x - 0)(x - 2)$ . Plugging these four points, we get 4 equations:

$$\begin{aligned} a_0 &= 0 \\ a_0 + a_1(0 + 1) &= 1 \\ a_0 + a_1(2 + 1) + a_2(2 + 1)(2 - 0) &= 1 \\ a_2 + a_2(1.5 + 1) + a_2(1.5 + 1)(1.5 - 0) + a_3(1.5 + 1)(1.5 - 0)(1.5 - 2) &= 3 \end{aligned} \quad (4.6)$$

Note that the first three equations are exactly the same as (4.5)! This means we can reuse the answers  $a_0, a_1, a_2$  solved earlier. We only need to solve for  $a_3$  as the additional work.

**Among all three methods presented, Newton's interpolation is the only one that is easy to add a node.**

#### 4.2.4 Flop count

	Vandermonde	Largange	Newton
Compute coefficients	$O(n^3)$	0	$O(n^2)$
Evaluate at one point	$O(n)$	$O(n^2)$	$O(n)$

Table 4.1: Flop count for different methods of polynomial interpolation

For the Vandermonde method, it takes  $O(n^3)$  flops to compute the coefficients, since this is how much you need for Gaussian elimination. To evaluate  $p(x) = \sum_{i=0}^n c_i x^i$  at a point, we mentioned at the end of Chapter 1 that we need to use the nested form, which takes  $2n$  flops.

For Lagrange interpolation, we simply write down the interpolation polynomial (4.3) so it takes no work to generate the polynomial. Evaluation at one point will take  $O(n^2)$  flops.

For Newton's interpolation, it takes  $O(n^2)$  flops to compute the  $a_i$ 's since we have a lower triangular system. Evaluation at one point only takes  $O(n)$  as well with the nested form. For example,  $a_0 + a_1(x - x_0) + a_2(x - x_0)(x - x_1) + a_3(x - x_0)(x - x_1)(x - x_2)$  should be rewritten as  $a_0 + (x - x_0)[a_1 + (x - x_1)(a_2 + a_3(x - x_2))]$ .

#### 4.2.5 Limitation of polynomial interpolant

Consider the function  $R(x) = \frac{1}{1 + 25x^2}$  on the interval  $I = [-1, 1]$ . This is called the Runge function (scaled). Suppose we are interpolating 11 equally spaced points that we get from sampling  $R(x)$  (Figure 4.1 (a)). Figure 4.1 (b) is the graph of the degree 10 polynomial that goes through these 11 points. Notice that there are big oscillations near the ends of the interval. In fact, as degree increases, the oscillations near the end will get even larger.

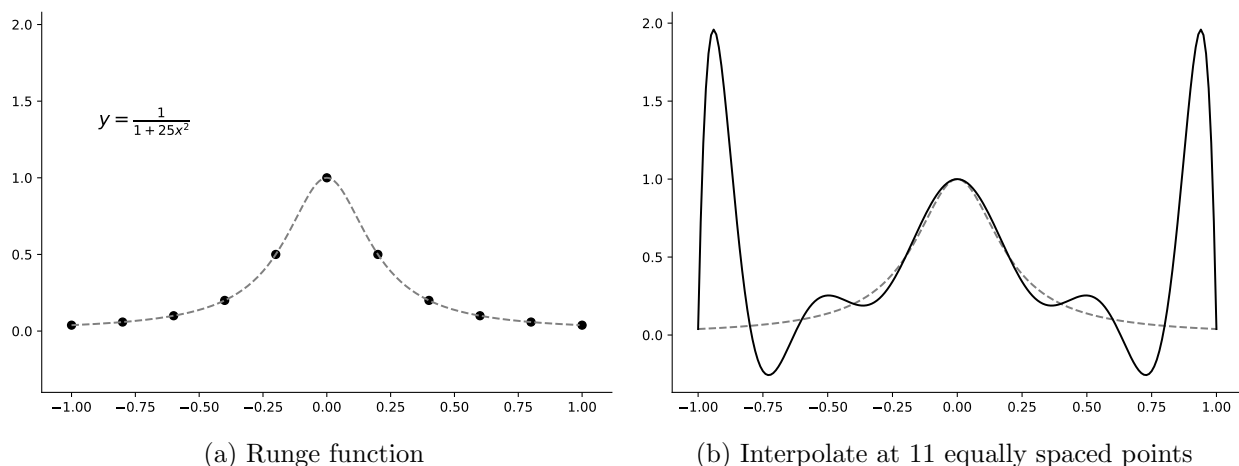


Figure 4.1: Runge function effect

In the next two sections, we introduce two approaches to reduce this Runge function effect. The first approach is to pick a better set of nodes. This is Section 4.3. The second approach, discussed in Section 4.4, is to use multiple lower degree polynomials.

### 4.3 Chebyshev nodes

The polynomial interpolant depends on the nodes  $\{x_i\}_{i=0}^n$ . In the Runge function example, we used equally spaced points, which is a common choice. However, we will show that there are better nodes in the sense that the corresponding polynomial interpolant is a better approximate. For this we will have to first find a way to measure distances between functions.

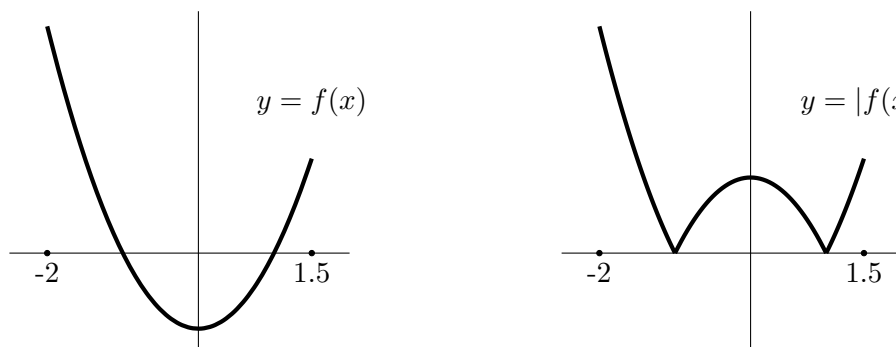
Given  $f$ , who is continuous on  $[a,b]$ , we define its  $\infty$ -norm by

$$\|f\|_{\infty} := \max_{x \in [a,b]} |f(x)|,$$

and its 2-norm by

$$\|f\|_2 := \sqrt{\int_a^b |f(x)|^2 dx}.$$

**Example 4.6.** Let  $f(x) = x^2 - 1$  on  $[-2, 1.5]$ . Since this function is easy to graph, we will figure out its norm by first graphing it. Note that to draw  $|f(x)|$ , we flip the negative portion up.



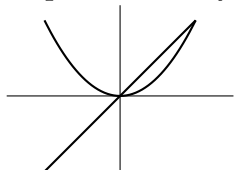
It is obvious that the maximum value of  $|f(x)|$  is obtained at  $x = -2$ , so  $\|f\|_{\infty} = |f(-2)| = 3$ .

$$\begin{aligned} \|f\|_2^2 &:= \int_{-2}^{1.5} (x^2 - 1)^2 dx = \int_{-2}^{1.5} x^4 - 2x^2 + 1 dx = \frac{1.5^5 + 2^5}{5} - 2 \frac{1.5^3 + 2^3}{3} + 3.5 \approx 3.835417. \\ \|f\|_2 &= \sqrt{3.835417} \approx 1.96. \end{aligned}$$

Correspondingly, we can define two different distances between two continuous functions  $f, g$  on  $[a, b]$ :

$$\begin{aligned} \infty\text{-distance: } \|f - g\|_{\infty} &= \max_{x \in [a,b]} |f(x) - g(x)| \\ 2\text{-distance: } \|f - g\|_2 &= \sqrt{\int_a^b |f(x) - g(x)|^2 dx} \end{aligned}$$

**Example 4.7.** Let  $f(x) = x^2, g(x) = x$  on  $[-1, 1]$ .



We can see from the graph on the left that  $f$  and  $g$  differ the most at  $x = -1$ , so  $\|f - g\|_{\infty} = |f(-1) - g(-1)| = 2$ .

$$\begin{aligned} \int_{-1}^1 (f - g)^2 dx &= \int_{-1}^1 (x^2 - x)^2 dx = \int_{-1}^1 x^4 + x^2 - 2x^3 dx = 16/15. \\ \|f - g\|_2 &= \sqrt{16/15}. \end{aligned}$$

The  $\infty$  distance between the Runge function  $R(x)$  and its equally spaced interpolant  $P_{10}(x)$  (as shown in Figure 4.1(b)) is  $\|R - P_{10}\|_{\infty} = 1.9157$ . (The maximum occurs at  $x = \pm 0.94$ .)

In general, it is not easy to find the maximum of a function. As learned in Calculus I, a common method is to find the critical points by taking its derivative, and then compare all values at the critical points and end points. This is why we also introduce the 2-distance, which is a lot easier to compute than the  $\infty$ -distance.

The following theorem illustrates the point-wise error of polynomial interpolation.

**Theorem 4.8.** Assume that  $f \in C^{n+1}[a, b]$  ( $n + 1$ st derivative is continuous) and that  $x_0, \dots, x_n$  are in  $[a, b]$ . Let  $p(x)$  be the polynomial of degree  $n$  that interpolates  $f$  at  $\{x_i\}_{i=0}^n$ . Then for any point  $x \in [a, b]$ ,

$$f(x) - p(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \prod_{i=0}^n (x - x_i) \quad (4.7)$$

for some point  $\xi \in [a, b]$ .

The proof only needs Rolle's theorem, see Theorem 8.4.1 of [2] for details.

Let's suppose the interval is  $[-1, 1]$  for now. This gives us an estimate on the  $L_{\infty}$  distance between  $f$  and  $p$  on the interval  $[-1, 1]$ :

$$\|f - p\|_{\infty} = \max_{x \in [-1, 1]} \frac{|f^{(n+1)}(\xi)|}{(n+1)!} \prod_{i=0}^n |x - x_i| \leq \max_{x \in [-1, 1]} \frac{\|f^{(n+1)}\|_{\infty}}{(n+1)!} \max_{x \in [-1, 1]} \prod_{i=0}^n |x - x_i|.$$

Suppose  $f$  is given (like the Runge function), in order to minimize the error bound, we need to solve

$$\min_{\{x_0, \dots, x_n\} \subset [-1, 1]} \max_{x \in [-1, 1]} \prod_{i=0}^n |x - x_i|. \quad (4.8)$$

The solution (not trivial) for this minimization problem is

$$\text{Chebyshev nodes: } x_i = \cos \frac{\pi i}{n}, i = 0, 1, \dots, n \quad (4.9)$$

which are the  $x$  coordinates of  $n$  equally spaced angles on the positive half unit circle, see Figure 4.2(a). These are the roots of the Chebyshev polynomial of the first kind. These nodes are not equally spaced, and are more clustered around the ends.

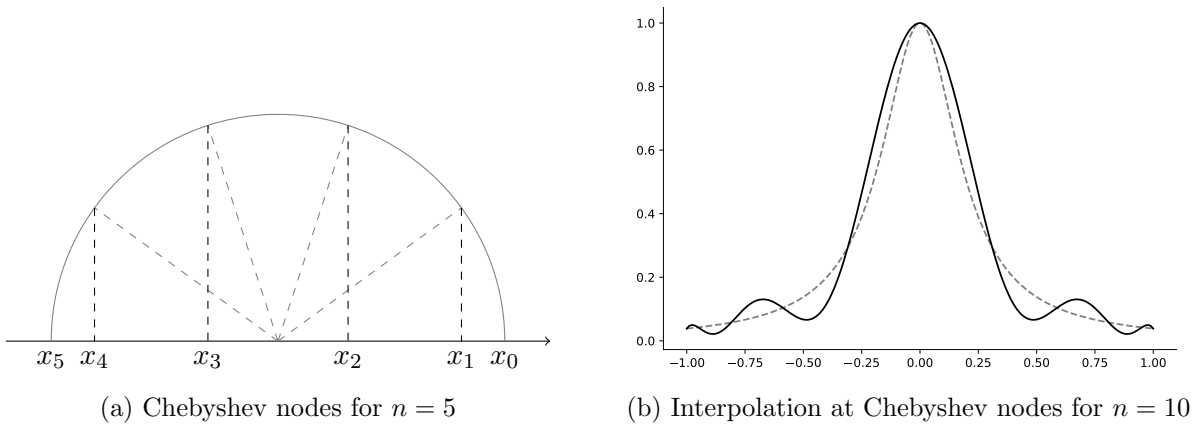


Figure 4.2: Interpolation with Chebyshev nodes

Figure 4.2 (b) is the Chebyshev interpolant, which approximates the Runge function much better. We call this polynomial  $C_{10}(x)$ . The  $\infty$  distance  $\|R - C_{10}\|_\infty = 0.1322$ , is only about 1/15 of  $\|R - P_{10}\|_\infty$ , when equally spaced points were used. Therefore if there is a choice on where to interpolate, picking the Chebyshev nodes is optimal.

The definition  $x_i = \cos \frac{\pi i}{n}$  only works for the interval  $[-1, 1]$ . What happens if we have  $[-2, 2]$ , or in general  $[a, b]$ ? This is left as an exercise.

## 4.4 Spline

Spline<sup>1</sup> means piecewise polynomials, meaning it is a piecewisely defined function and on each piece it is a different polynomial. Moreover, we generally require it to be continuous. This simply means that different pieces connect at the same point. There is no jump or drop. Linear interpolant from Section 4.1 is an instance of spline. As a spline, each piece does not have to be degree 1 polynomial, and their degrees can be different. The following function is another example of spline.

$$s(x) = \begin{cases} -x^2 + x + 1, & 0 \leq x \leq 1 \\ x^3, & 1 \leq x \leq 2 \\ 2x + 4, & 2 \leq x \leq 3 \end{cases}$$

The purpose of spline is to interpolate the same number of data points, with lower degrees, at the cost of multiple polynomials. Lower degree polynomials are not only faster to be evaluated by computers, they also tend to have less oscillations, which potentially reduces the Runge function effect. Figure 4.3(a) is to interpolate 4 points with one degree 3 polynomial (cubic), and Figure 4.3(b) is to interpolate the first two points with a line (degree 1) and the last three with a quadratic (degree 2).

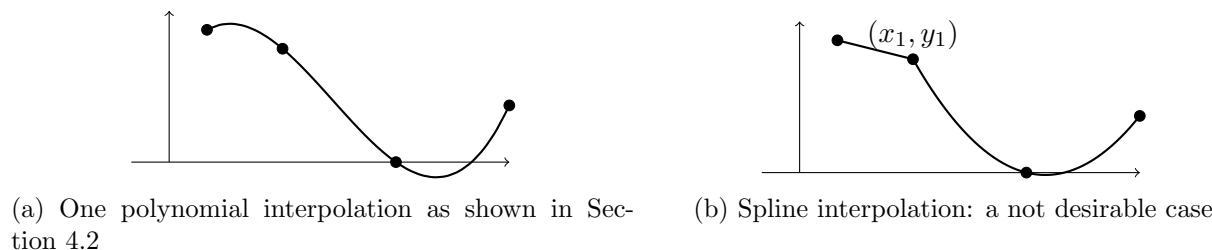


Figure 4.3

However, Figure 4.3 (b) is a situation that we want to avoid. It is not smooth at the second point  $x_1$ . To be precise, this piecewise function is not differentiable at  $x_1$  because the left derivative, which is the slope of the line, does not coincide with the right derivative.

Because spline is multiple polynomials piecing together, we need to make sure they are “well-connected” at  $x_i$ ’s. In real applications, we often require it to be twice differentiable. Below is an example of a (one time) differentiable spline.

**Example 4.9.** Given three points  $(0, 1)$ ,  $(1, -1)$ ,  $(2, 0)$ , find an interpolant

$$s(x) = \begin{cases} P_1(x), & 0 \leq x \leq 1 \\ P_2(x), & 1 \leq x \leq 2 \end{cases}$$

<sup>1</sup>The root of the word “spline” is the same as that of splint. It was originally a small strip of wood that could be used to join two boards. Later the word was used to refer to a long flexible strip, generally of metal, that could be used to draw continuous smooth curves by forcing the strip to pass through specified points and tracing along the curve.

such that  $s(x)$  is differentiable.  $P_1$  is a polynomial of degree 1, and  $P_2$  is a polynomial of degree 2.

This boils down to

$$(1) P_1(0) = 1, P_1(1) = -1; P_2(1) = -1, P_2(2) = 0$$

$$(2) P_1'(1) = P_2'(1).$$

Condition (1) already completely determines  $P_1$ :  $P_1(x) = -2x + 1$ . Therefore by (2) we know  $P_2'(1) = P_1'(1) = -2$ . Since  $P_2$  is of degree 2, we can let  $P_2'(x) = 2a(x - 1) - 2$ , which means that  $P_2(x) = a(x - 1)^2 - 2(x - 1) + b$ . Note that  $b = -1$  can be determined right away by  $P_2(1) = -1$ .

All that is left is to determine  $a$ . Plugging  $P_2(2) = 0$  yields  $a - 2 - 1 = 0$ . So  $a = 3$ .

In the end,  $P_1(x) = -2x + 1$ ,  $P_2(x) = 3(x - 1)^2 - 2(x - 1) - 1$ . It is graphed in Figure 4.4.

In order to simply fulfill the continuity requirement (1), there could be many other choices of  $P_2$  (like the gray dashed lines in Figure 4.4), but only the black one connects with the line smoothest.

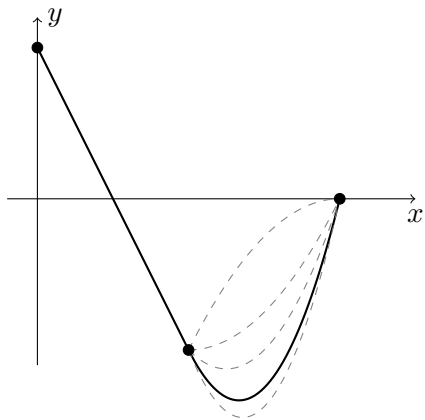


Figure 4.4: Spline interpolation: a somewhat desirable case

If a function is differentiable in higher order, then it is considered to be “smoother”. Let us consider another piecewisely defined function  $f(x) = q_2(x) := -x^2 + x$  on  $[0,1]$  and  $f(x) = q_3(x) := (x - 1)(x^2 - 3x + 1)$  on  $[1,2]$ , see Figure 4.5. The quadratic  $q_2$  and the cubic  $q_3$  share the 0th, 1st, and 2nd derivative at  $x = 1$  as  $q_2(1) = q_3(1) = 0$ ,  $q_2'(1) = q_3'(1) = -1$ ,  $q_2''(1) = q_3''(1) = -2$ . In fact, there are infinitely many other cubics will do this without specifying the value at  $x = 2$  (left as an exercise). As a result,  $f(x)$  is twice differentiable, and it transits even smoother at  $x = 1$  than  $s(x)$  in Example 4.9.

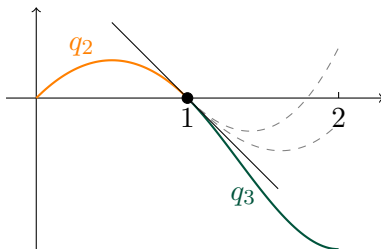


Figure 4.5: Twice differentiable spline: a desirable case

To show the difference between smoothness of once differentiable and twice differentiable functions, we draw in Figure 4.5 some other cubic pieces that will have the same first derivative as  $q_2$  at  $x = 1$ , but fails to have the second derivative equal (in gray dashed lines). Notice that the second derivative reveals the concavity of a function. In Figure 4.5,  $q_3$  is keeping the concavity near 1,

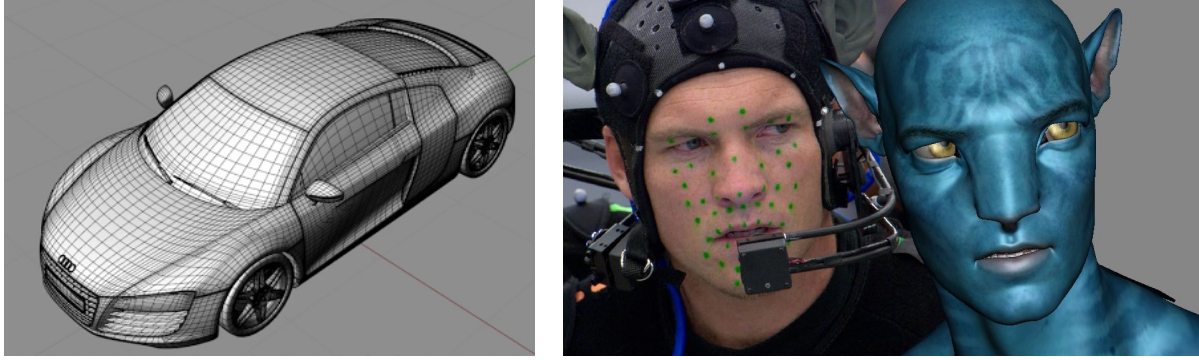


whereas the gray dashed lines are changing the concavity.

Next section will specifically talk about how to use multiple degree 3 polynomials (cubic) to interpolate given points smoothly (twice differentiable).

#### 4.4.1 Cubic Spline

Spline is used heavily in 2D and 3D computer graphics where one wishes to design a smooth curve/surface based on a collection of nodes. There are numerous examples including design of fonts, car manufacture, and animation.



For  $n \geq 2$ , given  $\{(x_i, y_i)\}_{i=0}^n$ , a Cubic spline is to find  $n$  cubics  $q_i(x)$  such that the interpolant

$$s_c(x) = \begin{cases} q_1(x), & x_0 \leq x \leq x_1 \\ q_2(x), & x_1 \leq x \leq x_2 \\ \vdots & \\ q_n(x) & x_{n-1} \leq x \leq x_n \end{cases}$$

is twice differentiable.

According to the discussion earlier, this precisely means

(a) Each cubic goes through two points:  $q_i(x_{i-1}) = y_{i-1}$ ,  $q_i(x_i) = y_i$ ,  $i = 1, 2, \dots, n$

(b) First derivative match:  $q'_i(x_i) = q'_{i+1}(x_i)$ ,  $i = 1, 2, \dots, n-1$

(c) Second derivative match:  $q''_i(x_i) = q''_{i+1}(x_i) = z_i$ ,  $i = 1, 2, \dots, n-1$ . ( $z_i$  to be determined.)

Each cubic has 4 coefficients, so there are a total of  $4(n+1)$  unknowns. We get  $2n + n - 1 + n - 1 = 4n - 2$  constraints from (a)(b)(c). This leaves 2 degrees of freedom, which can be used to enforce various conditions on the endpoints of the spline, resulting in different types of cubic splines.

It is not hard to derive a formula for  $s_c$ . Let  $h_i = x_i - x_{i-1}$  be the length of the  $i$ th subinterval. In the case of equally spaced points,  $h_i = h = \frac{b-a}{n} = \frac{x_n - x_0}{n}$ .

We start with constraint (c): For  $i = 1, \dots, n$ ,  $q''_i$  is a line, and it goes through  $(x_{i-1}, z_{i-1})$  and  $(x_i, z_i)$ , so

$$i = 1, \dots, n \quad q''_i(x) = \frac{z_i - z_{i-1}}{h_i}(x - x_{i-1}) + z_{i-1}. \quad (4.10)$$

Integrating (4.10) twice, we obtain

$$i = 1, \dots, n \quad q'_i(x) = \frac{z_i - z_{i-1}}{2h_i}(x - x_{i-1})^2 + z_{i-1}(x - x_{i-1}) + C_i. \quad (4.11)$$

$$i = 1, \dots, n \quad q_i(x) = \frac{z_i - z_{i-1}}{6h_i}(x - x_{i-1})^3 + \frac{z_{i-1}}{2}(x - x_{i-1})^2 + C_i(x - x_{i-1}) + y_{i-1}. \quad (4.12)$$

The constant term in (4.12) is  $y_{i-1}$  since we used the condition that  $q_i(x_{i-1}) = y_{i-1}$ . Constraint (a) generates  $\frac{z_i - z_{i-1}}{6}h_i^2 + \frac{z_{i-1}}{2}h_i^2 + C_i h_i + y_{i-1} = y_i$ , so

$$C_i = \frac{y_i - y_{i-1}}{h_i} - h_i\left(\frac{z_i}{6} + \frac{z_{i-1}}{3}\right), \quad i = 1, \dots, n. \quad (4.13)$$

Plugging in  $C_i$  in (4.11),

$$q'_i(x) = \frac{z_i - z_{i-1}}{2h_i}(x - x_{i-1})^2 + z_{i-1}(x - x_{i-1}) + \frac{y_i - y_{i-1}}{h_i} - h_i\left(\frac{z_i}{6} + \frac{z_{i-1}}{3}\right).$$

Finally we use constraint (b) to get,  $i = 1, \dots, n-1$ ,

$$\frac{z_i - z_{i-1}}{2h_i}h_i^2 + z_{i-1}h_i + \frac{y_i - y_{i-1}}{h_i} - h_i\left(\frac{z_i}{6} + \frac{z_{i-1}}{3}\right) = \frac{y_{i+1} - y_i}{h_{i+1}} - h_{i+1}\left(\frac{z_{i+1}}{6} + \frac{z_i}{3}\right)$$

which simplifies to  $n-1$  linear equations about  $z_i$ :

$$i = 1, \dots, n-1, \quad \frac{h_i}{6}z_{i-1} + \frac{h_i + h_{i+1}}{3}z_i + \frac{h_{i+1}}{6}z_{i+1} = \frac{y_{i+1} - y_i}{h_{i+1}} - \frac{y_i - y_{i-1}}{h_i}. \quad (4.14)$$

There are  $n-1$  equations, and  $n+1$  unknowns  $(z_0, z_1, \dots, z_n)$ , so we are free to choose  $z_0$  and  $z_n$ . The choice  $z_0 = z_n = 0$  gives what is called the *natural cubic spline*, in which case (4.14) becomes a symmetric tridiagonal system

$$\begin{bmatrix} \frac{h_1 + h_2}{3} & \frac{h_2}{6} & & & \\ \frac{h_2}{6} & \frac{h_2 + h_3}{3} & \frac{h_3}{6} & & \\ & \frac{h_3}{6} & \frac{h_3 + h_4}{3} & \ddots & \\ & & \ddots & \ddots & \frac{h_{n-1}}{6} \\ & & & \frac{h_{n-1}}{6} & \frac{h_{n-1} + h_n}{3} \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ \vdots \\ z_{n-1} \end{bmatrix} = \begin{bmatrix} \frac{y_2 - y_1}{h_2} - \frac{y_1 - y_0}{h_1} \\ \frac{y_3 - y_2}{h_3} - \frac{y_2 - y_1}{h_2} \\ \vdots \\ \vdots \\ \frac{y_n - y_{n-1}}{h_n} - \frac{y_{n-1} - y_{n-2}}{h_{n-1}} \end{bmatrix}. \quad (4.15)$$

Once  $z_i$ 's are solved, we use (4.13) to compute the constants  $C_i$ , and then find the  $n$  cubics using (4.12).

**Example 4.10.** Find the natural cubic spline that interpolates  $(0, 0), (1, 0), (2, 2), (3, 2), (4, -1)$ . Then use it to approximate  $f(0.5)$  and  $f'(0.5)$ .

In this example, we have  $n = 4$ . The points are equally spaced points so  $h_i = 1$ . We already know the boundary condition  $z_0 = z_4 = 0$ . The system (4.15) becomes a  $3 \times 3$  system

$$\begin{bmatrix} 2/3 & 1/6 & 0 \\ 1/6 & 2/3 & 1/6 \\ 0 & 1/6 & 2/3 \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} = \begin{bmatrix} 2 \\ -2 \\ -3 \end{bmatrix}$$

whose solution is  $z_1 = 15/4, z_2 = -3, z_3 = -15/4$ . We can use (4.12) and (4.13) to compute  $q_i$ . For example,  $C_1 = y_1 - y_0 - (z_1/6 + z_0/3) = -5/8$ , so

$$q_1(x) = \frac{z_1 - z_0}{6h_1}(x - x_0)^3 + 0 + C_1(x - x_0) + y_0 = \frac{5}{8}(x - 0)^3 - \frac{5}{8}(x - 0).$$

Similarly,  $C_2 = y_2 - y_1 - (z_2/6 + z_1/3) = 2 - (-1/2 + 5/4) = 5/4$ , so

$$q_2(x) = \frac{z_2 - z_1}{6h_2}(x - x_1)^3 + \frac{z_1}{2}(x - x_1)^2 + C_2(x - x_1) + y_1 = -\frac{9}{8}(x - 1)^3 + \frac{15}{8}(x - 1)^2 + \frac{5}{4}(x - 1).$$

To approximate  $f(0.5)$  and  $f'(0.5)$ , we only need  $q_1$  since 0.5 falls in the first subinterval.

$$f(0.5) \approx q_1(0.5) = -0.2344, \quad f'(0.5) \approx q'_1(0.5) = \frac{15}{8}(0.5)^2 - \frac{5}{8} = -0.1562.$$

**Remark 4.11.** It is the best that we compute the coefficients in Example 4.10 in a vectorized way. We have

index	$x$	$y$	$h$	$z$
0	0	0		0
1	1	0	1	15/4
2	2	2	1	-3
3	3	2	1	-15/4
4	4	-1	1	0

The formula (4.13),  $C_i = \frac{y_i - y_{i-1}}{h_i} - h_i(\frac{z_i}{6} + \frac{z_{i-1}}{3})$  translates to

$$C = \left( \begin{bmatrix} 0 \\ 2 \\ 2 \\ -1 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \\ 2 \\ 2 \end{bmatrix} \right) ./ \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} .* \left( \frac{1}{6} \begin{bmatrix} 15/4 \\ -3 \\ -15/4 \\ 0 \end{bmatrix} + \frac{1}{3} \begin{bmatrix} 0 \\ 15/4 \\ -3 \\ -15/4 \end{bmatrix} \right) = \begin{bmatrix} -5/8 \\ 5/4 \\ 13/8 \\ -7/4 \end{bmatrix}$$

The formula (4.12),  $q_i(x) = \frac{z_i - z_{i-1}}{6h_i}(x - x_{i-1})^3 + \frac{z_{i-1}}{2}(x - x_{i-1})^2 + C_i(x - x_{i-1}) + y_{i-1}$  translates to

$$\text{degree 3 coeff: } \frac{1}{6} \left( \begin{bmatrix} 15/4 \\ -3 \\ -15/4 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 15/4 \\ -3 \\ -15/4 \end{bmatrix} \right) ./ \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 5/8 \\ -9/8 \\ -1/8 \\ 5/8 \end{bmatrix} \quad (4.16)$$

$$\text{degree 2 coeff: } \frac{1}{2} \begin{bmatrix} 0 \\ 15/4 \\ -3 \\ -15/4 \end{bmatrix} = \begin{bmatrix} 0 \\ 15/8 \\ -3/2 \\ -15/8 \end{bmatrix} \quad (4.17)$$

$$\text{degree 1 coeff: } C = \begin{bmatrix} -5/8 \\ 5/4 \\ 13/8 \\ -7/4 \end{bmatrix} \quad (4.18)$$

$$\text{constant term: } \begin{bmatrix} 0 \\ 0 \\ 2 \\ 2 \end{bmatrix} \quad (4.19)$$

## Chapter 4 Exercises

1. Find the polynomial interpolant of (0,-1), (1,1), (2,1) using the Vandermonde matrix. Estimate the value at  $x = 0.5$ .

2. Find the polynomial interpolant of  $(1, 7/4)$ ,  $(2, 3/2)$ ,  $(3.5, 0)$ , and  $(5, 3/4)$  by
  - (a) Vandermonde method (feel free to use Python in finding the coefficients.)
  - (b) Lagrange method (handwritten, don't simplify)
3. Directly write down the polynomial interpolant for two points  $(x_0, y_0), (x_1, y_1)$  using Lagrange's interpolation.
4. Given four points  $p_1 = (-1, 0), p_2 = (0, -1), p_3 = (1, 1), p_4 = (2, 1)$ . Use Newton's interpolation to find the polynomial interpolant going through
  - (a)  $p_1, p_2, p_3$
  - (b)  $p_1, p_2, p_3, p_4$
5. Find  $\|f\|_\infty$  and  $\|f\|_2$  where  $f(x) = e^x - 2$  and the interval is  $[0, 1]$ .
6. Let  $f(x) = 3/2 - x/2, g(x) = \frac{1}{x}$ , defined on  $[1, 2]$ . Find  $\|f - g\|_\infty$  and  $\|f - g\|_2$ . (Hint: for  $\infty$ -norm, you can use derivative to find maximum value (material from Calc I).)
7. The Chebyshev nodes defined in (4.9) is only for interval  $[-1, 1]$ . You need to scale/shift your nodes correspondingly with a different interval. For example, for  $[-2, 2]$ , the interval is scaled by 2, so the nodes should be  $2 \cos \frac{\pi i}{n}$ . For  $[0, 2]$ , the nodes should be  $\cos \frac{\pi i}{n} + 1$ , shifted by 1. In general, we have to do both scaling and shifting.
  - (a) Write down the Chebyshev nodes for the interval  $[-1, 3]$ .  $n = 6$ .
  - (b) Write down the Chebyshev nodes for a general interval  $[a, b]$ .
8. Find the linear interpolant of the three points  $(0, -1), (1, 1), (2, 1)$ .
9. Determine a differentiable interpolant of  $(0, 1), (1, 0), (2, 0)$   $P$  such that
 
$$P(x) = \begin{cases} P_2(x) (\text{degree} \leq 2), & 0 \leq x \leq 1 \\ P_1(x) (\text{degree} \leq 1), & 1 \leq x \leq 2 \end{cases}$$
10. If  $s(x) = -x^2 + x$  when  $0 \leq x \leq 1$  and  $s(x) = P_3(x)$  on  $1 \leq x \leq 2$ .
  - (a) Find all possible cubic  $P_3$  that makes  $s(x)$  twice differentiable on  $[0, 2]$ .
  - (b) Among the ones in (a), find the specific cubic that goes through  $(2, -1)$ .
11. Given the data points  $f(-1) = 0, f(0) = 1, f(2) = 1$  (same data points as in Example 4.1 and Example 4.3).
  - (a) Determine the natural cubic spline that interpolates these points.
  - (b) Use this cubic spline to estimate  $f(-0.8)$ .
  - (c) Draw this cubic spline.
12. Given the five points  $(1, 0), (2, 1), (3, 0), (4, 1), (5, 0)$ ,
  - (a) find the natural cubic spline that interpolate . You should follow Remark 4.11 and your final answer should look like (4.16)-(4.19).

- (b) These five points are coming from an unknown function  $y = f(x)$ . Estimate  $f'(3.2)$ .

## Python Exercises

13. Draw the cubic spline found in Example 4.10. You can directly use the coefficients found in (4.16)-(4.19). You need to use the command `np.polyval`, and each piece needs to have a different color.
14. (a) Write a function `pp = myspline(x, y)` to return the natural cubic spline that interpolates the data  $\mathbf{x} = (x_0, x_1, \dots, x_n)$ ,  $\mathbf{y} = (y_0, y_1, \dots, y_n)$ . The output `pp` should be a  $4 \times n$  matrix where one column corresponds to the coefficients of one piece.

Remark: You will get extra credit if you do not use any loop (for/while) in your function.

- (b) Test `myspline` on Example 4.10. Your test is successful if your function returns

```
array([[ 0.625, -1.125, -0.125,  0.625],
       [ 0.    ,  1.875, -1.5   , -1.875],
       [-0.625,  1.25  ,  1.625, -1.75 ],
       [ 0.    ,  0.    ,  2.    ,  2.    ]])
```

Remark: You will need to use `myspline` in most of the subsequent problems. **If your function is not able to pass the test, use Python's `CubicSpline(x,y,bc_type='natural')` instead for the rest of the problems.**

15. Use `myspline` to construct the natural cubic spline for the following data. Print the coefficient matrix.

	$x$	$f(x)$		$x$	$f(x)$
(a)	-0.5	-0.0247500	(b)	0.1	0.90483742
	-0.25	0.3349375		0.2	0.81873075
	0	1.1010000		0.3	0.74081822
				0.5	0.60653066

16. The data in Exercise 15 were generated using the following functions. Use the cubic spline constructed to do the following approximation and **print the true value as well**.

(a)  $f(x) = x^3 + 4.001x^2 + 4.002x + 1.101$ ; approximate  $f(-1/3)$  and  $f'(-1/3)$ .

(b)  $f(x) = e^{-x}$ ; approximate  $f(0.25)$  and  $f'(0.25)$ .

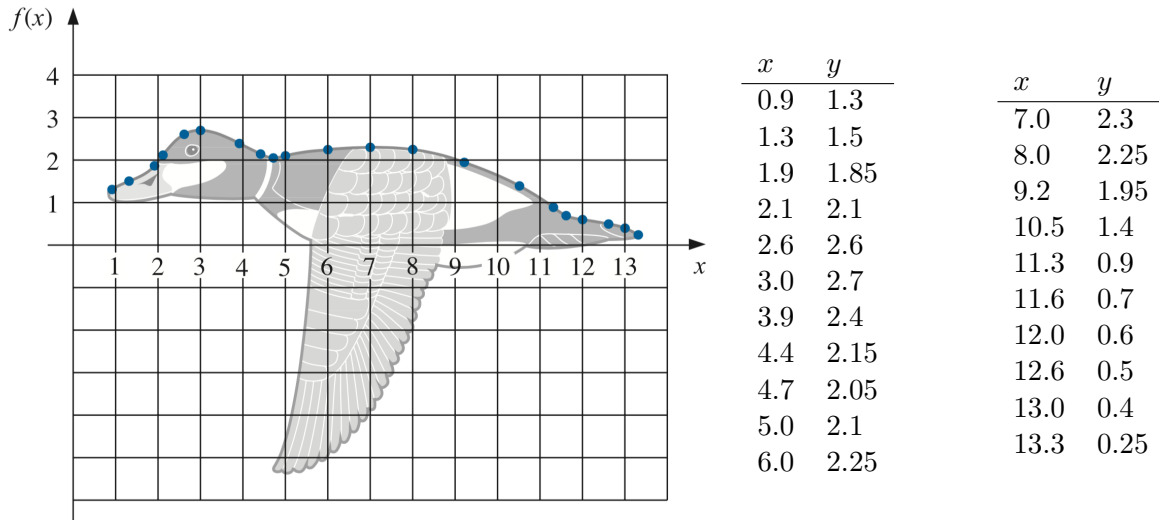
17. Interpolate the Runge function on  $[-1,1]$  at 11 equally spaced points using Python's `CubicSpline`. Plot the interpolant and the Runge function in ONE figure, but with Runge in dash and the spline in solid line. How is this spline performing compared to Figure 4.1(b) or Figure 4.2(b)?

18. The following table is the number of cumulative coronavirus cases in New Mexico in 2020.<sup>2</sup>

Date	3/23	3/24	3/25	3/26	3/27	3/28	3/30	3/31
Cases	83	100	112	136	191	208	281	315
Date	4/1	4/2	4/3	4/4	4/5	4/6	4/7	4/8
Cases	363	403	495	543	624	686	794	865

<sup>2</sup>source: covidtracking.com

- (a) Note that the data on 3/29 is missing. Find an estimation of the cumulative number of cases on 3/29 by interpolating these 16 points via `myspline`.
- (b) Redo (a) only using the data in March. How does your estimation compared to that of (a)? Comment on it.
19. (An application in computer graphics) Below is a picture of a ruddy duck in flight. Use `CubicSpline` to reconstruct its back with the 21 points listed. To save your work, coordinates are already saved in the data file `duck.npy`. If you execute `data = np.load("duck.npy")`, then the first row will be the  $x$  coordinates and the second row will be the  $y$  coordinates. Use the command `plt.axes().set_aspect('equal')` to show the true aspect ratio of a plot.



## Chapter 5

# Numerical Integration

It is crucial to be able to compute a definite integral  $\int_a^b f(x)dx$  in many scientific applications. Numerical integration is necessary when it is difficult or impossible to find an antiderivative of  $f$  that is an elementary function. An example of such an integrand is  $f(x) = e^{-x^2}$ , which is widely used in probability (normal distribution), signal processing, differential equations, etc.

A numerical integration formula is also called a *Quadrature* formula, which is a historical mathematical term that means calculating area. The name comes from the Latin word *quadratus* meaning square.

### 5.1 Basic Quadrature Rules

One way to approximate  $\int_a^b f(x)dx$  is to replace  $f$  by its polynomial interpolant and integrate the polynomial. Formulas derived in this way, using equally-spaced nodes, are called *Newton-Cotes formulas*.

When  $n = 1$ , then integrating  $f$  is replaced by integrating its degree 1 polynomial, i.e., linear interpolant. In other words, the area below the graph of  $y = f(x)$  on  $[a, b]$  is approximated by the trapezoid. See Figure 5.1(a). So the Newton-Cotes formula for two nodes is also called the *Trapezoidal rule*.

To analyze the error, we use formula (4.7) with  $n = 1$ , so

$$f(x) - p_1(x) = \frac{1}{2}f''(\xi_x)(x-a)(x-b).$$

Integrating both sides, we have

$$\begin{aligned}\int_a^b f(x)dx - \int_a^b p_1(x)dx &= \int_a^b \frac{1}{2}f''(\xi_x)(x-a)(x-b)dx \\ &= \frac{1}{2}f''(\eta) \int_a^b (x-a)(x-b)dx, \quad \text{for some } \eta \in [a, b] \\ &= -\frac{1}{12}(b-a)^3 f''(\eta).\end{aligned}$$

Note that the second line is justified by the Mean Value Theorem for integrals since  $(x-a)(x-b)$  has the same sign throughout  $[a, b]$ .

$$p_1(x) = \frac{f(b) - f(a)}{b - a}(x - a) + f(a), \text{ so in summary}$$

Trapezoidal Rule:

$$\int_a^b f(x)dx \approx \int_a^b p_1(x)dx = \frac{f(a) + f(b)}{2}(b - a). \quad (5.1)$$

$$\text{Error bound } \left| \int_a^b f(x)dx - \int_a^b p_1(x)dx \right| = \frac{(b-a)^3}{12} |f''(\eta)| \leq \frac{(b-a)^3}{12} \|f''\|_\infty. \quad (5.2)$$

**Example 5.1.** Use Trapezoidal rule to approximate  $\int_0^2 \sqrt{1+x^2}dx$  and provide an error bound.

Let  $f(x) = \sqrt{x^2 + 1}$ .

$$\int_0^2 \sqrt{1+x^2}dx \approx \frac{f(0) + f(2)}{2}(2 - 0) = 1 + \sqrt{5} \approx 3.2361.$$

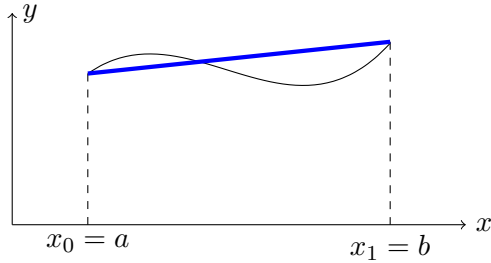
To find the error bound, we compute  $f'(x) = \frac{x}{(1+x^2)^{1/2}}$ ,  $f''(x) = \frac{1}{(1+x^2)^{3/2}}$ .

It is easy to see that  $|f''(x)| = f''(x)$  is decreasing on  $[0,2]$ , so  $\|f''\|_\infty = \max_{x \in [0,2]} |f''(x)| = f''(0) =$

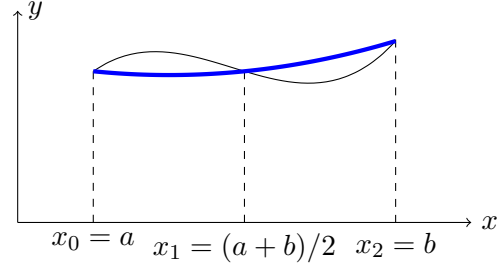
1. Therefore the error bound is  $\frac{(b-a)^3}{12} \|f''\|_\infty = \frac{2^3}{12} * 1 = \frac{2}{3}$ .

For the record, the exact absolute error is 0.2782, which is less than the error bound as expected. See Table 5.1.

Next in line is the Newton-Cotes formula for  $n = 2$ , which is approximating the function by a quadratic polynomial  $p_2$  that agrees with  $f$  at equally spaced points  $x_0 = a, x_1 = (a+b)/2, x_2 = b$ . This is also known as the Simpson's Rule. See Figure 5.1(b).



(a) Trapezoidal Rule



(b) Simpson's Rule

Figure 5.1: Basic Quadrature Rule

Simpson's Rule:

$$\int_a^b f(x)dx \approx \int_a^b p_2(x)dx = \frac{f(a) + 4f(\frac{a+b}{2}) + f(b)}{6}(b - a). \quad (5.3)$$

$$\text{Error bound } \left| \int_a^b f(x)dx - \int_a^b p_2(x)dx \right| = \frac{(b-a)^5}{2880} |f^{(4)}(\eta)| \leq \frac{(b-a)^5}{2880} \|f^{(4)}\|_\infty. \quad (5.4)$$

When use  $p_2$  to approximate  $f$ , it is reasonable to expect the error involves  $f'''$ . But the error term is actually a whole order better! This is because when use equally spaced points, we have  $\int_a^b p_2(x)dx = \int_a^b p_{x_0, x_1, x_2, x_3}(x)dx$ , where  $p_{x_0, x_1, x_2, x_3}$  is the cubic polynomial interpolant at nodes



$a, (a+b)/2, b$  and a fourth one  $x_3$  (can be anywhere in the interval). So the error term is actually 2 orders better than the Trapezoidal rule. The detailed derivation of the error term is omitted here, but can be found in [2, page 231].

**Example 5.2.** Use Simpson's rule to approximate  $\int_0^2 \sqrt{1+x^2} dx$  and provide an error bound.

$$\int_0^2 \sqrt{1+x^2} dx \approx \frac{f(0) + 4f(1) + f(2)}{6} (2-0) = \frac{1}{3} (1 + 4\sqrt{2} + \sqrt{5}) \approx 2.9643..$$

In order to provide an error bound, we have to compute the 4th derivative and find its maximum (in absolute value).  $f^{(4)}(x) = \frac{12x^2 - 3}{(1+x^2)^{7/2}}$  and  $f^{(5)}(x) = \frac{45x - 60x^3}{(1+x^2)^{9/2}}$ .

We have to compute the derivative of  $f^{(4)}$  also because that is typically how we optimize a function. Setting  $f^{(5)}(x) = 0$  gives  $x = 0, \pm\sqrt{3}/2$ , which are the critical points of  $f^{(4)}$ . The optimal value of a function occurs at either critical points or end points. Thus

$$\|f^{(4)}\|_\infty = \max_{x \in [0,2]} |f^{(4)}(x)| = \max\{|f^{(4)}(0)|, |f^{(4)}(\pm\sqrt{3}/2)|, |f^{(4)}(2)|\} = \max\{|-3|, \frac{768\sqrt{7}}{2401}, \frac{9\sqrt{5}}{125}\} = 3.$$

So the error bound is  $\frac{2^5}{2880} * 3 = 1/30$ . (The actual error is  $|2.9579 - 2.9643| = 0.0064$ .)

**Example 5.3.** Use both Trapezoidal and Simpson's rule to approximate  $\int_0^2 e^{-x^2} dx$ .

Let  $f(x) = e^{-x^2}$

$$\text{Trapezoidal: } \approx \frac{f(0) + f(2)}{2} (2-0) = 1 + e^{-4} \approx 1.0183$$

$$\text{Simpson's: } \approx \frac{2-0}{6} (f(0) + 4f(1) + f(2)) = \frac{1 + 4\sqrt{2} + e^{-4}}{6} \approx 0.8299.$$

$f(x)$	$\sqrt{1+x^2}$	$e^{-x^2}$	$x^3$	$x^4$
Exact value	2.9579	0.8821	4.0000	6.4000
Trapezoidal	3.2361	1.0183	8.0000	16.0000
Simpson's	2.9643	0.8299	4.0000	6.6667

Table 5.1: Integrate  $f(x)$  from 0 to 2.

The above table lists the results of the three examples, as well as  $x^3, x^4$ . There is no need to do numerical integration on polynomials, but column 3 and 4 from Table 5.1 can serve two purposes. First, we observe that Simpson's rule gives the exact integral value for  $\int_0^2 x^3 dx$ . This is because the error bound (5.4) has the term  $\|f^{(4)}\|_\infty$ , which is 0 for a polynomial of degree 3 or less. We say that the Simpson's rule has a precision of degree 3. Looking at the error bound for the Trapezoidal rule, we can say that it has a precision of degree 1. Second, we can see that the Trapezoidal (even Simpson's) rule can have big error in approximating integrals. We will address this problem in the next section.

## 5.2 Composite Quadrature Rules

Error bound (5.2) and (5.4) both have the term  $b-a$  raised to some power, so it is expected that we get a poor approximation when the interval is too wide. We circumvent the problem of

a large interval of integration by subdividing the interval  $[a, b]$  into a collection of intervals that are sufficiently small so that the error over each is kept under control. This is called *composite quadrature*.

**Example 5.4.** Use Composite Trapezoidal rule to approximate  $\int_0^2 \sqrt{1+x^2} dx$  by

- (a) subdividing  $[0, 2]$  into 2 subintervals.
- (b) subdividing  $[0, 2]$  into 4 subintervals.

Solution: (a) This is use Trapezoidal rule on  $[0, 1]$  and  $[1, 2]$ . Let  $h = 1$  be the width of each interval.

$$\int_0^2 f(x) dx \approx \frac{f(a) + f((a+b)/2)}{2} h + \frac{f((a+b)/2) + f(b)}{2} h = \frac{1 + \sqrt{2}}{2} + \frac{\sqrt{2} + \sqrt{5}}{2} \approx 3.0322$$

(b) This is to use Trapezoidal rule on 4 equally spaced subintervals. Let  $h = 1/2$  be the width of each interval, and  $0 = x_0 < x_1 < x_2 < x_3 < x_4 = 2$  be the 5 equally spaced nodes.

$$\int_0^2 f(x) dx \approx \sum_{i=1}^4 \frac{f(x_{i-1}) + f(x_i)}{2} h = \frac{1}{4} (f(x_0) + 2f(x_1) + 2f(x_2) + 2f(x_3) + f(x_4)) = \frac{1}{4} (1 + 2\sqrt{1+1/4} + 2\sqrt{2} + 2\sqrt{1+9/4} + \sqrt{5}) \approx 2.9765.$$

When we use composite trapezoidal rule, we divide  $[a, b]$  into  $n$  equally spaced subinterval and the nodes are labelled as  $a = x_0 < x_1 < \dots < x_{n-1} < x_n = b$ .

$$\int_a^b f(x) dx \approx \sum_{i=1}^n \frac{f(x_{i-1}) + f(x_i)}{2} \frac{b-a}{n} = \frac{b-a}{2n} (f(x_0) + 2f(x_1) + 2f(x_2) + \dots + 2f(x_{n-1}) + f(x_n)) \quad (5.5)$$

The error bound is the summation of the error bound on each subinterval. Let  $h = \frac{b-a}{n}$ ,

$$\text{composite trapezoidal error} \leq \sum_{i=1}^n \frac{h^3}{12} \|f''\|_{\infty} = n \frac{(b-a)^3}{12n^3} \|f''\|_{\infty} = \frac{(b-a)^3}{12n^2} \|f''\|_{\infty}. \quad (5.6)$$

The composite Simpson's rule is similar. We again divide  $[a, b]$  into  $n$  equally spaced subinterval and the nodes are labelled as  $a = x_0 < x_1 < \dots < x_{n-1} < x_n = b$ .

$$\int_a^b f(x) dx \approx \sum_{i=1}^n \frac{f(x_{i-1}) + 4f(\frac{x_{i-1}+x_i}{2}) + f(x_i)}{6} \frac{b-a}{n} \quad (5.7)$$

$$= \frac{b-a}{6n} (f(x_0) + 4f(\frac{x_0+x_1}{2}) + 2f(x_1) + \dots + 2f(x_{n-1}) + 4f(\frac{x_{n-1}+x_n}{2}) + f(x_n)) \quad (5.8)$$

The error bound is

$$\text{composite Simpson's error} \leq \sum_{i=1}^n \frac{h^5}{2880} \|f^{(4)}\|_{\infty} = n \frac{(b-a)^5}{2880n^5} \|f^{(4)}\|_{\infty} = \frac{(b-a)^5}{2880n^4} \|f^{(4)}\|_{\infty}. \quad (5.9)$$

**Example 5.5.** Use Composite Simpson's rule to approximate  $\int_0^2 \sqrt{1+x^2} dx$  by subdividing  $[0, 2]$  into 2 subintervals.

Solution:  $h = 1$

$$\int_0^2 f(x) dx \approx \frac{f(0) + 4f(1/2) + f(1)}{6} h + \frac{f(1) + 4f(3/2) + f(2)}{6} h = \frac{1}{6} (1 + 4\sqrt{5/4} + 2\sqrt{2} + 4\sqrt{13/4} + \sqrt{5}) \approx 2.9680.$$

### 5.3 Gaussian Quadrature

All the quadrature rules we have presented so far are in the form of

$$\int_a^b f(x)dx \approx c_0 f(x_0) + c_1 f(x_1) + \cdots + c_n f(x_n), \quad \text{where } \sum_{i=0}^n c_i = b - a.$$

In the trapezoidal rule, we have  $n = 1$  and  $x_0 = a, x_1 = b, c_0 = c_1 = (b - a)/2$ , and we know that this rule has a degree of precision 1. However, we can actually have a better degree of precision if we pick the two nodes differently.

We will assume the interval is  $[-1, 1]$  for now. We want to develop a new quadrature formula  $\int_a^b f(x)dx \approx c_0 f(x_0) + c_1 f(x_1)$  such that this formula is exact for polynomials of degree up to 3. This is to pick  $c_0, c_1, x_0, x_1$  such that the formula will be exact for  $f(x) = 1, f(x) = x, f(x) = x^2, f(x) = x^3$ . So we will get four equations

$$\begin{aligned} 2 &= c_0 + c_1 \\ 0 &= c_0 x_0 + c_1 x_1 \\ 2/3 &= c_0 x_0^2 + c_1 x_1^2 \\ 0 &= c_0 x_0^3 + c_1 x_1^3 \end{aligned}$$

Solving this system of nonlinear equations, we get  $c_1 = c_2 = 1, x_0 = -\sqrt{3}/3, x_1 = \sqrt{3}/3$ .

With these two new nodes (same number of nodes as the trapezoidal rule), we are able to increase the degree of precision from 1 to 3. It turns out that in general, we can find  $n + 1$  nodes  $r_0, r_1, \dots, r_n$  along with  $n + 1$  coefficients  $\{c_i\}_{i=0}^n$  such that the quadrature rule

$$\int_{-1}^1 f(x)dx \approx \sum_{i=0}^n c_i f(r_i) \tag{5.10}$$

has degree of precision  $2n + 1$ , which means that (5.10) is an equality if  $f$  is a polynomial of degree less than or equal to  $2n + 1$ . This is called the *Gaussian Quadrature*.

In order to specify  $r_i$ , we have to first introduce *Legendre polynomials*, a collection  $\{L_0(x), L_1(x), \dots, L_n(x), \dots\}$  that has the following properties:

- For each  $k$ ,  $L_k(x)$  is a polynomial of degree  $k$ .
- $\langle L_m, L_n \rangle := \int_{-1}^1 L_m(x) L_n(x) dx = 0$  whenever  $m \neq n$ . Another way of describing this is that these polynomials are **orthogonal**.

The first few Legendre polynomials are  $L_0(x) = 1, L_1(x) = x, L_2(x) = x^2 - \frac{1}{3}, L_3(x) = x^3 - \frac{3}{5}x$ .

With  $n$  fixed, one can show (proof can be found in [2, Theorem 10.3.1]) that  $\{r_i\}_{i=0}^n$  are exactly the roots of the Legendre polynomial  $L_{n+1}(x)$ . For example, we solved the nodes earlier for  $n = 1$ , and they are the roots of  $L_2(x)$ .

The coefficients  $c_i$  in (5.10) can be found by letting  $f(x) = \frac{\prod_{j \neq i} (x - r_j)}{\prod_{j \neq i} (r_i - r_j)}$  (used in Lagrange interpolation), so

$$c_i = \int_{-1}^1 \frac{(x - r_0) \cdots (x - r_{i-1})(x - r_{i+1}) \cdots (r_i - r_n)}{(r_i - r_0) \cdots (r_i - r_{i-1})(r_i - r_{i+1}) \cdots (r_i - r_n)} dx.$$

For example, when  $n = 2$ , the nodes  $r_i$  can be found by solving  $L_3(x) = x^3 - \frac{3}{5}x = 0$ , which is  $r_0 = -\sqrt{3/5} = -0.7746, r_1 = 0, r_2 = \sqrt{3/5} = 0.7746$ . Then  $c_0 = \int_{-1}^1 \frac{(x - r_1)(x - r_2)}{(r_0 - r_1)(r_0 - r_2)} dx = 0.5556$ .

The roots of the Legendre polynomials are always fixed, consequently so are the coefficients  $c_i$ . So they can be all precomputed and tabulated. See the table on the right for  $n = 1, 2, 3, 4$ .

**Example 5.6.** Approximate  $\int_{-1}^1 e^{-x} dx$  using Gaussian quadrature with 3 nodes.

This means that  $n = 2$ .

Let  $f(x) = e^{-x}$ . Using the table,

$$\int_{-1}^1 e^{-x} dx \approx 0.5556f(0.7746) + 0.8889f(0) + 0.5556f(-0.7746) = 0.5556e^{-0.7746} + 0.8889e^{-0} + 0.5556e^{0.7746} = 2.3505$$

(The actual value is  $-e^{-1} + e = 2.3504$ , so we have a very good approximation with only 3 nodes.)

$n$	Roots $r_i$	Coefficients $c_i$
1	0.5774	1.0000
	-0.5774	1.0000
2	0.7746	0.5556
	0.0000	0.8889
	-0.7746	0.5556
3	0.8611	0.3479
	0.3400	0.6521
	-0.3400	0.6521
	-0.8611	0.3479
4	0.9062	0.2369
	0.5385	0.4786
	0.0000	0.5689
	-0.5385	0.4786
	-0.9062	0.2369

## Gaussian quadrature on arbitrary intervals

By definition, the Legendre polynomials are orthogonal polynomials on the interval  $[-1, 1]$  and the roots are all in  $[-1, 1]$  (need some proof). For an arbitrary interval  $[a, b]$ , we can just do a linear transformation to change it to  $[-1, 1]$ . This is ultimately the same as Exercise 7 of Chapter 4. Let

$$r = \frac{2x - a - b}{b - a} \implies x = [a, b] \text{ corresponds to } r = [-1, 1]$$

So we do a change of variable with  $x = \frac{1}{2}(r(b - a) + a + b)$ , and

$$\int_a^b f(x) dx = \frac{b - a}{2} \int_{-1}^1 f\left(\frac{r(b - a) + a + b}{2}\right) dr \quad (5.11)$$

$$= \frac{b - a}{2} \sum_{i=0}^n c_i f\left(\frac{r_i(b - a) + a + b}{2}\right) \quad (5.12)$$

**Example 5.7.** Approximate  $\int_0^2 \sqrt{1 + x^2} dx$  using Gaussian quadrature with  $n = 2$ .

Let  $f(x) = \sqrt{x^2 + 1}$ . The new nodes are  $x_i = \frac{1}{2}(r_i(2 - 0) + 2 + 0) = r_i + 1$ . Using (5.11),  $\int_0^2 \sqrt{1 + x^2} dx = \frac{2 - 0}{2} \int_{-1}^1 f(r + 1) dr \approx 0.5556f(0.7746 + 1) + 0.8889f(0 + 1) + 0.5556f(-0.7746 + 1) = 2.9584$ .

This is better than the result in Example 5.5 (Composite Simpson with 2 subintervals), which involved 5 nodes.

**Example 5.8.** Approximate  $\int_0^1 e^{-x^2} dx$  using Gaussian quadrature with  $n = 3$ .

The new nodes are  $x_i = \frac{1}{2}(r_i + 1)$ .

$$\begin{aligned} & \int_0^1 e^{-x^2} dx \\ & \approx \frac{1-0}{2} \left[ 0.3479f\left(\frac{0.8611+1}{2}\right) + 0.6521f\left(\frac{0.34+1}{2}\right) + 0.6521f\left(\frac{-0.34+1}{2}\right) + 0.3479f\left(\frac{-0.8611+1}{2}\right) \right] \\ & = 0.7468. \end{aligned}$$

## Chapter 5 Exercises

1. Approximate the following integrals using the Trapezoidal rule.

(a)  $\int_{1/2}^1 x^4 dx$

(b)  $\int_0^{1/2} \frac{2}{x-4} dx$

(c)  $\int_0^4 xe^{-x} dx$

2. Find the actual value of the integrals in Exercise 1. [Use Calc II knowledge]
3. Find the error bound for each integral in Exercise 1, and compare the error bound to the actual error.
4. Repeat Exercise 1 using Simpson's rule.
5. Simpson's rule is to use a quadratic to approximate a function, and therefore it makes sense that this formula has a degree 2 precision. As mentioned, it actually has a degree 3 precision. To verify this, we only need to

(a) Compute  $\int_a^b x^3 dx$ .

(b) Use Simpson's rule to approximate  $\int_a^b x^3 dx$ .

(c) Verify that the values in (a) and (b) are equal.

6. Given an unknown function  $f$ . We know that the trapezoidal rule applied to  $\int_0^2 f(x) dx$  gives the value 5, and the Simpson's rule gives the value 2. What is  $f(1)$ ?
7. Repeat Exercise 1(a)(b) using composite trapezoidal rule with  $n = 4$ . Find the error bounds.
8. Repeat Exercise 1(b)(c) using composite Simpson's rule with  $n = 2$ .
9. What is the degree of precision for Gaussian quadrature when 4 nodes are used?
10. The quadrature formula  $\int_{-1}^1 f(x) dx = c_0 f(-1) + c_1 f(0) + c_2 f(1)$  is exact for all polynomials of degree less than or equal to 2. Determine  $c_0$ ,  $c_1$ , and  $c_2$ .

11. Approximate  $\int_{-1}^1 e^{-x^2} dx$  using Gaussian quadrature with  $n = 2$ . (Extra credit: how can you take advantage of this even function to produce a better approximation with the same  $n$ ?)
12. Repeat Exercise 1(a) using Gaussian quadrature with  $n = 2$ . Remember that a change of variable is needed.
13. Repeat Exercise 1(c) using Gaussian quadrature with 4 nodes. Remember that a change of variable is needed.

## Chapter 6

# Least Squares and Singular Value Decomposition

### 6.1 Positive definite matrices

Given a **symmetric** matrix  $A$ , there are 3 equivalent criteria of positive definite matrices.

- I All eigenvalues of  $A$  are positive.
- II For any vector  $x$ , we have  $x^T Ax \geq 0$  and equality holds only when  $x = 0$ .
- III All leading principal minors are positive.

Leading principal minors are determinants of all upper left submatrices of the original matrix. For an  $n \times n$  matrix  $A$ , let  $A_i$  be the  $i \times i$  submatrix of  $A$  that only takes the first  $i$  rows and first  $i$  columns of  $A$ . Then all principal minors of  $A$  are  $\det(A_1), \det(A_2), \dots, \det(A_n)$ . This is illustrated in Example 6.1.

It is also useful to introduce the concept of positive semidefinite matrices. Given a symmetric matrix  $A$ , there are 2 equivalent criteria of positive semidefinite matrix.

- I All eigenvalues of  $A$  are nonnegative.
- II For any vector  $x$ , we have  $x^T Ax \geq 0$ .

Note: All leading principal minors are nonnegative is not a characterization of positive semidefinite matrix.

We did in Chapter 3 Exercise 3 that

$$\begin{aligned} x^T Ax &= \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{12} & a_{22} & a_{23} \\ a_{13} & a_{23} & a_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \\ &= a_{11}x_1^2 + a_{22}x_2^2 + a_{33}x_3^2 + 2a_{12}x_1x_2 + 2a_{13}x_1x_3 + 2a_{23}x_2x_3. \end{aligned}$$

The pattern is the same for  $x \in \mathbb{R}^2$ :

$$x^T Ax = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} \\ a_{12} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = a_{11}x_1^2 + a_{22}x_2^2 + 2a_{12}x_1x_2.$$

**Example 6.1.** Verify that  $A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 1 & 2 \end{bmatrix}$  is positive definite using all three criteria.

$$\text{I. } \det \begin{bmatrix} 1-\lambda & 0 & 0 \\ 0 & 2-\lambda & 1 \\ 0 & 1 & 2-\lambda \end{bmatrix} = (1-\lambda)[(2-\lambda)^2 - 1] = (1-\lambda)(2-\lambda+1)(2-\lambda-1)$$

eigenvalues are 1, 3, all positive.

$$\text{II. } x^T A x = x_1^2 + 2x_2^2 + 2x_3^2 + 2x_2x_3 = x_1^2 + x_2^2 + x_3^2 + (x_2 + x_3)^2 \geq 0.$$

For the equality to hold, we must have  $x_1 = x_2 = x_3 = x_2 + x_3 = 0$ , which is  $x = 0$ .

$$\text{III. } \det[1] = 1 > 0, \det \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} = 2 > 0, \det A = 1(4-1) = 3 > 0.$$

## 6.2 Cholesky Factorization

For a positive definite matrix  $A$ , we do not need to exchange rows (no pivoting). Moreover, its LU decomposition can be made in such a way that  $L$  is the transpose of  $U$ .

**Definition 6.2.** For a positive definite matrix  $A$ , we can decompose it as  $A = LL^T$  where  $L$  is a lower triangular matrix. This is called the Cholesky factorization.

Cholesky factorization is the LU factorization for PD matrices. We will still use the same matrix in Example 3.2 because it happens to be PD. We first do the regular LU decomposition, and get to

$$\begin{bmatrix} 1 & 2 & 2 \\ 2 & 7 & 7 \\ 2 & 7 & 9 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 2 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 2 \\ 0 & 3 & 3 \\ 0 & 0 & 2 \end{bmatrix} = LU.$$

Currently these two matrices are not transpose of each other. To fix it, we can factor out 3 (row 2) and 2 (row 3) from the  $U$  matrix:

$$\begin{aligned} \begin{bmatrix} 1 & 2 & 2 \\ 2 & 7 & 7 \\ 2 & 7 & 9 \end{bmatrix} &= \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 2 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 2 \\ 0 & 3 & 3 \\ 0 & 0 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 2 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 2 \end{bmatrix} \begin{bmatrix} 1 & 2 & 2 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 2 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & \sqrt{3} & 0 \\ 0 & 0 & \sqrt{2} \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & \sqrt{3} & 0 \\ 0 & 0 & \sqrt{2} \end{bmatrix} \begin{bmatrix} 1 & 2 & 2 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & \sqrt{3} & 0 \\ 2 & \sqrt{3} & \sqrt{2} \end{bmatrix} \begin{bmatrix} 1 & 2 & 2 \\ 0 & \sqrt{3} & \sqrt{3} \\ 0 & 0 & \sqrt{2} \end{bmatrix} \end{aligned}$$

In this example, the diagonals of the  $U$  matrix are all positive, and this is guaranteed by the fact that  $A$  is positive definite. If the diagonal entries of  $U$  are not positive, then we wouldn't be able to take their squared roots.

Cholesky factorization takes about  $\frac{1}{3}n^3$  flops, half of the LU factorization due to symmetry.

## 6.3 Least Squares

Chapter 3 deals with the case that  $A$  is invertible, therefore  $Ax = b$  has a unique solution. But real life linear systems often have no solutions due to different reasons (too many unknowns, noise/perturbation, rounding off error...)

When  $Ax = b$  does not have a solution, we seek the least squares solution

$$x_{LS} = \arg \min_x \|Ax - b\|^2. \quad (6.1)$$

$x_{LS}$  is called the *least squares solution* of  $Ax = b$ .



**Example 6.3.** Let  $A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$  and  $b = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}$ . The system  $Ax = b$  is inconsistent (no solution) because the third equation is  $0=1$ .

We compute  $\|Ax - b\|^2 = (x_1 - 2)^2 + (x_2 - 1)^2 + 1$ , which is minimized when  $x_1 = 2, x_2 = 1$ . The least squares solution is therefore  $(2,1)$ .

In general, we do not need to do an optimization problem to solve the least squares solution. It can be proven (one way is to use geometrical approach using projections) that  $x_{LS}$  is the solution of the *normal equation*

$$A^T Ax = A^T b. \quad (6.2)$$

The normal equation is guaranteed to have one or infinitely many solutions for any matrix  $A$ .

Following Example 6.3, its normal equation is  $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} x = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$ , and we clearly get  $x_{LS} = (2, 1)$ , the same answer as before.

When  $A$  has full column rank (columns are independent),  $A^T A$  is invertible hence there is only one unique solution for (6.2). In this case,  $x_{LS}$  can be mathematically expressed as

$$x_{LS} = (A^T A)^{-1} A^T b, \quad (6.3)$$

even though this is not how it's calculated in a computer (or by hand).

The matrix  $A^T A$  is always positive semidefinite so we can use Cholesky factorization to solve the normal equation.

---

**Least Squares via Cholesky factorization** ( $A \in \mathbb{R}^{m \times n}$  is full column rank,  $m > n$ )

---

Step 1:	Form $B = A^T A, f = A^T b$	$O(mn^2)$
Step 2:	Find Cholesky factorization $B = LL^T$ .	$O(n^3)$
Step 2:	Solve $Ly = f$ . ( $y$ will equal to $L^T x$ .)	$O(n^2)$
Step 3:	Solve $L^T x = y$ .	$O(n^2)$

---

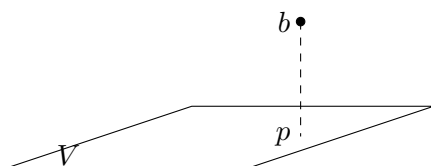


Figure 6.1: Projection onto  $V$

Least squares problem is closely related to projections. If we let  $A = [v_1, v_2, \dots, v_n]$ , then  $\{Ax : x \in \mathbb{R}^n\} = \text{all linear combinations of } v_i\text{'s} = \text{span}\{v_1, v_2, \dots, v_n\} := C(A)$ . By definition (6.1), we are looking for a vector in  $C(A)$  that is closest to  $b$ , and this vector is called the projection of  $b$  onto  $C(A)$ . If the least squares solution  $x_{LS}$  has been found, then the projection of  $b$  onto  $C(A)$  is  $Ax_{LS}$ .

**Example 6.4.** Find the projection of  $b = (2, 1, 1)$  onto  $V = \text{span}\{(1, 0, 0), (0, 1, 0)\}$ .

Solution: We let  $A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$ , so  $\text{span}\{(1, 0, 0), (0, 1, 0)\} = C(A)$ . We have found earlier that  $x_{LS} = (2, 1)$ , so the projection of  $b$  onto  $C(A)$  is  $p = Ax_{LS} = (2, 1, 0)$ .

This makes sense geometrically because  $V$  is the  $xy$ -plane. The projection of any point onto the  $xy$ -plane is to set the  $z$  component to 0.

**Example 6.5** (Line Fitting). Find the best fitting line going through three points  $(-1, 7)$ ,  $(1, 7)$ , and  $(2, 21)$  on the  $xy$  plane.

Let the line be  $y = c_0 + c_1x$ . We plug in all 3 points to get 3 equations:

$$7 = c_0 - c_1$$

$$7 = c_0 + c_1 \quad . \text{ There is likely no solution for this system, so we find}$$

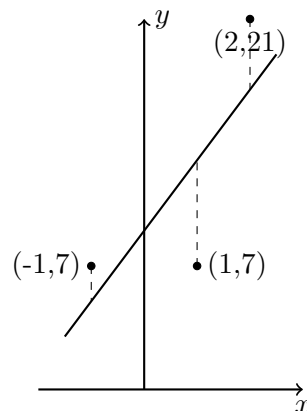
$$21 = c_0 + 2c_1$$

the least square solution of  $\begin{bmatrix} 1 & -1 & | & 7 \\ 1 & 1 & | & 7 \\ 1 & 2 & | & 21 \end{bmatrix}$ .

$$A^T A = \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix}, A^T b = \begin{bmatrix} 35 \\ 42 \end{bmatrix}, 3c_0 + 2c_1 = 35, 2c_0 + 6c_1 = 42 \Rightarrow$$

$$c_0 = 9, c_1 = 4.$$

So the best fitting line is  $y = 4x + 9$ .



By (6.1), this is the best fitting line in the sense that it minimizes  $ey_1^2 + ey_2^2 + ey_3^2$ , where  $ey_i$  is the  $y$  distance from point  $i$  to the line, as drawn in dashed lines.

## 6.4 Orthogonal matrices

**Definition 6.6.** A square matrix  $U = [u_1, u_2, \dots, u_n]$  is an orthogonal matrix if  $u_i^T u_j = 0$  for any  $i \neq j$ , and  $u_i^T u_i = 1$ .

Simply put, an orthogonal matrix is a matrix whose columns are orthogonal (perpendicular) to each other, and all unit norm.

**Example 6.7.**  $P = [e_2, e_3, e_1] = \begin{bmatrix} & & 1 \\ 1 & & \\ & 1 & \end{bmatrix}$  is a permutation matrix. Its columns all have norm 1 and their mutual dot products are all 0, so  $P$  is orthogonal.

**Example 6.8.**  $U = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}$  is orthogonal.

**Example 6.9.** We conclude  $V = [v_1, v_2, v_3] = \begin{bmatrix} 1 & 1 & -1 \\ 1 & -1 & -1 \\ 1 & 0 & 2 \end{bmatrix}$  is not orthogonal because the first column's norm is not 1. Upon further examination, we see that all columns are perpendicular to each other, so if we normalize each column, the resulting matrix  $\tilde{V} = \left[ \frac{v_1}{\|v_1\|}, \frac{v_2}{\|v_2\|}, \frac{v_3}{\|v_3\|} \right] =$

$\begin{bmatrix} \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & 0 & \frac{2}{\sqrt{6}} \end{bmatrix}$  is orthogonal.

When  $U = [u_1, u_2, \dots, u_n]$ , by definition of matrix multiplication,  $U^T U = \begin{bmatrix} u_1^T u_1 & u_1^T u_2 & \cdots & u_1^T u_n \\ u_2^T u_1 & u_2^T u_2 & \cdots & u_2^T u_n \\ \vdots & \vdots & \ddots & \vdots \\ u_n^T u_1 & u_n^T u_2 & \cdots & u_n^T u_n \end{bmatrix}$ .

Therefore we can immediately conclude that

$$U \text{ is orthogonal if and only if } U^T U = I, \text{ the identity matrix.} \quad (6.4)$$

This also implies that if  $U$  is orthogonal, then its transpose is its inverse!

**Example 6.10.** Verify that all the matrices below are orthogonal matrices using (6.4).

$$Q_1 = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}, \quad Q_2 = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad Q_3 = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & 1 \\ -1 & 1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ -1 & 1 & -1 & 1 \end{bmatrix}.$$

## 6.5 Singular Value Decomposition

**Definition 6.11.** Given any matrix  $A$ , whose size is  $m \times n$ , it can be proven that (we will not prove it here)  $A$  can be factored as

$$A = U \Sigma V^T, \quad (6.5)$$

where  $U$  is an  $m \times m$  orthogonal matrix,  $V$  is an  $n \times n$  orthogonal matrix, and  $\Sigma$  is an  $m \times n$  diagonal matrix. The entries on the diagonal of  $\Sigma$  are called the *singular values* of  $A$ . The singular values are always nonnegative. The singular value are often listed in descending order:  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min\{m,n\}} \geq 0$ . The columns of  $U$  and the columns of  $V$  are called the *left-singular vectors* and *right-singular vectors* of  $A$  respectively.

**Example 6.12.** 
$$\begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \boxed{2} & 0 & 0 \\ 0 & \boxed{0} & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \\ 0 & 0 & 1 \end{bmatrix}^T.$$

The singular values are 2,0.

**Example 6.13.** 
$$\begin{bmatrix} -2 & 2.5 & -1 \\ -2 & 2.5 & -1 \\ 0 & 1.5 & -3 \\ 0 & 1.5 & -3 \end{bmatrix} = \begin{bmatrix} 0.5 & 0.5 & -0.5 & 0.5 \\ 0.5 & 0.5 & 0.5 & -0.5 \\ 0.5 & -0.5 & 0.5 & 0.5 \\ 0.5 & -0.5 & -0.5 & -0.5 \end{bmatrix} \begin{bmatrix} \boxed{6} & 0 & 0 \\ 0 & \boxed{3} & 0 \\ 0 & 0 & \boxed{0} \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} -1/3 & -2/3 & -2/3 \\ 2/3 & 1/3 & -2/3 \\ -2/3 & 2/3 & -1/3 \end{bmatrix}^T.$$

The singular values are 6,3,0.

**Example 6.14.** 
$$\begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} = \begin{bmatrix} -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \boxed{\sqrt{2}} & 0 \\ 0 & \boxed{\sqrt{2}} \end{bmatrix} \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}^T.$$

The singular values are  $\sqrt{2}, \sqrt{2}$ , which are quite different from its eigenvalues (they are complex numbers).

In general, given a square matrix, the singular values are not related to its eigenvalues unless it is a symmetric matrix. In particular, if  $A$  is positive semidefinite, then its eigenvalues are the singular values; moreover, its SVD will look like  $A = U \Sigma U^T$  where the columns of  $U$  are the eigenvectors.

## Chapter 6 Exercises

1. Find the determinant of

$$(a) \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \quad (b) \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$$

2. Are the vectors  $v_1 = (3, -1, 1)$ ,  $v_2 = (0, 2, 1)$  linearly independent?

3. Find the rank of  $\begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 1 \\ 2 & 3 & 2 \end{bmatrix}$ . Are the columns of this matrix linearly independent?

4. According to what we have learned, what do you think is the most efficient way to determine whether a matrix is invertible?

5. Is the statement  $\text{span}\{(1, 1, 1), (2, 2, 2)\} = \text{span}\{(1, 1, 1)\}$  true? Explain your answer.

6. If  $Au = \lambda u$ , simplify  $(3A^3 - 4A^2 - A + 2I)u$ .

7. Given  $x = (1, -2, 3)$ ,  $y = (-2, 0, 2)$ ,

(a) Find  $x^T y, \|x\|$ .

(b) Normalize the vector  $y$ .

8.  $B = \begin{bmatrix} 1 & 4 \\ 4 & 7 \end{bmatrix}$ .

(a) Find the eigenvalues and corresponding eigenvectors of  $B$ .

(b) Is this matrix positive definite?

9. Check that  $C = \begin{bmatrix} 2 & 0 & -1 \\ 0 & 2 & -1 \\ -1 & -1 & 3 \end{bmatrix}$  is positive definite using criteria III.

10. Determine if  $D = \begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix}$  is positive definite using criteria II.

11. In equation (6.3), can we simplify the answer as  $(A^T A)^{-1} A^T b = A^{-1} (A^T)^{-1} A^T b = A^{-1} b$ ? Justify your answer.

12. Compute the Cholesky factorization of  $\begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$ .

13. Find the associated normal equation of  $\begin{bmatrix} 2 & -1 & 0 \\ 3 & 1 & 2 \\ -1 & 4 & 5 \\ 1 & 2 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 1 \\ 2 \end{bmatrix}$ . (no need to solve.)

14. Find the least squares solution of the equation  $\begin{bmatrix} 2 & -2 \\ 1 & 1 \\ 2 & 1 \end{bmatrix} x = \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix}$ .

15. Find the projection of  $b = (2, -1, 1)$  onto the space spanned by  $(2, 1, 2)$  and  $(-2, 1, 1)$ .
16. Find the least squares straight line fit to the data points  $(0,0)$ ,  $(1,2)$ ,  $(2,7)$ . Plot the 3 points and the fitted line.
17. Prove that  $A^T A$  is positive semidefinite for any matrix  $A$ .
18. Check if the following matrices are orthogonal. If not, can it be turned into an orthogonal matrix by normalization?

$$(a) A = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \quad (b) B = \frac{1}{3} \begin{bmatrix} 2 & 2 & 1 \\ -2 & 1 & 2 \\ 1 & -2 & 2 \end{bmatrix} \quad (c) C = \begin{bmatrix} 1 & -1 & 2 & -1 \\ -2 & 2 & 3 & 2 \\ 1 & 2 & 0 & -1 \\ 1 & 3 & 2 & 0 \end{bmatrix}$$

19. Compute the singular values of  $A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 2 & 2 \end{bmatrix}$ .
20. Rewrite the SVD of Example 6.12 and Example 6.14 as sum of rank 1 matrices (as how it's done in (??)), and check that they do add up to the original matrix.
21. The singular value decomposition of  $B$  is

$$B = \frac{1}{\sqrt{2}} \begin{bmatrix} \sin \frac{\pi}{4} & \sin \frac{2\pi}{4} & \sin \frac{3\pi}{4} \\ \sin \frac{2\pi}{4} & \sin \frac{4\pi}{4} & \sin \frac{6\pi}{4} \\ \sin \frac{3\pi}{4} & \sin \frac{4\pi}{4} & \sin \frac{9\pi}{4} \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 0.02 & 0 & 0 \\ 0 & 0 & 0.01 & 0 \end{bmatrix} \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & -1 \\ 1 & 1 & -1 & 1 \\ 1 & -1 & 1 & 1 \\ 1 & -1 & -1 & -1 \end{bmatrix}^T$$

- (a) What is the rank of  $B$ ?
- (b) What are the singular values of  $B$ ?
- (c) What is the rank-1 approximation of  $B$ ?
- (d) What is the rank-2 approximation of  $B$ ?
- (e) On which line do the columns of  $B$  approximately lie? On which line do the rows of  $B$  approximately lie?

## Python Exercises

22. First create a random  $10 \times 5$  Gaussian matrix  $A$  and a random  $10 \times 1$  vector  $b$ . It is likely that  $Ax = b$  has no solution. Use Python's `np.linalg.lstsq` command to find its least squares solution, name it `x_python`. Now use Cholesky factorization (see 4 steps on page 54) to find its least squares solution, name it `x_chol` (you can use Python's `np.linalg.cholesky` and `scipy.linalg.solve_triangular`). Finally check that these two solutions are equal.
23. (Linear regression) Use the data in Table 6.1 to predict other college's scores.
  - (a) Suppose the score depends on the 8 features linearly as  $y = w_0 + w_1x_1 + w_2x_2 + \dots + w_8x_8$ . Find the least squares solution of the weights  $w = (w_0, \dots, w_8)$  using the given data.

Table 6.1: Comprehensive score for 14 liberal arts colleges

College (2013)	score	Acc	<20	>50	Rat.	FT	GR	FR	Alum
	$y$	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$
Williams	100	17	71	4	7	93	1	97	58
Amherst	98	13	70	2	9	94	1	98	57
Swarthmore	96	15	74	2	8	93	4	97	46
Middlebury	94	18	68	1	9	94	11	96	55
Pomona	94	14	70	1	8	94	1	98	43
Bowdoin	93	16	68	1	10	93	6	96	50
Wellsley	93	31	69	1	8	93	14	95	46
Carleton	92	31	65	1	9	97	4	97	58
Haverford	91	25	79	1	8	94	6	96	44
Cl. Mc.	90	14	86	2	9	94	11	96	43
Vassar	90	23	68	0.3	8	95	6	97	33
Davidson	89	28	69	0	11	99	6	96	53
Harvey Mudd	89	22	67	2	8	97	21	98	33
US Naval	88	7	61	0	9	94	25	97	21

(b) Second estimate the scores of the 4 colleges below.

College	score	Acc	<20	>50	Rat.	FT	GR	FR	Alum
	$y$	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$
Wash. & Lee	?	18	74	0.2	9	91	14	94	46
Hamilton	?	27	74	1	9	94	21	95	47
Welsleyan	?	24	68	5	9	97	6	96	49
Colby	?	29	69	2	10	93	14	95	41

24. Line fitting using both least squares and SVD on (0,1), (1,3), (2,4), (3,4).

- (a) Find the equation of the best fitted line  $y = ax + b$  using least squares.
- (b) Find the equation of the best fitted line  $y = cx + d$  using SVD.
- (c) Plot these 4 points and both lines in the same graph. Make sure you have a legend.

## Chapter 7

# Preliminaries

### 7.1 Calculus

**Theorem 7.1** (Intermediate Value Theorem). *If  $f(x)$  is continuous on  $[a, b]$ , then for any  $m$  that is in between  $f(a)$  and  $f(b)$ , there exists a number  $c \in [a, b]$  such that  $f(c) = m$ .*

**Theorem 7.2** (Taylor's Theorem). *Suppose  $f \in C^n[a, b]$  and  $f^{(n+1)}$  exists on  $[a, b]$ . Let  $x_0 \in [a, b]$ . For every  $x \in [a, b]$ , there exists a number  $\xi$  between  $x_0$  and  $x$  such that*

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \cdots + \frac{f^{(n)}(x_0)}{n!}(x - x_0)^n + \frac{f^{(n+1)}(\xi)}{(n+1)!}(x - x_0)^{n+1}. \quad (7.1)$$

**Theorem 7.3** (Mean Value Theorem). *If  $f \in C[a, b]$  and  $f$  is differentiable on  $(a, b)$ , then a number  $c$  in  $(a, b)$  exists such that*

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

MVT is often in the form of  $f(b) - f(a) = f'(c)(b - a)$ , which can be rewritten as  $\int_a^b f'(x)dx = f'(c)(b - a)$ . This becomes the MVT for integrals:

$$\int_a^b h(x)dx = h(c)(b - a),$$

but we will present a more general version below.

**Theorem 7.4** (Mean Value Theorem for Integrals). *If  $f \in C[a, b]$ ,  $g$  is integrable on  $[a, b]$  and  $g(x)$  does not change sign on  $[a, b]$ , then there exists a number  $c$  in  $(a, b)$  such that*

$$\int_a^b f(x)g(x)dx = f(c) \int_a^b g(x)dx.$$

## 7.2 Linear Algebra

### 7.2.1 Vector/Matrix algebra

The notation for a matrix will be a capital letter, such as  $A$ , with its corresponding lower letter (with double subscripts), such as  $a_{ij}$  to refer to  $i$ th row,  $j$ th column entry; that is

$$A = [a_{ij}] = [a_{ij}]_{m \times n} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

is an  $m \times n$  matrix. The *transpose* of  $A$  is

$$A^T = [a_{ji}]_{n \times m} = \begin{bmatrix} a_{11} & a_{21} & \cdots & a_{m1} \\ a_{12} & a_{22} & \cdots & a_{m2} \\ \vdots & \vdots & \vdots & \vdots \\ a_{1n} & a_{2n} & \cdots & a_{mn} \end{bmatrix}.$$

The letters towards the end of the alphabet, like  $x, y, z$ , usually indicate vectors. **They are considered as column vectors unless otherwise stated.** For a vector  $x$ ,  $x_i$  will denote its  $i$ th coordinate. The square brackets  $[\cdot]$  is for matrices. Since vectors are a special kind of matrices, we have

- $x = [x_1, x_2, \dots, x_n]$  is a row vector.
- $x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$  is a column vector.
- $x = [x_1, x_2, \dots, x_n]^T$  is a column vector.
- $x = (x_1, x_2, \dots, x_n)$  is a column vector. Yes, if we use parenthesis, it is a column vector!

In consideration of saving space, you will see the last two expressions often.

Given  $x = (x_1, x_2, \dots, x_n)$ ,  $y = (y_1, y_2, \dots, y_n)$ , and a scalar  $a \in \mathbb{R}$ .

$$x + ay = (x_1 + ay_1, x_2 + ay_2, \dots, x_n + ay_n).$$

Let  $A$  be an  $m \times n$  matrix and  $x = (x_1, x_2, \dots, x_n)$ . There are two equivalent ways to compute  $Ax$  (a matrix times a vector). First we write

$$A = [a_{ij}]_{m \times n} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} = \begin{bmatrix} r_1^T \\ r_2^T \\ \vdots \\ r_m^T \end{bmatrix} = [c_1, c_2, \dots, c_n],$$

where  $r_i^T$  is the  $i$ th row of this matrix and  $c_j$  is the  $j$ th column of this matrix.



$$1. Ax = \begin{bmatrix} r_1^T \\ r_2^T \\ \vdots \\ r_m^T \end{bmatrix} x = \begin{bmatrix} r_1^T x \\ r_2^T x \\ \vdots \\ r_m^T x \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + a_{2n}x_n \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + a_{mn}x_n \end{bmatrix}.$$

$$2. Ax = [c_1, c_2, \dots, c_n] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \sum_{i=1}^n x_i c_i. \quad Ax \text{ is a linear combination of columns of } A.$$

A very important special case is when  $A$  is reduced to a  $1 \times n$  matrix (a row vector). Given  $y = (y_1, y_2, \dots, y_n)$ ,

$$y^T x = [y_1, y_2, \dots, y_n] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \sum_{i=1}^n x_i y_i = x^T y = x \cdot y = \langle x, y \rangle. \quad (7.2)$$

The number in (7.2) is called the *inner product* or *dot product* of  $x$  and  $y$ .

For two general matrices  $A = [a_{ij}]_{m \times n}$ ,  $B = [b_{ij}]_{n \times p}$ , their product  $AB$  is  $m \times p$ , and is defined as

$$(AB)_{ij} = \sum_{k=1}^n a_{ik} b_{kj} = a_{i1}b_{1j} + a_{i2}b_{2j} + \dots + a_{in}b_{nj}.$$

For example, given  $y = (y_1, y_2, \dots, y_n)$ ,  $z = (z_1, z_2, \dots, z_m)$ .  $y$  can be viewed as  $n \times 1$  matrix and  $z^T$  is  $1 \times m$ , then  $yz^T$  is  $n \times m$  and

$$yz^T = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} [z_1, z_2, \dots, z_m] = [y_i z_j]_{n \times m} = \begin{bmatrix} y_1 z_1 & y_1 z_2 & \cdots & y_1 z_m \\ y_2 z_1 & y_2 z_2 & \cdots & y_2 z_m \\ \vdots & \vdots & \ddots & \vdots \\ y_n z_1 & y_n z_2 & \cdots & y_n z_m \end{bmatrix}$$

For more examples, please refer to [3, Section 1.2].

A matrix  $U = [u_{ij}]$  is called *upper (lower) triangular* if  $u_{ij} = 0$  for  $i > j$  ( $i < j$ ). A matrix  $D = [d_{ij}]$  is *diagonal* if  $d_{ij} = 0$  for  $i \neq j$ . The *identity* matrix of order  $n$ ,  $I_n = [\delta_{ij}]$ , is an  $n \times n$  diagonal matrix with entries  $\delta_{ii} = 1$ . Moreover, the columns of  $I_n$  are often denoted  $e_i$ , that is  $I_n = [e_1, e_2, \dots, e_n]$ .

A square matrix  $A$  is called *symmetric* if  $A = A^T$ . Visually, it appears that the entries of such matrices are symmetric about the diagonal.

### 7.2.2 Matrix Inversion

An  $n \times n$  matrix  $A$  is *invertible* if there exists  $B$  such that  $AB = I_n$  or  $BA = I_n$ .<sup>1</sup> The matrix  $B$  is called the *inverse* of  $A$ , denoted by  $A^{-1}$ . A matrix is called *singular* if it is not invertible; an invertible matrix is also called *nonsingular*.

Gaussian Elimination is the standard way to find out whether a matrix is invertible, and to compute its inverse when invertible. **If we only need to decide invertibility, we stop at the**

<sup>1</sup>The common definite requires  $AB = BA = I_n$ , but it can be shown that either  $AB = I$  or  $BA = I$  is sufficient.

**end of forward elimination. If number of pivots equals  $n$ , then it's invertible; if number**

**of pivots is less than  $n$ , then it's singular.** Taking  $A = \begin{bmatrix} 1 & 2 & 2 \\ 2 & 7 & 7 \\ 2 & 7 & 9 \end{bmatrix}$  for example, we already

know from Example 3.1 that it's invertible because there are 3 pivots. To find its inverse, we want to find a matrix  $B = [b_1, b_2, b_3]$  such that  $AB = I_3$ , which is  $[Ab_1, Ab_2, Ab_3] = [e_1, e_2, e_3]$ . This is solving three linear systems with the same coefficient matrix, so the row operations can be done in parallel, i.e., perform Gaussian elimination on  $[A|e_1, e_2, e_3]$ . We recycle all the row operations done in Example 3.1:

$$\begin{aligned} & \left[ \begin{array}{ccc|ccc} 1 & 2 & 2 & 1 & 0 & 0 \\ 2 & 7 & 7 & 0 & 1 & 0 \\ 2 & 7 & 9 & 0 & 0 & 1 \end{array} \right] \xrightarrow[\rho 1(-2)+\rho 3]{\rho 1(-2)+\rho 2} \left[ \begin{array}{ccc|ccc} 1 & 2 & 2 & 1 & 0 & 0 \\ 0 & 3 & 3 & -2 & 1 & 0 \\ 0 & 3 & 5 & -2 & 0 & 1 \end{array} \right] \xrightarrow{\rho 2(-1)+\rho 3} \left[ \begin{array}{ccc|ccc} 1 & 2 & 2 & 1 & 0 & 0 \\ 0 & 3 & 3 & -2 & 1 & 0 \\ 0 & 0 & 2 & 0 & -1 & 1 \end{array} \right] \\ & \xrightarrow[\rho 2/3]{\rho 3/2} \left[ \begin{array}{ccc|ccc} 1 & 2 & 2 & 1 & 0 & 0 \\ 0 & 1 & 1 & -\frac{2}{3} & \frac{1}{3} & 0 \\ 0 & 0 & 1 & 0 & -\frac{1}{2} & \frac{1}{2} \end{array} \right] \xrightarrow[\rho 3(-2)+\rho 1]{\rho 3(-1)+\rho 2} \left[ \begin{array}{ccc|ccc} 1 & 2 & 0 & 1 & 1 & -1 \\ 0 & 1 & 0 & -\frac{2}{3} & \frac{5}{6} & -\frac{1}{2} \\ 0 & 0 & 1 & 0 & -\frac{1}{2} & \frac{1}{2} \end{array} \right] \xrightarrow{\rho 2(-2)+\rho 1} \left[ \begin{array}{ccc|ccc} 1 & 0 & 0 & \frac{7}{3} & -\frac{2}{3} & 0 \\ 0 & 1 & 0 & -\frac{2}{3} & \frac{5}{6} & -\frac{1}{2} \\ 0 & 0 & 1 & -1 & -\frac{1}{2} & \frac{1}{2} \end{array} \right]. \end{aligned}$$

**Remark 7.5.** Cancellation DOES NOT hold for matrix algebra, that is  $AB = AC$  DOES NOT imply  $B = C$ . What is true in matrix multiplication is that we can cancel invertible matrices: if  $A$  is invertible and  $AB = AC$ , then we can apply  $A^{-1}$  on both sides to get  $B = C$ .

If  $A, B$  are both invertible, then  $AB$  is also invertible and

$$(AB)^{-1} = B^{-1}A^{-1}.$$

### 7.2.3 Determinant

The determinant is a number associated to a square matrix and can be intuitively understood as a signed area/volume of the parallelepiped associated with the matrix. For our purpose, we need to know

- For an upper or lower triangular matrix, its determinant is the product of diagonal entries.
- The best way (for humans or computers) to compute determinant is to use forward elimination. In Example 3.1, the determinant of the original coefficient matrix is the same as the determinant of  $\begin{bmatrix} 1 & 2 & 2 \\ 0 & 3 & 3 \\ 0 & 0 & 2 \end{bmatrix}$ , which is just the product of the pivots, which is 6. Scaling changes the determinant and row exchange makes the determinant the opposite.

- $\det \begin{bmatrix} a & b \\ c & d \end{bmatrix} = ad - bc.$

- Determinant of a  $3 \times 3$  matrix can be found by row/column expansion. For example, if we expand along first row, then

$$\begin{vmatrix} a & b & c \\ d & e & f \\ g & h & i \end{vmatrix} = a \begin{vmatrix} e & f \\ h & i \end{vmatrix} - b \begin{vmatrix} d & f \\ g & i \end{vmatrix} + c \begin{vmatrix} d & e \\ g & h \end{vmatrix}.$$

### 7.2.4 Linear Independence and span

Given vectors  $v_1, v_2, \dots, v_n$  and let  $A = [v_1, v_2, \dots, v_n]$ , the following are equivalent:

1.  $v_1, v_2, \dots, v_n$  are linearly independent;
2. The solution of  $Ax = 0$  is unique (only  $x = 0$ : the trivial solution).
3.  $\text{rank} A = \text{number of pivots after doing forward elimination} = n$ .

If in addition that  $A$  is square, then there are 2 additional equivalent conditions:

4.  $A$  is invertible.
5.  $\det A \neq 0$ .

Given a set of vectors  $\{v_1, v_2, \dots, v_m\} \subset \mathbb{R}^n$ ,

$$\text{span}\{v_1, v_2, \dots, v_m\} = \text{all linear combinations of } v_1, \dots, v_m = \{c_1 v_1 + \dots + c_m v_m : c_i \in \mathbb{R}\}$$

### 7.2.5 Eigenvalues

Given  $A \in \mathbb{R}^{n \times n}$ , if  $Ax = \lambda x$  for some  $x \neq 0$ , then  $x$  is an eigenvector of  $A$ , and  $\lambda$  is the corresponding eigenvalue. For example, we have  $\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 3 \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ , which means that

$[x, 1]^T$  is an eigenvector of  $A = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$ , with the corresponding eigenvalue 3.

To figure out all the eigenvalues and its corresponding eigenvectors of  $A$ , we need to first find roots of  $\det(A - \lambda I)$ , and second solve  $(A - \lambda I)x = 0$  for each root  $\lambda$ .

**Example 7.6.** Find eigenvalues and eigenvectors of  $A = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$ .

$$\det \begin{bmatrix} 1-\lambda & 2 \\ 2 & 1-\lambda \end{bmatrix} = (1-\lambda)^2 - 4 = (1-\lambda-2)(1-\lambda+2) = (-1-\lambda)(3-\lambda).$$

$$\lambda_1 = -1, \text{ solve } \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} x = 0, u_1 = [1, -1]^T.$$

$$\lambda_2 = 3, \text{ solve } \begin{bmatrix} -2 & 2 \\ 2 & -2 \end{bmatrix} x = 0, u_2 = [1, 1]^T.$$

### 7.2.6 Norm

Let  $x \in \mathbb{R}^n$ , then its (Euclidean) *norm* is defined as

$$\|x\| := \left( \sum_{i=1}^n x_i^2 \right)^{1/2} = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}. \quad (7.3)$$

Note that  $x^T x = \langle x, x \rangle = \|x\|^2$ .  $x$  is called a unit (norm) vector if  $\|x\| = 1$ .

Given an arbitrary vector  $x$ ,  $\frac{x}{\|x\|}$  is a unit vector that has the same direction as  $x$ . This is called normalization of  $x$ .

### 7.3 Complexity

Big O notation is frequently used in mathematics and computer science.

Let  $f(x), g(x)$  be both real functions. We write

$$f(x) = O(g(x)) \text{ as } x \rightarrow \infty$$

if and only if there exist  $M$  and  $x_0$  such that  $|f(x)| \leq Mg(x)$  for all  $x \geq x_0$ .

This notation is used a lot in algorithm analysis. For example, for a polynomial, we have

$$\sum_{k=1}^M a_k x^k = O(x^M).$$

# Bibliography

- [1] Richard L. Burden, and J. Douglas Faires. *Numerical analysis*. Cengage Learning, 9th (2010).
- [2] Anne, Greenbaum, and Timothy P. Chartier. *Numerical methods: design, analysis, and computer implementation of algorithms*. Princeton University Press, 2012.
- [3] Ab Mooijaart, Matthijs Warrens, and Eeke van der Burg, *An introduction to matrix algebra with Matlab*, manuscript, 2006
- [4] Lloyd N. Trefethen, and David Bau III. *Numerical linear algebra*. Vol. 50. Siam, 1997.