

# MATH50003 Numerical Analysis

Sheehan Olver

January 16, 2026



# Contents

<b>I Calculus on a Computer</b>	<b>5</b>
I.1 Rectangular rule . . . . .	6
I.1.1 Lab and problem sheet . . . . .	8
I.2 Divided Differences . . . . .	8
I.2.1 Lab and problem sheet . . . . .	9
I.3 Dual Numbers . . . . .	9
I.3.1 Differentiating polynomials . . . . .	10
I.3.2 Differentiating other functions . . . . .	11
I.3.3 Lab and problem sheet . . . . .	13
I.4 Newton's method . . . . .	13
I.4.1 Lab and problem sheet . . . . .	14
<b>II Representing Numbers</b>	<b>15</b>
II.1 Reals . . . . .	16
II.1.1 Real numbers in binary . . . . .	16
II.1.2 Floating-point numbers . . . . .	17
II.1.3 IEEE floating-point numbers . . . . .	18
II.1.4 Sub-normal and special numbers . . . . .	19
II.2 Floating Point Arithmetic . . . . .	19
II.2.1 Bounding errors in floating point arithmetic . . . . .	21
II.2.2 Idealised floating point . . . . .	22
II.2.3 Divided differences floating point error bound . . . . .	23
<b>A Asymptotics and Computational Cost</b>	<b>27</b>
A.1 Asymptotics as $n \rightarrow \infty$ . . . . .	27
A.2 Asymptotics as $x \rightarrow x_0$ . . . . .	28
A.3 Computational cost . . . . .	29

<b>B Integers</b>	<b>31</b>
B.0.1 Unsigned Integers . . . . .	31
B.0.2 Signed integer . . . . .	32
B.0.3 Hexadecimal format . . . . .	33

# Chapter I

## Calculus on a Computer

In this first chapter we explore the basics of mathematical computing and numerical analysis. In particular we investigate the following mathematical problems which can not in general be solved exactly:

1. Integration. General integrals have no closed form expressions. Can we instead use a computer to approximate the values of definite integrals? Numerical integration underpins much of modern scientific computing and simulations of physical systems modelled by partial differential equations.
2. Differentiation. Differentiating a formula as in calculus is usually algorithmic, however, it is often needed to compute derivatives without access to an underlying formula, eg, a function defined only in code. Can we use a computer to approximate derivatives? A very important application is in Machine Learning, where there is a need to compute gradients in training neural networks.
3. Root finding. There is no general formula for finding roots (zeros) of arbitrary functions, or even polynomials that are of degree 5 (quintics) or higher. Can we compute roots of general functions using a computer?

Each chapter is divided into sections that roughly correspond to individual lectures. In this chapter we investigate solving the above computational problems:

1. I.1 Rectangular rule: we review the rectangular rule for integration and deduce the *convergence rate* of the approximation. In the lab/problem sheet we investigate its implementation as well as extensions to the Trapezium rule.
2. I.2 Divided differences: we investigate approximating derivatives by a divided difference and again deduce the convergence rates. In the lab/problem sheet we extend the approach to the central differences formula and computing second derivatives. We also observe a mystery: the approximations may have significant errors in practice, and there is a limit to the accuracy.
3. I.3 Dual numbers: we introduce the algebraic notion of a *dual number* which allows the implementation of *forward-mode automatic differentiation*, a high accuracy alternative to divided differences for computing derivatives.

4. I.4 Newton's method: Newton's method is a basic approach for computing roots/zeros of a function. We use dual numbers to implement this algorithm.

Each week there are labs and problem sheets that further explore the mathematical material introduced in each section. The labs generally explore practical implementation and the impact of implementing methods in computer arithmetic. The problem sheets dig deeper into analysis of other methods and phenomena observed in the labs. The material introduced in the labs and problem sheets is also examinable so it's important to study these as well.

## I.1 Rectangular rule

One possible definition for an integral is the limit of a Riemann sum, for example:

$$\int_a^b f(x)dx = \lim_{n \rightarrow \infty} h \sum_{j=1}^n f(x_j)$$

where  $x_j = a + jh$  are evenly spaced points dividing up the interval  $[a, b]$ , that is with the *step size*  $h = (b - a)/n$ . This suggests an algorithm known as the (*right-sided*) *rectangular rule* for approximating an integral: choose  $n$  large so that

$$\int_a^b f(x)dx \approx h \sum_{j=1}^n f(x_j).$$

We will show that the error in approximation is bounded by  $C/n$  for some constant  $C$ . This can be expressed using “Big-O” notation:

$$\int_a^b f(x)dx = h \sum_{j=1}^n f(x_j) + O(1/n).$$

In these notes we consider the “Analysis” part of “Numerical Analysis”: we want to *prove* the convergence rate of the approximation, including finding an explicit expression for the constant  $C$ .

To tackle this question we consider the error incurred on a single panel  $(x_{j-1}, x_j)$ , then sum up the errors on rectangles.

Now for a secret. There are only so many tools available in analysis (especially at this stage of your career), and one can make a safe bet that the right tool in any analysis proof is either (1) integration-by-parts, (2) geometric series or (3) Taylor series. In this case we use (1):

**Lemma 1** ((Right-sided) Rectangular Rule error on one panel). *Assuming  $f$  is differentiable on  $[a, b]$  and its derivative is integrable we have*

$$\int_a^b f(x)dx = (b - a)f(b) + \delta$$

where  $|\delta| \leq M(b - a)^2$  for  $M = \sup_{a \leq x \leq b} |f'(x)|$ .

**Proof** We write

$$\begin{aligned} \int_a^b f(x)dx &= \int_a^b (x - a)' f(x)dx = [(x - a)f(x)]_a^b - \int_a^b (x - a)f'(x)dx \\ &= (b - a)f(b) + \underbrace{\left( - \int_a^b (x - a)f'(x)dx \right)}_{\delta}. \end{aligned}$$

Recall that we can bound the absolute value of an integral by the supremum of the integrand times the width of the integration interval:

$$\left| \int_a^b g(x)dx \right| \leq (b-a) \sup_{a \leq x \leq b} |g(x)|.$$

The lemma thus follows since

$$\begin{aligned} \left| \int_a^b (x-a)f'(x)dx \right| &\leq (b-a) \sup_{a \leq x \leq b} |(x-a)f'(x)| \\ &\leq (b-a) \sup_{a \leq x \leq b} |x-a| \sup_{a \leq x \leq b} |f'(x)| \\ &\leq M(b-a)^2. \end{aligned}$$

■

Now summing up the errors in each panel gives us the error of using the Rectangular rule:

**Theorem 1** (Rectangular Rule error). *Assuming  $f$  is differentiable on  $[a, b]$  and its derivative is integrable we have*

$$\int_a^b f(x)dx = h \sum_{j=1}^n f(x_j) + \delta$$

where  $|\delta| \leq M(b-a)h$  for  $M = \sup_{a \leq x \leq b} |f'(x)|$ ,  $h = (b-a)/n$  and  $x_j = a + jh$ .

**Proof** We split the integral into a sum of smaller integrals:

$$\int_a^b f(x)dx = \sum_{j=1}^n \int_{x_{j-1}}^{x_j} f(x)dx = \sum_{j=1}^n [(x_j - x_{j-1})f(x_j) + \delta_j] = h \sum_{j=1}^n f(x_j) + \underbrace{\sum_{j=1}^n \delta_j}_{\delta}$$

where  $\delta_j$ , the error on each panel as in the preceding lemma, satisfies

$$|\delta_j| \leq (x_j - x_{j-1})^2 \sup_{x_{j-1} \leq x \leq x_j} |f'(x)| \leq Mh^2.$$

Thus using the triangular inequality we have

$$|\delta| = \left| \sum_{j=1}^n \delta_j \right| \leq \sum_{j=1}^n |\delta_j| \leq Mn h^2 = M(b-a)h.$$

■

Note a consequence of this lemma is that the approximation converges as  $n \rightarrow \infty$  (i.e.  $h \rightarrow 0$ ). In the labs and problem sheets we will consider the left-sided rule:

$$\int_a^b f(x)dx \approx h \sum_{j=0}^{n-1} f(x_j).$$

We also consider the *Trapezium rule*. Here we approximate an integral by an affine function:

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

Subdividing an interval  $a = x_0 < x_1 < \dots < x_n = b$  and applying this approximation separately on each subinterval  $[x_{j-1}, x_j]$ , where  $h = (b-a)/n$  and  $x_j = a + jh$ , leads to the approximation

$$\int_a^b f(x)dx \approx \frac{h}{2} f(a) + h \sum_{j=1}^{n-1} f(x_j) + \frac{h}{2} f(b)$$

We shall see both experimentally and provably that this approximation converges faster than the rectangular rule.

### I.1.1 Lab and problem sheet

In the lab, we explore the practical implementation of the right-sided rectangular rule and extensions to other rules like the left-sided rectangular rule and trapezium rule. We also see how linear convergence ( $O(h) = O(1/n)$ ) can be deduced *experimentally*: by comparing an implementation of the rule to specific integrals with known formulæ we can compute the error, and determine its rate of decay visually by plotting it. In particular, we deduce that the Trapezium rule converges much faster to the true value of the integral than the other rules. In the problem sheet we explore the *analysis* of these other rules, proving that the Trapezium rule converges to the true integral at a faster quadratic ( $O(h^2)$ ) error rate. This is a guarantee that the integral can be computed much more accurately for the same amount of work by taking into account the analysis, highlighting the important contribution of analysis in the construction of algorithms.

## I.2 Divided Differences

Given a function, how can we approximate its derivative at a point? We consider an intuitive approach to this problem using *(Right-sided) Divided Differences*:

$$f'(x) \approx \frac{f(x+h) - f(x)}{h}$$

Note by the definition of the derivative we know that this approximation will converge to the true derivative as  $h \rightarrow 0$ . But in numerical approximations we also need to consider the rate of convergence.

Now in the previous section I mentioned there are three basic tools in analysis: (1) integration-by-parts, (2) geometric series or (3) Taylor series. In this case we use (3):

**Proposition 1** (divided differences error). *Suppose that  $f$  is twice-differentiable on the interval  $[x, x+h]$ . The error in approximating the derivative using divided differences is*

$$f'(x) = \frac{f(x+h) - f(x)}{h} + \delta$$

where  $|\delta| \leq Mh/2$  for  $M = \sup_{x \leq t \leq x+h} |f''(t)|$ .

**Proof** Follows immediately from Taylor's theorem: recall that

$$f(x+h) = f(x) + f'(x)h + \frac{f''(t)}{2}h^2$$

for some  $t \in [x, x+h]$ . Rearranging we get

$$f'(x) = \frac{f(x+h) - f(x)}{2} + \underbrace{\left(-\frac{f''(t)}{2h^2}\right)}_{\delta}.$$

We then bound:

$$|\delta| \leq \left| \frac{f''(t)}{2}h \right| \leq \frac{Mh}{2}.$$



Unlike the rectangular rule, the computational cost of computing the divided difference is independent of  $h$ ! We only need to evaluate a function  $f$  twice and do a single division. Here we are assuming that the computational cost of evaluating  $f$  is independent of the point of evaluation. Later we will investigate the details of how computers work with numbers via floating point, and confirm that this is a sensible assumption.

In the lab we investigate the convergence rate of these approximations (in particular, that central differences is more accurate than standard divided differences) and observe that they too suffer from unexplained (for now) loss of accuracy as  $h \rightarrow 0$ . In the problem sheet we prove the theoretical convergence rate, which is never realised because of these errors.

### I.2.1 Lab and problem sheet

In the labs and problem sheets we explore alternative versions of divided differences. Left-side divided differences evaluates to the left of the point where we wish to know the derivative:

$$f'(x) \approx \frac{f(x) - f(x-h)}{h}$$

and central differences evaluates both left and right:

$$f'(x) \approx \frac{f(x+h) - f(x-h)}{2h}.$$

We can further arrive at an approximation to the second derivative by composing a left- and right-sided finite difference:

$$f''(x) \approx \frac{f'(x+h) - f'(x)}{h} \approx \frac{\frac{f(x+h)-f(x)}{h} - \frac{f(x)-f(x-h)}{h}}{h} = \frac{f(x+h) - 2f(x) + f(x-h)}{h^2}.$$

The lab explores these approximations *experimentally*, and we will observe that central differences converges much faster to the true value of the derivative as  $h$  becomes moderately small.

An important distinction between rectangular rules and divided difference is that the computational cost of divided differences is independent of  $h$ : we can choose  $h$  arbitrarily and the approximation will take the same amount of time. This raises a question: why not just set  $h$  ridiculously small so that the approximation is extremely accurate? Unfortunately, we will observe in the lab a serious issue: if  $h$  becomes too small, the error mysteriously starts to grow, and hence these rules do not actually converge to the true value of the derivatives! Thus there is a limitation to how accurate one can approximate a derivative using divided differences, an issue we will overcome in the next section by re-thinking derivatives in an algebraic way.

The problem sheet explores the *analysis* of divided difference rules, proving the precise theoretical convergence rates observed for moderately small  $h$ . This presents a bit of a conundrum: why does the theory say the method converges but in practice it diverges, and spectacularly so! This is a mystery that we will return to later, by understanding how computer arithmetic with real numbers works.

## I.3 Dual Numbers

In this section we introduce a mathematically beautiful alternative to divided differences for computing derivatives: *dual numbers*. These are a commutative ring that *exactly* compute

derivatives, which when implemented on a computer gives very high-accuracy approximations to derivatives. They underpin forward-mode [automatic differentiation](#). Automatic differentiation is a basic tool in Machine Learning for computing gradients necessary for training neural networks.

**Definition 1** (Dual numbers). Dual numbers  $\mathbb{D}$  are a commutative ring (over  $\mathbb{R}$ ) generated by 1 and  $\epsilon$  such that  $\epsilon^2 = 0$ , that is,

$$\mathbb{D} := \{a + b\epsilon \quad : \quad a, b \in \mathbb{R}, \quad \epsilon^2 = 0\}.$$

This is very much analogous to complex numbers, which are a field generated by 1 and  $i$  such that  $i^2 = -1$ , that is,

$$\mathbb{C} := \{a + bi \quad : \quad a, b \in \mathbb{R}, \quad i^2 = -1\}.$$

Compare multiplication of each number type which falls out of the rules of the generators:

$$\begin{aligned} (a + bi)(c + di) &= ac + (bc + ad)i + bd i^2 = ac - bd + (bc + ad)i, \\ (a + b\epsilon)(c + d\epsilon) &= ac + (bc + ad)\epsilon + bd\epsilon^2 = ac + (bc + ad)\epsilon. \end{aligned}$$

And just as we view  $\mathbb{R} \subset \mathbb{C}$  by equating  $a \in \mathbb{R}$  with  $a + 0i \in \mathbb{C}$ , we can view  $\mathbb{R} \subset \mathbb{D}$  by equating  $a \in \mathbb{R}$  with  $a + 0\epsilon \in \mathbb{D}$ .

Conceptually, dual numbers can be thought of as introducing an infinitesimally small  $\epsilon$ , where  $\epsilon^2$  is so small it is treated as zero. This is the intuitive reason they allow for differentiation of functions. But we do not need to appeal to this calculus-like interpretation, instead, their construction and relationship to differentiation can be accomplished using purely algebraic reasoning.

### I.3.1 Differentiating polynomials

Polynomials evaluated on dual numbers are well-defined as they depend only on the operations  $+$  and  $*$ . From the formula for multiplication of dual numbers we deduce that evaluating a polynomial at a dual number  $a + b\epsilon$  tells us the derivative of the polynomial at  $a$ :

**Theorem 2** (polynomials on dual numbers). *Suppose  $p$  is a polynomial. Then*

$$p(a + b\epsilon) = p(a) + bp'(a)\epsilon$$

#### Proof

First consider  $p(x) = x^n$  for  $n \geq 0$ . The cases  $n = 0$  and  $n = 1$  are immediate. For  $n > 1$  we have by induction:

$$(a + b\epsilon)^n = (a + b\epsilon)(a + b\epsilon)^{n-1} = (a + b\epsilon)(a^{n-1} + (n-1)b a^{n-2}\epsilon) = a^n + b n a^{n-1}\epsilon.$$

For a more general polynomial

$$p(x) = \sum_{k=0}^n c_k x^k$$

the result follows from linearity:

$$p(a + b\epsilon) = \sum_{k=0}^n c_k (a + b\epsilon)^k = c_0 + \sum_{k=1}^n c_k (a^k + kba^{k-1}\epsilon) = \sum_{k=0}^n c_k a^k + b \sum_{k=1}^n c_k k a^{k-1}\epsilon = p(a) + bp'(a)\epsilon.$$



**Example 1** (differentiating polynomial). Consider computing  $p'(2)$  where

$$p(x) = (x - 1)(x - 2) + x^2.$$

We can use dual numbers to differentiate, avoiding expanding in monomials or applying rules of differentiating:

$$p(2 + \epsilon) = (1 + \epsilon)\epsilon + (2 + \epsilon)^2 = \epsilon + 4 + 4\epsilon = 4 + \underbrace{5\epsilon}_{p'(2)}.$$

### I.3.2 Differentiating other functions

We can extend real-valued differentiable functions to dual numbers in a similar manner. First, consider a standard function with a Taylor series (e.g. cos, sin, exp, etc.)

$$f(x) = \sum_{k=0}^{\infty} f_k x^k$$

so that  $a$  is inside the radius of convergence. This leads naturally to a definition on dual numbers:

$$\begin{aligned} f(a + b\epsilon) &= \sum_{k=0}^{\infty} f_k (a + b\epsilon)^k = f_0 + \sum_{k=1}^{\infty} f_k (a^k + ka^{k-1}b\epsilon) = \sum_{k=0}^{\infty} f_k a^k + \sum_{k=1}^{\infty} f_k k a^{k-1} b\epsilon \\ &= f(a) + bf'(a)\epsilon. \end{aligned}$$

More generally, given a differentiable function (which may not have a Taylor series) we can extend it to dual numbers:

**Definition 2** (dual extension). Suppose a real-valued function  $f : \Omega \rightarrow \mathbb{R}$  is differentiable in  $\Omega \subset \mathbb{R}$ . We can construct the *dual extension*  $\underline{f} : \Omega + \epsilon\mathbb{R} \rightarrow \mathbb{D}$  by defining

$$\underline{f}(a + b\epsilon) := f(a) + bf'(a)\epsilon.$$

By viewing  $\mathbb{R} \subset \mathbb{D}$ , it is natural to reuse the notation  $f$  for the dual extension, hence when there's no chance of confusion we will identify  $f(a + b\epsilon) \equiv \underline{f}(a + b\epsilon)$ .

Thus, for basic functions we have natural extensions:

$$\begin{aligned} \exp(a + b\epsilon) &:= \exp(a) + b \exp(a)\epsilon & (a, b \in \mathbb{R}) \\ \sin(a + b\epsilon) &:= \sin(a) + b \cos(a)\epsilon & (a, b \in \mathbb{R}) \\ \cos(a + b\epsilon) &:= \cos(a) - b \sin(a)\epsilon & (a, b \in \mathbb{R}) \\ \log(a + b\epsilon) &:= \log(a) + \frac{b}{a}\epsilon & (a \in (0, \infty), b \in \mathbb{R}) \\ \sqrt{a + b\epsilon} &:= \sqrt{a} + \frac{b}{2\sqrt{a}}\epsilon & (a \in (0, \infty), b \in \mathbb{R}) \\ |a + b\epsilon| &:= |a| + b \operatorname{sign} a \epsilon & (a \in \mathbb{R} \setminus \{0\}, b \in \mathbb{R}) \end{aligned}$$

provided the function is differentiable at  $a$ . Note the last example does not have a convergent Taylor series (at 0) but we can still extend it where it is differentiable.

Going further, we can add, multiply, and compose such dual-extensions. And the beauty is these automatically satisfy the right properties to be dual-extensions themselves, thus

allowing for differentiation of complicated functions built from basic differentiable building blocks.

The following lemma shows that addition and multiplication in some sense “commute” with the dual-extension, hence we can recover the product rule from dual number multiplication:

**Lemma 2** (addition/multiplication). *Suppose  $f, g : \Omega \rightarrow \mathbb{R}$  are differentiable for  $\Omega \subset \mathbb{R}$  and  $c \in \mathbb{R}$ . Then for  $a \in \Omega$  and  $b \in \mathbb{R}$  we have*

$$\begin{aligned}\underline{f+g}(a+b\epsilon) &= \underline{f}(a+b\epsilon) + \underline{g}(a+b\epsilon) \\ \underline{cf}(a+b\epsilon) &= c\underline{f}(a+b\epsilon) \\ \underline{fg}(a+b\epsilon) &= \underline{f}(a+b\epsilon)\underline{g}(a+b\epsilon)\end{aligned}$$

**Proof** The first two are immediate due to linearity:

$$\begin{aligned}\underline{(f+g)}(a+b\epsilon) &= (f+g)(a) + b(f+g)'(a)\epsilon \\ &= (f(a) + bf'(a)\epsilon) + (g(a) + bg'(a)\epsilon) = \underline{f}(a+b\epsilon) + \underline{g}(a+b\epsilon), \\ \underline{cf}(a+b\epsilon) &= (cf)(a) + b(cf)'(a)\epsilon = c(f(a) + bf'(a)\epsilon) = c\underline{f}(a+b\epsilon).\end{aligned}$$

The last property essentially captures the product rule of differentiation:

$$\begin{aligned}\underline{fg}(a+b\epsilon) &= f(a)g(a) + b(f(a)g'(a) + f'(a)g'(a))\epsilon \\ &= (f(a) + bf'(a)\epsilon)(g(a) + bg'(a)\epsilon) = \underline{f}(a+b\epsilon)\underline{g}(a+b\epsilon).\end{aligned}$$

■

Furthermore composition recovers the chain rule:

**Lemma 3** (composition). *Suppose  $f : \Gamma \rightarrow \mathbb{R}$  and  $g : \Omega \rightarrow \Gamma$  are differentiable in  $\Omega, \Gamma \subset \mathbb{R}$ . Then*

$$\underline{(f \circ g)}(a+b\epsilon) = \underline{f}(\underline{g}(a+b\epsilon))$$

**Proof** Again it falls out of the properties of dual numbers:

$$\underline{(f \circ g)}(a+b\epsilon) = f(g(a)) + bg'(a)f'(g(a))\epsilon = \underline{f}(g(a) + bg'(a)\epsilon) = \underline{f}(g(a+b\epsilon))$$

■

A simple corollary is that any function defined in terms of addition, multiplication, composition, etc. of basic functions with dual-extensions will be differentiable via dual numbers. In this following example we see a practical realisation of this, where we differentiate a function by just evaluating it on dual numbers, implicitly, using the dual-extension for the basic build blocks:

**Example 2** (differentiating non-polynomial). Consider differentiating  $f(x) = \exp(x^2 + \cos x)$  at the point  $a = 1$ , where we automatically use the dual-extension of  $\exp$  and  $\cos$ . We can differentiate  $f$  by simply evaluating on the duals:

$$f(1+\epsilon) = \exp(1+2\epsilon + \cos 1 - \sin 1\epsilon) = \exp(1+\cos 1) + \exp(1+\cos 1)(2-\sin 1)\epsilon.$$

Therefore we deduce that

$$f'(1) = \exp(1+\cos 1)(2-\sin 1).$$

### I.3.3 Lab and problem sheet

In the lab we explore how one can turn this mathematical idea into a practical implementation on a computer, giving a basic version of *forward-mode automatic differentiation*. This is a concept that underpins machine learning, which uses *reverse-mode automatic differentiation* to compute gradients when performing stochastic gradient descent. In order to implement dual numbers, we will introduce the concept of a *type*: a data structure with fields. For example, we will implement a type `Rat` for representing rationals  $p/q$ , where the type has two fields (`p` and `q`). Basic arithmetic operations like `+` and `*` can be implemented to correctly do rational arithmetic. We will then create a new type that can represent a dual number  $a + b\epsilon$ , where the type has two fields (`a` and `b`). By implementing basic arithmetic operations as well as more complicated functions like `exp` we can efficiently, and extremely accurately, compute derivatives of quite general functions.

In the problem sheet, we explore how dual numbers can also be used for pen-and-paper calculations of derivatives. This gives an alternative to traditional differentiation rules like chain and product rule, that while it is mathematically equivalent feels very different in practice. (I prefer it because it is much more algorithmic!) Make sure when doing the problem sheet to only use dual numbers and not fall back to the more traditional rules. We also see that one can extend the concept to a 2D-analogue of dual numbers, which allows for computation of gradients.

## I.4 Newton's method

In school you may recall learning Newton's method: a way of approximating zeros/roots to a function by using a local approximation by an affine function. That is, approximate a function  $f(x)$  locally around an initial guess  $x_0$  by its first order Taylor series:

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0)$$

and then find the root of the right-hand side which is

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

We can then repeat using this root as the new initial guess. In other words we have a sequence of *hopefully* more accurate approximations:

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

Thus if we can compute derivatives, we can (sometimes) compute roots.

In terms of analysis, we can guarantee convergence provided our initial guess is accurate enough. The first step is the bound the error of an iteration in terms of the previous error:

**Theorem 3** (Newton error). *Suppose  $f$  is twice-differentiable in a neighbourhood  $B$  of  $r$  such that  $f(r) = 0$ , and  $f'$  does not vanish in  $B$ . Denote the error of the  $k$ -th Newton iteration as  $\varepsilon_k := r - x_k$ . If  $x_k \in B$  then*

$$|\varepsilon_{k+1}| \leq M|\varepsilon_k|^2$$

where

$$M := \frac{1}{2} \sup_{x \in B} |f''(x)| \sup_{x \in B} \left| \frac{1}{f'(x)} \right|.$$

**Proof** Using Taylor's theorem we find that

$$0 = f(r) = f(x_k + \varepsilon_k) = f(x_k) + f'(x_k)\varepsilon_k + \frac{f''(t)}{2}\varepsilon_k^2.$$

for some  $t \in B$  between  $r$  and  $x_k$ . Rearranging this we get an expression for  $f(x_k)$  that tells us that

$$\varepsilon_{k+1} = r - \underbrace{x_{k+1}}_{x_k - f(x_k)/f'(x_k)} = \varepsilon_k + \frac{f(x_k)}{f'(x_k)} = -\frac{f''(t)}{2f'(x_k)}\varepsilon_k^2.$$

Taking absolute values of each side gives the result.

■

This result says that the error decays *quadratically*, which in this case means that the number of digits roughly doubles each iteration. That is, if the error at one step is about  $10^{-3}$  then the error at the next step is about  $10^{-6}$  and the step after about  $10^{-12}$ : this is a drastic improvement! Hidden in this result is a guarantee of convergence provided  $x_0$  is sufficiently close to  $r$ .

**Corollary 1** (Newton convergence). *If  $x_0 \in B$  is sufficiently close to  $r$  then  $x_k \rightarrow r$ .*

**Proof**

Suppose  $x_k \in B$  satisfies  $|\varepsilon_k| = |r - x_k| \leq M^{-1}$ . Then

$$|\varepsilon_{k+1}| \leq M|\varepsilon_k|^2 \leq |\varepsilon_k|,$$

hence  $x_{k+1} \in B$ . Thus from induction if  $x_0$  satisfies the condition  $|\varepsilon_0| < M^{-1}$  condition then  $x_k \in B$  for all  $k$  and satisfies  $|\varepsilon_k| \leq M^{-1}$ . Thus we find (for large enough  $k$ )

$$|\varepsilon_k| \leq M|\varepsilon_{k-1}|^2 \leq M^3|\varepsilon_{k-2}|^4 \leq M^7|\varepsilon_{k-3}|^8 \leq \dots \leq M^{2^k-1}|\varepsilon_0|^{2^k} = \frac{1}{M}(M|\varepsilon_0|)^{2^k}.$$

Provided  $x_0$  satisfies the strict inequality  $|\varepsilon_0| < M^{-1}$  this will go to zero as  $k \rightarrow \infty$ .

■

### I.4.1 Lab and problem sheet

In the lab we explore using Newton's method for some simple root finding problems. We also see that automatic differentiation via dual numbers can be used effectively to compute the derivatives. This is in some sense a baby version of how Machine Learning algorithms train neural networks; but whilst Newton uses derivatives (or in higher-dimensions, gradients) to find roots of functions Machine Learning uses gradients to (very roughly) minimise functions that represent the error between a neural network and training data. Minimisation problems are very closely related to root finding problems (essentially the minima are associated with roots of the gradient) and there are specialised training algorithms in ML built on a randomised version of Newton's method.

In the problem sheet we see how the error bound for Newton iteration can be extended to the degenerate case where the second derivative also vanishes, but now we no longer achieve quadratic convergence, but it still decays exponentially with the number of iterations (which is called *linear convergence*).

# Chapter II

## Representing Numbers

In this chapter we aim to answer the question: when can we rely on computations done on a computer? Why are some computations (differentiation via divided differences), extremely inaccurate whilst others (integration via rectangular rule) accurate up to about 16 digits? In order to address these questions we need to dig deeper and understand at a basic level what a computer is actually doing when manipulating numbers.

Before we begin it is important to have a basic model of how a computer works. Our simplified model of a computer will consist of a [Central Processing Unit \(CPU\)](#)—the brains of the computer—and [Memory](#)—where data is stored. Inside the CPU there are [registers](#), where data is temporarily stored after being loaded from memory, manipulated by the CPU, then stored back to memory. Memory is a sequence of bits: 1s and 0s, essentially “on/off” switches, and memory is *finite*. Finally, if one has a  $p$ -bit CPU (eg a 32-bit or 64-bit CPU), each register consists of exactly  $p$ -bits. Most likely  $p = 64$  on your machine.

Thus representing numbers on a computer must overcome three fundamental limitations:

1. CPUs can only manipulate data  $p$ -bits at a time.
2. Memory is finite (in particular at most  $2^p$  bytes).
3. There is no such thing as an “error”: if anything goes wrong in the computation we must use some of the  $p$ -bits to indicate this.

This is clearly problematic: there are an infinite number of integers and an uncountable number of reals! Each of which we need to store in precisely  $p$ -bits. Moreover, some operations are simply undefined, like division by 0. This chapter discusses the solution used to this problem, alongside the mathematical analysis that is needed to understand the implications, in particular, that computations have *error*.

In particular we discuss:

1. II.1 Reals: real numbers are approximated by floating point numbers, which are a computers version of scientific notation.
2. II.2 Floating Point Arithmetic: arithmetic with floating point numbers is exact up-to-rounding, which introduces small-but-understandable errors in the computations. We explain how these errors can be analysed mathematically to get rigorous bounds.

3. II.3 Interval Arithmetic: rounding can be controlled in order to implement *interval arithmetic*, a way to compute rigorous bounds for computations. In the lab, we use this to compute up to 15 digits of  $e \equiv \exp 1$  rigorously with precise bounds on the error.

## II.1 Reals

In this chapter, we introduce the [IEEE Standard for Floating-Point Arithmetic](#). There are multiple ways of representing real numbers on a computer, as well as the precise behaviour of operations such as addition, multiplication, etc. One can use

1. [Fixed-point arithmetic](#): essentially representing a real number as an integer where a decimal point is inserted at a fixed position. This turns out to be impractical in most applications, e.g., due to loss of relative accuracy for small numbers.
2. [Floating-point arithmetic](#): essentially scientific notation where an exponent is stored alongside a fixed number of digits. This is what is used in practice.
3. [Level-index arithmetic](#): stores numbers as iterated exponents. This is the most beautiful mathematically but unfortunately is not as useful for most applications and is not implemented in hardware.

Before the 1980s each processor had potentially a different representation for floating-point numbers, as well as different behaviour for operations. IEEE introduced in 1985 standardised this across processors so that algorithms would produce consistent and reliable results.

This chapter may seem very low level for a mathematics course but there are two important reasons to understand the behaviour of floating-point numbers in details:

1. Floating-point arithmetic is very precisely defined, and can even be used in rigorous computations as we shall see in the labs. But it is not exact and its important to understand how errors in computations can accumulate.
2. Failure to understand floating-point arithmetic can cause catastrophic issues in practice, with the extreme example being the [explosion of the Ariane 5 rocket](#).

### II.1.1 Real numbers in binary

Integers can be written in binary as follows:

**Definition 3** (binary format). For  $B_0, \dots, B_p \in \{0, 1\}$  denote an integer in *binary format* by:

$$\pm(B_p \dots B_1 B_0)_2 := \pm \sum_{k=0}^p B_k 2^k$$

Reals can also be presented in binary format, that is, a sequence of 0s and 1s alongside a decimal point:

**Definition 4** (real binary format). For  $b_1, b_2, \dots \in \{0, 1\}$ , Denote a non-negative real number in *binary format* by:

$$(B_p \dots B_0.b_1b_2b_3\dots)_2 := (B_p \dots B_0)_2 + \sum_{k=1}^{\infty} \frac{b_k}{2^k}.$$

**Example 3** (rational in binary). Consider the number  $1/3$ . In decimal recall that:

$$1/3 = 0.3333\dots = \sum_{k=1}^{\infty} \frac{3}{10^k}$$

We will see that in binary

$$1/3 = (0.010101\dots)_2 = \sum_{k=1}^{\infty} \frac{1}{2^{2k}}$$

Both results can be proven using the geometric series:

$$\sum_{k=0}^{\infty} z^k = \frac{1}{1-z}$$

provided  $|z| < 1$ . That is, with  $z = \frac{1}{4}$  we verify the binary expansion:

$$\sum_{k=1}^{\infty} \frac{1}{4^k} = \frac{1}{1 - 1/4} - 1 = \frac{1}{3}$$

A similar argument with  $z = 1/10$  shows the decimal case.

### II.1.2 Floating-point numbers

Floating-point numbers are a subset of real numbers that are representable using a fixed number of bits.

**Definition 5** (floating-point numbers). Given integers  $\sigma$  (the *exponential shift*),  $Q$  (the number of *exponent bits*) and  $S$  (the *precision*), define the set of *Floating-point numbers* by dividing into *normal*, *sub-normal*, and *special number* subsets:

$$F_{\sigma,Q,S} := F_{\sigma,Q,S}^{\text{normal}} \cup F_{\sigma,Q,S}^{\text{sub}} \cup F_{\sigma,Q,S}^{\text{special}}.$$

The *normal numbers*  $F_{\sigma,Q,S}^{\text{normal}} \subset \mathbb{R}$  are

$$F_{\sigma,Q,S}^{\text{normal}} := \{\pm 2^{\textcolor{green}{q}-\sigma} \times (1.\textcolor{blue}{b}_1\textcolor{blue}{b}_2\textcolor{blue}{b}_3\dots\textcolor{blue}{b}_S)_2 : 1 \leq q < 2^Q - 1\}.$$

The *sub-normal numbers*  $F_{\sigma,Q,S}^{\text{sub}} \subset \mathbb{R}$  are

$$F_{\sigma,Q,S}^{\text{sub}} := \{\pm 2^{\textcolor{red}{1}-\sigma} \times (0.\textcolor{blue}{b}_1\textcolor{blue}{b}_2\textcolor{blue}{b}_3\dots\textcolor{blue}{b}_S)_2\}.$$

The *special numbers*  $F_{\sigma,Q,S}^{\text{special}} \not\subset \mathbb{R}$  are

$$F_{\sigma,Q,S}^{\text{special}} := \{\infty, -\infty, \text{NaN}\}$$

where NaN is a special symbol representing “not a number”, essentially an error flag.

Note this set of real numbers has no nice *algebraic structure*: it is not closed under addition, subtraction, etc. On the other hand, we can control errors effectively hence it is extremely useful for analysis.

Floating-point numbers are stored in  $1 + Q + S$  total number of bits, in the format

$$\textcolor{red}{s} \textcolor{teal}{q}_{Q-1} \dots q_0 \textcolor{blue}{b}_1 \dots b_S$$

The first bit ( $s$ ) is the *sign bit*: 0 means positive and 1 means negative. The bits  $q_{Q-1} \dots q_0$  are the *exponent bits*: they are the binary digits of the unsigned integer  $q$ :

$$q = (\textcolor{teal}{q}_{Q-1} \dots q_0)_2.$$

Finally, the bits  $b_1 \dots b_S$  are the *significand bits*. If  $1 \leq q < 2^Q - 1$  then the bits represent the normal number

$$x = \pm 2^{\textcolor{teal}{q}-\sigma} \times (1.\textcolor{blue}{b}_1 b_2 b_3 \dots b_S)_2.$$

If  $q = 0$  (i.e. all bits are 0) then the bits represent the sub-normal number

$$x = \pm 2^{1-\sigma} \times (0.\textcolor{blue}{b}_1 b_2 b_3 \dots b_S)_2.$$

If  $q = 2^Q - 1$  (i.e. all bits are 1) then the bits represent a special number. If all significand bits are 0 then it represents  $\pm\infty$ . Otherwise if any significand bit is 1 then it represents NaN.

### II.1.3 IEEE floating-point numbers

**Definition 6** (IEEE floating-point numbers). IEEE has 3 standard floating-point formats: 16-bit (half precision), 32-bit (single precision) and 64-bit (double precision) defined by (you *do not* need to memorise these):

$$F_{16} := F_{15,5,10}$$

$$F_{32} := F_{127,8,23}$$

$$F_{64} := F_{1023,11,52}$$

**Example 4** (interpreting 16-bits as a float). Consider the number with bits

$$\textcolor{red}{0} \textcolor{teal}{10000} \textcolor{blue}{1010000000}$$

assuming it is a half-precision float ( $F_{16}$ ). Since the sign bit is 0 it is positive. The exponent bits encode

$$q = (10000)_2 = 2^4$$

hence the exponent is

$$q - \sigma = 2^4 - 15 = 1$$

and the number is:

$$2^1(1.1010000000)_2 = 2(1 + 1/2 + 1/8) = 3 + 1/4 = 3.25.$$

**Example 5** (rational to 16-bits). How is the number  $1/3$  stored in  $F_{16}$ ? Recall that

$$1/3 = (0.010101\dots)_2 = 2^{-2}(1.0101\dots)_2 = 2^{13-15}(1.0101\dots)_2$$

and since  $13 = (1101)_2$  the exponent bits are 01101. For the significand we round the last bit to the nearest element of  $F_{16}$ , (the exact rule for rounding is explained in detail later), so we have

$$1.010101010101010101\dots \approx 1.01010101 \in F_{16}$$

and the significand bits are 0101010101. Thus the stored bits for  $1/3$  are:

$$\textcolor{red}{0} \textcolor{teal}{01101} \textcolor{blue}{0101010101}$$

### II.1.4 Sub-normal and special numbers

For sub-normal numbers, the simplest example is zero, which has  $q = 0$  and all significand bits zero: 0 00000 0000000000. Unlike integers, we also have a negative zero, which has bits: 1 00000 0000000000. This is treated as identical to positive 0 (except for degenerate operations as explained in the lab).

**Example 6** (subnormal in 16-bits). Consider the number with bits

$$\textcolor{red}{1} \textcolor{green}{00000} \textcolor{blue}{1100000000}$$

assuming it is a half-precision float ( $F_{16}$ ). Since all exponent bits are zero it is sub-normal. Since the sign bit is 1 it is negative. Hence this number is:

$$-2^{1-\sigma}(0.1100000000)_2 = -2^{-14}(2^{-1} + 2^{-2}) = -3 \times 2^{-16}$$

The special numbers extend the real line by adding  $\pm\infty$  but also a notion of “not-a-number” NaN. Whenever the bits of  $q$  of a floating-point number are all 1 then they represent an element of  $F^{\text{special}}$ . If all  $b_k = 0$ , then the number represents either  $\pm\infty$ . All other special floating-point numbers represent NaN.

**Example 7** (special in 16-bits). The number with bits

$$\textcolor{red}{1} \textcolor{green}{11111} \textcolor{blue}{0000000000}$$

has all exponent bits equal to 1, and significand bits 0 and sign bit 1, hence represents  $-\infty$ . On the other hand, the number with bits

$$\textcolor{red}{1} \textcolor{green}{11111} \textcolor{blue}{0000000001}$$

has all exponent bits equal to 1 but does not have all significand bits equal to 0, hence is one of many representations for NaN.

## II.2 Floating Point Arithmetic

Arithmetic operations on floating-point numbers are *exact up to rounding*. There are three basic rounding strategies: round up/down/nearest. Mathematically we introduce a function to capture the notion of rounding:

**Definition 7** (rounding). The function  $\text{fl}_{\sigma,Q,S}^{\text{up}} : \mathbb{R} \rightarrow F_{\sigma,Q,S}$  rounds a real number up to the nearest floating-point number that is greater or equal:

$$\text{fl}_{\sigma,Q,S}^{\text{up}}(x) := \min\{y \in F_{\sigma,Q,S} : y \geq x\}.$$

The function  $\text{fl}_{\sigma,Q,S}^{\text{down}} : \mathbb{R} \rightarrow F_{\sigma,Q,S}$  rounds a real number down to the nearest floating-point number that is less or equal:

$$\text{fl}_{\sigma,Q,S}^{\text{down}}(x) := \max\{y \in F_{\sigma,Q,S} : y \leq x\}.$$

The function  $\text{fl}_{\sigma,Q,S}^{\text{nearest}} : \mathbb{R} \rightarrow F_{\sigma,Q,S}$  denotes the function that rounds a real number to the nearest floating-point number. In case of a tie, it returns the floating-point number whose least significant bit is equal to zero. We use the notation fl when  $\sigma, Q, S$  and the rounding mode are implied by context, with  $\text{fl}^{\text{nearest}}$  being the default rounding mode.

In more detail on the behaviour of nearest mode, if a positive number  $x$  is between two normal floats  $x_- \leq x \leq x_+$  we can write its expansion as

$$x = 2^{q-\sigma}(1.b_1 b_2 \dots b_S b_{S+1} \dots)_2$$

where

$$\begin{aligned} x_- &:= \text{fl}^{\text{down}}(x) = 2^{q-\sigma}(1.b_1 b_2 \dots b_S)_2 \\ x_+ &:= \text{fl}^{\text{up}}(x) = x_- + 2^{q-\sigma-S} \end{aligned}$$

Write the half-way point as:

$$x_h := \frac{x_+ + x_-}{2} = x_- + 2^{q-\sigma-S-1} = 2^{q-\sigma}(1.b_1 b_2 \dots b_S 1)_2$$

If  $x_- \leq x < x_h$  then  $\text{fl}(x) = x_-$  and if  $x_h < x \leq x_+$  then  $\text{fl}(x) = x_+$ . If  $x = x_h$  then it is exactly half-way between  $x_-$  and  $x_+$ . The rule is if  $b_S = 0$  then  $\text{fl}(x) = x_-$  and otherwise  $\text{fl}(x) = x_+$ .

In IEEE arithmetic, the arithmetic operations  $+$ ,  $-$ ,  $*$ ,  $/$  are defined by the property that they are exact up to rounding. Mathematically we denote these operations as  $\oplus, \ominus, \otimes, \oslash : F_{\sigma,Q,S} \times F_{\sigma,Q,S} \rightarrow F_{\sigma,Q,S}$  as follows:

$$\begin{aligned} x \oplus y &:= \text{fl}(x + y) \\ x \ominus y &:= \text{fl}(x - y) \\ x \otimes y &:= \text{fl}(x * y) \\ x \oslash y &:= \text{fl}(x / y) \end{aligned}$$

Note also that  $\wedge$  and `sqrt` are similarly exact up to rounding. Also, note that when we convert a Julia command with constants specified by decimal expansions we first round the constants to floats, e.g., `1.1 + 0.1` is actually reduced to

$$\text{fl}(1.1) \oplus \text{fl}(0.1)$$

This includes the case where the constants are integers (which are normally exactly floats but may be rounded if extremely large).

**Example 8** (decimal is not exact). On a computer `1.1+0.1` is close to but not exactly the same thing as `1.2`. This is because  $\text{fl}(1.1) \neq 1 + 1/10$  and  $\text{fl}(0.1) \neq 1/10$  since their expansion in *binary* is not finite. For  $F_{16}$  we have:

$$\begin{aligned} \text{fl}(1.1) &= \text{fl}((1.00011001100111\dots)_2) = (1.0001100110)_2 \\ \text{fl}(0.1) &= \text{fl}(2^{-4}(1.10011001100111\dots)_2) = 2^{-4} * (1.1001100110)_2 = (0.00011001100110)_2 \end{aligned}$$

Thus when we add them we get

$$\text{fl}(1.1) + \text{fl}(0.1) = (1.00110011000111\dots)_2$$

where the red digits indicate those beyond the 10 significant digits representable in  $F_{16}$ . In this case we round down and get

$$\text{fl}(1.1) \oplus \text{fl}(0.1) = (1.0011001100)_2$$

On the other hand,

$$\text{fl}(1.2) = \text{fl}((1.001100110011001100\dots)_2) = (1.0011001101)_2$$

which differs by 1 bit.

**WARNING (non-associative)** These operations are not associative! E.g.  $(x \oplus y) \oplus z$  is not necessarily equal to  $x \oplus (y \oplus z)$ . Commutativity is preserved, at least.

### II.2.1 Bounding errors in floating point arithmetic

When dealing with normal numbers there are some important constants that we will use to bound errors.

**Definition 8** (machine epsilon/smallest positive normal number/largest normal number). *Machine epsilon* is denoted

$$\epsilon_{m,S} := 2^{-S}.$$

When  $S$  is implied by context we use the notation  $\epsilon_m$ . The *smallest positive normal number* is  $q = 1$  and  $b_k$  all zero:

$$\min |F_{\sigma,Q,S}^{\text{normal}}| = 2^{1-\sigma}$$

where  $|A| := \{|x| : x \in A\}$ . The *largest (positive) normal number* is

$$\max F_{\sigma,Q,S}^{\text{normal}} = 2^{2^Q-2-\sigma}(1.11\dots)_2 = 2^{2^Q-2-\sigma}(2 - \epsilon_m)$$

We can bound the error of basic arithmetic operations in terms of machine epsilon, provided a real number is close to a normal number:

**Definition 9** (normalised range). The *normalised range*  $\mathcal{N}_{\sigma,Q,S} \subset \mathbb{R}$  is the subset of real numbers that lies between the smallest and largest normal floating-point number:

$$\mathcal{N}_{\sigma,Q,S} := \{x : \min |F_{\sigma,Q,S}^{\text{normal}}| \leq |x| \leq \max |F_{\sigma,Q,S}^{\text{normal}}|\}$$

When  $\sigma, Q, S$  are implied by context we use the notation  $\mathcal{N}$ .

We can use machine epsilon to determine bounds on rounding:

**Proposition 2** (round bound). *If  $x \in \mathcal{N}$  then*

$$\text{fl}^{\text{mode}}(x) = x(1 + \delta_x^{\text{mode}})$$

where the relative error is bounded by:

$$\begin{aligned} |\delta_x^{\text{nearest}}| &\leq \frac{\epsilon_m}{2} \\ |\delta_x^{\text{up/down}}| &< \epsilon_m. \end{aligned}$$

#### Proof

We will show this result for the nearest rounding mode. Note first that

$$\text{fl}(-x) = -\text{fl}(x)$$

and hence it suffices to prove the result for positive  $x$ . Write

$$x = 2^{q-\sigma}(1.b_1b_2\dots b_Sb_{S+1}\dots)_2.$$

Define

$$\begin{aligned} x_- &:= \text{fl}^{\text{down}}(x) = 2^{q-\sigma}(1.b_1b_2\dots b_S)_2 \\ x_+ &:= \text{fl}^{\text{up}}(x) = x_- + 2^{q-\sigma-S} \\ x_h &:= \frac{x_+ + x_-}{2} = x_- + 2^{q-\sigma-S-1} = 2^{q-\sigma}(1.b_1b_2\dots b_S1)_2 \end{aligned}$$

so that  $x_- \leq x \leq x_+$ . We consider two cases separately.

**(Round Down)** First consider the case where  $x$  is such that we round down:  $\text{fl}(x) = x_-$ . Since  $2^{q-\sigma} \leq x_- \leq x \leq x_h$  we have

$$|\delta_x| = \frac{x - x_-}{x} \leq \frac{x_h - x_-}{x_-} \leq \frac{2^{q-\sigma-S-1}}{2^{q-\sigma}} = 2^{-S-1} = \frac{\epsilon_m}{2}.$$

**(Round Up)** If  $\text{fl}(x) = x_+$  then  $2^{q-\sigma} \leq x_- < x_h \leq x \leq x_+$  and hence

$$|\delta_x| = \frac{x_+ - x}{x} \leq \frac{x_+ - x_h}{x_-} \leq \frac{2^{q-\sigma-S-1}}{2^{q-\sigma}} = 2^{-S-1} = \frac{\epsilon_m}{2}.$$

■

This immediately implies relative error bounds on all IEEE arithmetic operations, e.g., if  $x + y \in \mathcal{N}$  then we have

$$x \oplus y = (x + y)(1 + \delta_1)$$

where (assuming the default nearest rounding)  $|\delta_1| \leq \frac{\epsilon_m}{2}$ .

### II.2.2 Idealised floating point

With a complicated formula it is mathematically inelegant to work with normalised ranges: one cannot guarantee apriori that a computation always results in a normal float. Extending the bounds to subnormal numbers is tedious, rarely relevant, and beyond the scope of this module. Thus to avoid this issue we will work with an alternative mathematical model:

**Definition 10** (idealised floating point). An idealised mathematical model of floating point numbers for which the only subnormal number is zero can be defined as:

$$F_{\infty,S} := \{\pm 2^q \times (1.b_1 b_2 b_3 \dots b_S)_2 : q \in \mathbb{Z}\} \cup \{0\}$$

Note that  $F_{\sigma,Q,S}^{\text{normal}} \subset F_{\infty,S}$  for all  $\sigma, Q \in \mathbb{N}$ . The definition of rounding  $\text{fl}_{\infty,S}^{\text{mode}} : \mathbb{R} \rightarrow F_{\infty,S}$  naturally extend to  $F_{\infty,S}$  and hence we can consider bounds for floating point operations such as  $\oplus, \ominus, \text{etc}$ . And in this model the round bound is valid for all real numbers (including  $x = 0$ ).

**Example 9** (bounding a simple computation). We show how to bound the error in computing  $(1.1 + 1.2) * 1.3 = 2.99$  and we may assume idealised floating-point arithmetic  $F_{\infty,S}$ . First note that  $1.1$  on a computer is in fact  $\text{fl}(1.1)$ , and we will always assume nearest rounding unless otherwise stated. Thus this computation becomes

$$(\text{fl}(1.1) \oplus \text{fl}(1.2)) \otimes \text{fl}(1.3)$$

We will show the *absolute error* is given by

$$(\text{fl}(1.1) \oplus \text{fl}(1.2)) \otimes \text{fl}(1.3) = 2.99 + \delta$$

where  $|\delta| \leq 23\epsilon_m$ . First we find

$$\begin{aligned} \text{fl}(1.1) \oplus \text{fl}(1.2) &= (1.1(1 + \delta_1) + 1.2(1 + \delta_2))(1 + \delta_3) \\ &= 2.3 + \underbrace{1.1\delta_1 + 1.2\delta_2 + 2.3\delta_3 + 1.1\delta_1\delta_3 + 1.2\delta_2\delta_3}_{\varepsilon_1}. \end{aligned}$$

While  $\delta_1\delta_3$  and  $\delta_2\delta_3$  are absolutely tiny in practice we will bound them rather naïvely by eg.

$$|\delta_1\delta_3| \leq \epsilon_m^2/4 \leq \epsilon_m/4.$$

Further we round up constants to integers in the bounds for simplicity. We thus have the bound

$$|\varepsilon_1| \leq (2 + 2 + 3 + 1 + 1) \frac{\epsilon_m}{2} \leq 5\epsilon_m.$$

Writing  $\text{fl}(1.3) = 1.3(1 + \delta_4)$  and also incorporating an error from the rounding in  $\otimes$  we arrive at

$$\begin{aligned} (\text{fl}(1.1) \oplus \text{fl}(1.2)) \otimes \text{fl}(1.3) &= (2.3 + \varepsilon_1)1.3(1 + \delta_4)(1 + \delta_5) \\ &= 2.99 + \underbrace{1.3(\varepsilon_1 + 2.3\delta_4 + 2.3\delta_5 + \varepsilon_1\delta_4 + \varepsilon_1\delta_5 + 2.3\delta_4\delta_5 + \varepsilon_1\delta_4\delta_5)}_{\delta} \end{aligned}$$

We use the bounds

$$\begin{aligned} |\varepsilon_1\delta_4|, |\varepsilon_1\delta_5| &\leq 5\epsilon_m^2/2 \leq 5\epsilon_m/2, \\ |\delta_4\delta_5| &\leq \epsilon_m^2/4 \leq \epsilon_m/4, \\ |\varepsilon_1\delta_4\delta_5| &\leq 5\epsilon_m^3/4 \leq 5\epsilon_m/4. \end{aligned}$$

Thus the *absolute error* is bounded (bounding 1.3 by 3/2) by

$$|\delta| \leq (3/2)(5 + 3/2 + 3/2 + 5/2 + 5/2 + 3/4 + 5/4)\epsilon_m \leq 23\epsilon_m.$$

### II.2.3 Divided differences floating point error bound

We can use the bound on floating point arithmetic to deduce a bound on divided differences that captures the phenomena we observed where the error of divided differences became large as  $h \rightarrow 0$ . We assume that the function we are attempting to differentiate is computed using floating point arithmetic in a way that has a small absolute error.

**Theorem 4** (divided difference error bound). *Assume we are working in idealised floating-point arithmetic  $F_{\infty,S}$ . Let  $f$  be twice-differentiable in a neighbourhood of  $x \in F_{\infty,S}$  and assume that*

$$f(x) = f^{\text{FP}}(x) + \delta_x^f$$

where  $f^{\text{FP}} : F_{S,\infty} \rightarrow F_{S,\infty}$  has uniform absolute accuracy in that neighbourhood, that is:

$$|\delta_x^f| \leq c\epsilon_m$$

for a fixed constant  $c \geq 0$ . The divided difference approximation partially implemented with floating point satisfies

$$\frac{f^{\text{FP}}(x+h) \ominus f^{\text{FP}}(x)}{h} = f'(x) + \delta_{x,h}^{\text{FD}}$$

where

$$|\delta_{x,h}^{\text{FD}}| \leq \frac{|f'(x)|}{2}\epsilon_m + Mh + \frac{4c\epsilon_m}{h}$$

for  $M = \sup_{x \leq t \leq x+h} |f''(t)|$ .

## Proof

We have

$$\begin{aligned}(f^{\text{FP}}(x+h) \ominus f^{\text{FP}}(x))/h &= \frac{f(x+h) - \delta_{x+h}^f - f(x) + \delta_x^f}{h}(1 + \delta_1) \\ &= \frac{f(x+h) - f(x)}{h}(1 + \delta_1) + \frac{\delta_x^f - \delta_{x+h}^f}{h}(1 + \delta_1)\end{aligned}$$

where  $|\delta_1| \leq \epsilon_m/2$ . Applying Taylor's theorem we get

$$(f^{\text{FP}}(x+h) \ominus f^{\text{FP}}(x))/h = f'(x) + f'(x)\delta_1 + \underbrace{\frac{f''(t)}{2}h(1 + \delta_1)}_{\delta_{x,h}^{\text{FD}}} + \frac{\delta_x^f - \delta_{x+h}^f}{h}(1 + \delta_1)$$

The bound then follows, using the very pessimistic bound  $|1 + \delta_1| \leq 2$ .

■

The previous theorem neglected some errors due to rounding, which was done for simplicity. This is justified under fairly general restrictions:

**Corollary 2** (divided differences in practice). *We have*

$$(f^{\text{FP}}(x \oplus h) \ominus f^{\text{FP}}(x)) \oslash h = \frac{f^{\text{FP}}(x+h) \ominus f^{\text{FP}}(x)}{h}$$

whenever  $h = 2^{j-n}$  for  $0 \leq n \leq S$  and the last binary place of  $x \in F_{\infty,S}$  is zero, that is  $x = \pm 2^j(1.b_1 \dots b_{S-1}0)_2$ .

## Proof

We first confirm  $x \oplus h = x + h$ . If  $b_S = 0$  the worst possible case is that we increase the exponent by one as we are just adding 1 to one of the digits  $b_1, \dots, b_S$ . This would cause us to lose the last digit. But if that is zero no error is incurred when we round.

Now write  $y := (f^{\text{FP}}(x \oplus h) \ominus f^{\text{FP}}(x)) = \pm 2^\nu(1.c_1 \dots c_S)_2 \in F_{\infty,S}$ . We have

$$y/h = \pm 2^{\nu+n-j}(1.c_1 \dots c_S)_2 \in F_{\infty,S} \Rightarrow y/h = y \oslash h.$$

■

The three-terms of this bound tell us a story: the first term is a fixed (small) error, the second term tends to zero as  $h \rightarrow 0$ , while the last term grows like  $\epsilon_m/h$  as  $h \rightarrow 0$ . Thus we observe convergence while the second term dominates, until the last term takes over. Of course, a bad upper bound is not the same as a proof that something grows, but it is a good indication of what happens *in general* and suffices to choose  $h$  so that these errors are balanced (and thus minimised). Since in general we do not have access to the constants  $c$  and  $M$  we employ the following heuristic to balance the two sources of errors:

**Heuristic (divided difference with floating-point step)** Choose  $h$  proportional to  $\sqrt{\epsilon_m}$  in divided differences so that  $Mh$  and  $\frac{4c\epsilon_m}{h}$  are (roughly) the same magnitude.

In the case of double precision  $\sqrt{\epsilon_m} \approx 1.5 \times 10^{-8}$ , which is close to when the observed error begins to increase in the examples we saw before.

**Remark** While divided differences is of debatable utility for computing derivatives, it is extremely effective in building methods for solving differential equations, as we shall see

later. It is also very useful as a “sanity check” if one wants something to compare with other numerical methods for differentiation.

**Remark** It is also possible to deduce an error bound for the rectangular rule showing that the error caused by round-off is on the order of  $n\epsilon_m$ , that is it does in fact grow but the error without round-off which was bounded by  $M/n$  will be substantially greater for all reasonable values of  $n$ .



# Appendix A

## Asymptotics and Computational Cost

We introduce Big-O, little-o and asymptotic notation and see how they can be used to describe computational cost.

### A.1 Asymptotics as $n \rightarrow \infty$

Big-O, little-o, and “asymptotic to” are used to describe behaviour of functions at infinity.

**Definition 11** (Big-O).

$$f(n) = O(\phi(n)) \quad (\text{as } n \rightarrow \infty)$$

means  $\left| \frac{f(n)}{\phi(n)} \right|$  is bounded for sufficiently large  $n$ . That is, there exist constants  $C$  and  $N_0$  such that, for all  $n \geq N_0$ ,  $\left| \frac{f(n)}{\phi(n)} \right| \leq C$ .

**Definition 12** (little-O).

$$f(n) = o(\phi(n)) \quad (\text{as } n \rightarrow \infty)$$

means  $\lim_{n \rightarrow \infty} \frac{f(n)}{\phi(n)} = 0$ .

**Definition 13** (asymptotic to).

$$f(n) \sim \phi(n) \quad (\text{as } n \rightarrow \infty)$$

means  $\lim_{n \rightarrow \infty} \frac{f(n)}{\phi(n)} = 1$ .

**Example 10** (asymptotics with  $n$ ). 1.

$$\frac{\cos n}{n^2 - 1} = O(n^{-2})$$

as

$$\left| \frac{\frac{\cos n}{n^2 - 1}}{n^{-2}} \right| \leq \left| \frac{n^2}{n^2 - 1} \right| \leq 2$$

for  $n \geq N_0 = 2$ .

2.

$$\log n = o(n)$$

as  $\lim_{n \rightarrow \infty} \frac{\log n}{n} = 0$ .

3.

$$n^2 + 1 \sim n^2$$

as  $\frac{n^2+1}{n^2} \rightarrow 1$ .

Note we sometimes write  $f(O(\phi(n)))$  for a function of the form  $f(g(n))$  such that  $g(n) = O(\phi(n))$ .

We have some simple algebraic rules:

**Proposition 3** (Big-O rules).

$$\begin{aligned} O(\phi(n))O(\psi(n)) &= O(\phi(n)\psi(n)) && (\text{as } n \rightarrow \infty) \\ O(\phi(n)) + O(\psi(n)) &= O(|\phi(n)| + |\psi(n)|) && (\text{as } n \rightarrow \infty). \end{aligned}$$

**Proof** See any standard book on asymptotics, eg [F.W.J. Olver, Asymptotics and Special Functions](#). ■

## A.2 Asymptotics as $x \rightarrow x_0$

We also have Big-O, little-o and "asymptotic to" at a point:

**Definition 14** (Big-O).

$$f(x) = O(\phi(x)) \quad (\text{as } x \rightarrow x_0)$$

means  $|\frac{f(x)}{\phi(x)}|$  is bounded in a neighbourhood of  $x_0$ . That is, there exist constants  $C$  and  $r$  such that, for all  $0 \leq |x - x_0| \leq r$ ,  $|\frac{f(x)}{\phi(x)}| \leq C$ .

**Definition 15** (little-O).

$$f(x) = o(\phi(x)) \quad (\text{as } x \rightarrow x_0)$$

means  $\lim_{x \rightarrow x_0} \frac{f(x)}{\phi(x)} = 0$ .

**Definition 16** (asymptotic to).

$$f(x) \sim \phi(x) \quad (\text{as } x \rightarrow x_0)$$

means  $\lim_{x \rightarrow x_0} \frac{f(x)}{\phi(x)} = 1$ .

**Example 11** (asymptotics with  $x$ ).

$$\exp x = 1 + x + O(x^2) \quad \text{as } x \rightarrow 0$$

since  $\exp x = 1 + x + \frac{\exp t}{2}x^2$  for some  $t \in [0, x]$  and

$$\left| \frac{\frac{\exp t}{2}x^2}{x^2} \right| \leq \frac{3}{2}$$

provided  $x \leq 1$ .

## A.3 Computational cost

We will use Big-O notation to describe the computational cost of algorithms. Consider the following simple sum

$$\sum_{k=1}^n x_k^2$$

which we might implement as:

```
function sumsq(x)
    n = length(x)
    ret = 0.0
    for k = 1:n
        ret = ret + x[k]^2
    end
    ret
end

sumsq (generic function with 1 method)
```

Each step of this algorithm consists of one memory look-up ( $z = x[k]$ ), one multiplication ( $w = z*z$ ) and one addition ( $ret = ret + w$ ). We will ignore the memory look-up in the following discussion. The number of CPU operations per step is therefore 2 (the addition and multiplication). Thus the total number of CPU operations is  $2n$ . But the constant 2 here is misleading: we didn't count the memory look-up, thus it is more sensible to just talk about the asymptotic complexity, that is, the *computational cost* is  $O(n)$ .

Now consider a double sum like:

$$\sum_{k=1}^n \sum_{j=1}^k x_j^2$$

which we might implement as:

```
function sumsq2(x)
    n = length(x)
    ret = 0.0
    for k = 1:n
        for j = 1:k
            ret = ret + x[j]^2
        end
    end
    ret
end

sumsq2 (generic function with 1 method)
```

Now the inner loop is  $O(1)$  operations (we don't try to count the precise number), which we do  $k$  times for  $O(k)$  operations as  $k \rightarrow \infty$ . The outer loop therefore takes

$$\sum_{k=1}^n O(k) = O\left(\sum_{k=1}^n k\right) = O\left(\frac{n(n+1)}{2}\right) = O(n^2)$$

operations.



# Appendix B

## Integers

In this appendix we discuss the following:

1. Unsigned integers: how computers represent non-negative integers using only  $p$ -bits, via [modular arithmetic](#).
2. Signed integers: how negative integers are handled using the [Two's-complement](#) format.

Mathematically, CPUs only act on  $p$ -bits at a time, with  $2^p$  possible sequences. That is, essentially all functions  $f$  are either of the form  $f : \mathbb{Z}_{2^p} \rightarrow \mathbb{Z}_{2^p}$  or  $f : \mathbb{Z}_{2^p} \times \mathbb{Z}_{2^p} \rightarrow \mathbb{Z}_{2^p}$ , where we use the following notation:

**Definition 17** (finite integers). Denote the set of the first  $m$  non-negative integers as  $\mathbb{Z}_m := \{0, 1, \dots, m - 1\}$ .

To translate between integers and bits we will need to write integers in binary format. That is, as sequence of 0s and 1s:

**Example 12** (integers in binary). A simple integer example is  $5 = 2^2 + 2^0 = (101)_2$ . On the other hand, we write  $-5 = -(101)_2$ . Another example is  $258 = 2^8 + 2 = (100000010)_2$ .

### B.0.1 Unsigned Integers

Computers represent integers by a finite number of  $p$ -bits, with  $2^p$  possible combinations of 0s and 1s. Denote these  $p$ -bits as  $B_{p-1} \dots B_1 B_0$  where  $B_k \in \{0, 1\}$ . For *unsigned integers* (non-negative integers) these bits dictate the first  $p$  binary digits:  $(B_{p-1} \dots B_1 B_0)_2$ . Integers represented with  $p$ -bits on a computer are interpreted as representing elements of  $\mathbb{Z}_{2^p}$  and integer arithmetic on a computer is equivalent to arithmetic modulo  $2^p$ . We denote modular arithmetic with  $m = 2^p$  as follows:

$$\begin{aligned}x \oplus_m y &:= (x + y) \pmod{m} \\x \ominus_m y &:= (x - y) \pmod{m} \\x \otimes_m y &:= (x * y) \pmod{m}\end{aligned}$$

When  $m$  is implied by context we just write  $\oplus, \ominus, \otimes$ . Note that the  $(\text{mod } m)$  function simply drops all bits except for the first  $p$ -bits when writing a number in binary.

**Example 13** (arithmetic with 8-bit unsigned integers). If the result of an operation lies between 0 and  $m = 2^8 = 256$  then arithmetic works exactly like standard integer arithmetic. For example,

$$\begin{aligned} 17 \oplus_{256} 3 &= 20 \pmod{256} = 20 \\ 17 \ominus_{256} 3 &= 14 \pmod{256} = 14 \end{aligned}$$

**Example 14** (overflow with 8-bit unsigned integers). If we go beyond the range the result “wraps around”. For example, with true integers we have

$$255 + 1 = (11111111)_2 + (00000001)_2 = (100000000)_2 = 256$$

However, the result is impossible to store in just 8-bits! So as mentioned instead it treats the integers as elements of  $\mathbb{Z}_{256}$  by dropping any extra digits:

$$255 \oplus_{256} 1 = 255 + 1 \pmod{256} = (100000000)_2 \pmod{256} = 0.$$

On the other hand, if we go below 0 we wrap around from above:

$$3 \ominus_{256} 5 = -2 \pmod{256} = 254 = (11111110)_2$$

**Example 15** (multiplication of 8-bit unsigned integers). Multiplication works similarly: for example,

$$254 \otimes_{256} 2 = 254 * 2 \pmod{256} = (11111110)_2 * 2 \pmod{256} = (111111100)_2 \pmod{256} = 252.$$

Note that multiplication by 2 is the same as shifting the binary digits left by one, just as multiplication by 10 shifts base-10 digits left by 1.

## B.0.2 Signed integer

Signed integers use the [Two's complement](#) convention. The convention is if the first bit is 1 then the number is negative: in this case if the bits had represented the unsigned integer  $2^p - y$  then they represent the signed integer  $-y$ . Thus for  $p = 8$  we are interpreting  $2^7$  through  $2^8 - 1$  as negative numbers. More precisely:

**Definition 18** (signed integers). Denote the finite signed integers as

$$\mathbb{Z}_{2^p}^s := \{-2^{p-1}, \dots, -1, 0, 1, \dots, 2^{p-1} - 1\}.$$

**Definition 19** (Shifted mod). Define for  $y = x \pmod{2^p}$

$$x \pmod{s}{2^p} := \begin{cases} y & 0 \leq y \leq 2^{p-1} - 1 \\ y - 2^p & 2^{p-1} \leq y \leq 2^p - 1 \end{cases}$$

Note that if  $R_p(x) = x \pmod{s}{2^p}$  then it can be viewed as a map  $R_p : \mathbb{Z} \rightarrow \mathbb{Z}_{2^p}^s$  or a one-to-one map  $R_p : \mathbb{Z}_{2^p} \rightarrow \mathbb{Z}_{2^p}^s$  whose inverse is  $R_p^{-1}(x) = x \pmod{2^p}$ . It can also be viewed as the identity map on signed integers  $R_p : \mathbb{Z}_{2^p}^s \rightarrow \mathbb{Z}_{2^p}^s$ , that is,  $R_p(x) = x$  if  $x \in \mathbb{Z}_{2^p}^s$ .

Arithmetic works precisely the same for signed and unsigned integers up to the mapping  $R_p$ , e.g. we have for  $m = 2^p$

$$\begin{aligned} x \oplus_m^s y &:= (x + y) \pmod{s}{m} \\ x \ominus_m^s y &:= (x - y) \pmod{s}{m} \\ x \otimes_m^s y &:= (x * y) \pmod{s}{m} \end{aligned}$$

**Example 16** (addition of 8-bit signed integers). Consider  $(-1) + 1$  in 8-bit arithmetic:

$$-1 \oplus_{256}^s 1 = -1 + 1 \pmod{s} 256 = 0$$

On the bit level this computation is exactly the same as unsigned integers. We represent the number  $-1$  using the same bits as the unsigned integer  $2^8 - 1 = 255$ , that is using the bits **11111111** (i.e., we store it equivalently to  $(11111111)_2 = 255$ ) and the number  $1$  is stored using the bits **00000001**. When we add this with true integer arithmetic we have

$$\begin{aligned} & (01111111)_2 + \\ & (00000001)_2 = \\ & (10000000)_2 \end{aligned}$$

Modular arithmetic drops the leading 1 and we are left with all zeros.

**Example 17** (signed overflow with 8-bit signed integers). If we go above  $2^{p-1}-1 = 2^7-1 = 127$  we have perhaps unexpected results:

$$127 \oplus_{256}^s 1 = 128 \pmod{s} 256 = 128 - 256 = -128.$$

Again on the bit level this computation is exactly the same as unsigned integers. We represent the number  $127$  using the bits **01111111** and the number  $1$  is stored using the bits **00000001**. When we add this with true integer arithmetic we have

$$\begin{aligned} & (01111111)_2 + \\ & (00000001)_2 = \\ & (10000000)_2 \end{aligned}$$

Because the first bit is **1** we interpret this as a negative number using the formula:

$$(10000000)_2 \pmod{s} 256 = 128 \pmod{s} 256 = -128.$$

**Example 18** (multiplication of 8-bit signed integers). Consider computation of  $(-2) * 2$ :

$$(-2) \otimes_{2^p}^s 2 = -4 \pmod{s} 2^p = -4$$

On the bit level, the bits of  $-2$  (which is one less than  $-1$ ) are **11111110**. Multiplying by 2 is like multiplying by 10 in base-10, that is, we shift the bits. Hence in true arithmetic we have

$$\begin{aligned} & (01111110)_2 * 2 = \\ & (111111100)_2 \end{aligned}$$

We drop the leading 1 due to modular arithmetic. We still have a leading 1 hence the number is viewed as negative. In particular we have

$$\begin{aligned} (111111100)_2 \pmod{s} 256 &= (11111100)_2 \pmod{s} 256 = 2^7 + 2^6 + 2^5 + 2^4 + 2^3 + 2^2 \pmod{s} 256 \\ &= 252 \pmod{s} 256 = -4. \end{aligned}$$

### B.0.3 Hexadecimal format

In coding it is often convenient to use base-16 as it is a power of 2 but uses less characters than binary. The digits used are 0 through 9 followed by  $a = 10$ ,  $b = 11$ ,  $c = 12$ ,  $d = 13$ ,  $e = 14$ , and  $f = 15$ .

**Example 19** (Hexadecimal number). We can interpret a number in format as follows:

$$(a5f2)_{16} = a * 16^3 + 5 * 16^2 + f * 16 + 2 = 10 * 16^3 + 5 * 16^2 + 15 * 16 + 2 = 42,482$$

We will see in the labs that unsigned integers are displayed in base-16.