

# Chapter 1

## Floating point arithmetic

### Introduction

When we study numerical algorithms in the next chapters, we assume implicitly that the operations involved are performed exactly. On a computer, however, only a subset of the real numbers can be stored and, consequently, many arithmetic operations are performed only approximately. This is the source of the so-called *round-off errors*. The rest of this chapter is organized as follows.

- In [Section 1.1](#), we discuss the binary representation of real numbers.
- In [Section 1.2](#), we describe the set of floating point numbers that can be represented in the usual floating point formats;
- In [Section 1.3](#) we explain how arithmetic operations between floating point numbers behave. We insist in particular on the fact that, in a calculation involving several successive arithmetic operations, the result of each intermediate operation is stored as a floating point number, with a possible error.
- In [Section 1.4](#), we briefly present how floating point numbers are encoded according to the IEEE 754 standard, which is widely adopted today. We discuss also the encoding of special values such as `Inf`, `-Inf` and `NaN`.
- Finally, in [Section 1.5](#), we present the standard integer formats and their encoding.

In order to completely describe computer arithmetic, one would in principle need to also discuss the conversion mechanisms between different number formats, as well as a number of edge cases. A comprehensive discussion of the subject is beyond the scope of this course; our aim in this chapter is only to introduce the key concepts.

## 1.1 Binary representation of real numbers

Given any integer number  $\beta > 0$ , called the *base*, a real number  $x$  can always be expressed as a finite or infinite series of the form

$$x = \pm \sum_{k=-n}^{\infty} a_k \beta^{-k}, \quad a_k \in \{0, \dots, \beta - 1\}. \quad (1.1)$$

The number  $x$  may then be denoted as  $\pm(a_{-n}a_{-n+1} \dots a_{-1}a_0.a_1a_2 \dots)_\beta$ , where the subscript  $\beta$  indicates the base. This numeral system is called the *positional notation* and is universally used today, both by humans (usually with  $\beta = 10$ ) and machines (usually with  $\beta = 2$ ). If the base  $\beta$  is omitted, it is always assumed in this course that  $\beta = 10$  unless otherwise specified – this is the *decimal* representation. The *digits*  $a_{-n}, a_{-n+1}, \dots$  are also called *bits* if  $\beta = 2$ . In computer science, several bases other than 10 are regularly employed, for example the following:

- Base 2 (binary) is the usual choice for storing numbers on a machine. The binary format is convenient because the digits have only two possible values, 0 or 1, and so they can be stored using simple electrical circuits with two states. We employ the binary notation extensively in the rest of this chapter. Notice that, just like multiplying and dividing by 10 is easy in base 10, multiplying and dividing by 2 is very simple in base 2: these operations amount to shifting all the bits by one position to the left or right, respectively.
- Base 16 (hexadecimal) is sometimes convenient to represent numbers in a compact manner. In order to represent the values 0-15 with a single digit, 16 different symbols are required, which are conventionally denoted by  $\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F\}$ . With this notation, we have  $(FF)_{16} = (255)_{10}$ , for example.

The hexadecimal notation is often used in programming languages for describing colors specified by a triplet  $(r, g, b)$  of values between 0 and 255, corresponding to the primary colors *red*, *green* and *blue*. The number of possible values for each component is  $256 = 16^2$ , and so only 2 digits are required to represent these in the hexadecimal notation. Hexadecimal numbers are also employed in IPv6 addresses, which are used to identify computers connected to a network.

### 1.1.1 Conversion between binary and decimal formats

Obtaining the decimal representation of a binary number can be achieved from (1.1), using the decimal representations of the powers of 2. Since all the positive and negative powers of 2 have a finite decimal representation, any real number with a finite representation in base 2 has a finite representation also in base 10. For example,  $(0.01)_2 = (0.25)_{10}$  and  $(0.111)_2 = (0.875)_{10}$ .

*Example 1.1* (Converting a binary number to decimal notation). Let us calculate the decimal representation of  $x = (0.\overline{10})_2$ , where the horizontal line indicates repetition:  $x =$

$(0.101010\dots)_2$ . By definition, it holds that

$$x = \sum_{k=0}^{\infty} a_k 2^{-k},$$

where  $a_k = 0$  if  $k$  is even and 1 otherwise. Thus, the series may be rewritten as

$$x = \sum_{k=0}^{\infty} 2^{-(2k+1)} = \frac{1}{2} \sum_{k=0}^{\infty} (2^{-2})^k.$$

We recognize on the right-hand side a geometric series with common ratio  $r = 2^{-2} = \frac{1}{4}$ , and so we obtain

$$x = \frac{1}{2} \left( \frac{1}{1-r} \right) = \frac{2}{3} = (0.\overline{6})_{10}.$$

Obtaining the binary representation of a decimal number is more difficult, because negative powers of 10 have *infinite* binary representations, as [Exercise 1.4](#) demonstrates. There is, however, a simple procedure to perform the conversion, which we present for the specific case of a real number  $x$  with decimal representation of the form  $x = (0.a_1 \dots a_n)_{10}$ . In this setting, the bits  $(b_1, b_2, \dots)$  in the binary representation of  $x = (0.b_1 b_2 b_2 \dots)_2$  may be obtained as follows:

---

**Algorithm 1** Conversion of a number to binary format

---

```

1:  $i \leftarrow 1$ 
2: while  $x \neq 0$  do
3:    $x \leftarrow 2x$ 
4:   if  $x \geq 1$  then
5:      $b_i \leftarrow 1$ 
6:   else
7:      $b_i \leftarrow 0$ 
8:   end if
9:    $x \leftarrow x - b_i$ 
10:   $i \leftarrow i + 1$ 
11: end while

```

---

*Example 1.2* (Converting a decimal number to binary notation). Let us calculate the binary representation of  $x = \frac{1}{3} = (0.\overline{3})_{10}$ . We apply [Algorithm 1](#) and collate the values of  $i$  and  $x$  obtained at the beginning of each iteration, i.e. just before [Line 3](#), in the table below.

$i$	$x$	Result
1	$\frac{1}{3}$	0.0000...
2	$\frac{2}{3}$	0.0100...
3	$\frac{1}{3}$	0.0000...

Since  $x$  in the last row is again  $\frac{1}{3}$ , successive bits alternate between 0 and 1, and the binary representation of  $x$  is given by  $(0.\overline{01})_2$ . This is not surprising since  $2x = (0.66)_{10} = (0.\overline{10})_2$ , as we saw in [Example 1.1](#).

## 1.2 Set of values representable in floating point formats

We mentioned in the introduction that, because of memory limitations, only a subset of the real numbers can be stored exactly in a computer. Nowadays, the vast majority of programming languages and software comply with the IEEE 754 standard, which requires that the set of representable numbers be of the form

$$\mathbf{F}(p, E_{\min}, E_{\max}) = \left\{ (-1)^s 2^E (b_0.b_1b_2 \dots b_{p-1})_2 : \right. \\ \left. s \in \{0, 1\}, b_i \in \{0, 1\} \text{ and } E_{\min} \leq E \leq E_{\max} \right\}. \quad (1.2)$$

In addition to these, floating number formats provide the special entities `Inf`, `-Inf` and `NaN`, the latter being an abbreviation for *Not a Number*. Three parameters appear in the set definition (1.2). The parameter  $p \in \mathbf{N}_{>0}$  is the number of significant bits (also called the precision), and  $(E_{\min}, E_{\max}) \in \mathbf{Z}^2$  are respectively the minimum and maximum exponents. From the precision, the *machine epsilon* is defined as  $\varepsilon_M = 2^{-(p-1)}$ ; its significance is discussed in Section 1.2.2.

For a number  $x \in \mathbf{F}(p, E_{\min}, E_{\max})$ ,  $s$  is called the *sign*,  $E$  is the *exponent* and  $b_0.b_1b_2 \dots b_{p-1}$  is the *significand*. The latter can be divided into a *leading bit*  $b_0$  and the *fraction*  $b_1b_2 \dots b_{p-1}$ , to the right of the binary point. The most widely used floating point formats are the *single* and *double precision* formats, which are called respectively **Float32** and **Float64** in Julia. Their parameters, together with those of the lesser-known half-precision format, are summarized in Table 1.1. In the rest of this section we use the shorthand notation  $\mathbf{F}_{16}$ ,  $\mathbf{F}_{32}$  and  $\mathbf{F}_{64}$ . Note that  $\mathbf{F}_{16} \subset \mathbf{F}_{32} \subset \mathbf{F}_{64}$ .

	Half precision	Single precision	Double precision
$p$	11	24	53
$E_{\min}$	-14	-126	-1022
$E_{\max}$	15	127	1023

Table 1.1: Floating point formats. The first column corresponds to the *half-precision* format. This format, which is available through the **Float16** type in Julia, is more recent than the single and double precision formats. It was introduced in the 2008 revision to the IEEE 754 standard of 1985, a revision known as IEEE 754-2008.

*Remark 1.1.* Some definitions, notably that in [10, Section 2.5.2], include a general base  $\beta$  instead of the base 2 as an additional parameter in the definition of the number format (1.2). Since the binary format ( $\beta = 2$ ) is always employed in practice, we focus on this case for simplicity in most of this chapter.

*Remark 1.2.* Given a real number  $x \in \mathbf{F}(p, E_{\min}, E_{\max})$ , the exponent  $E$  and significand are generally not uniquely defined. For example, the number  $2.0 \in \mathbf{F}_{64}$  may be expressed as  $(-1)^0 2^1 (1.00 \dots 00)_2$  or, equivalently, as  $(-1)^0 2^2 (0.100 \dots 00)_2$ .

In Julia, non-integer number literals are interpreted as **Float64** by default, which can be verified by using the `typeof` function. For example, the instruction “`a = 0.1`” is equivalent

to “ $a = \text{Float64}(0.1)$ ”. In order to define a number of type **Float32**, the suffix `f0` must be appended to the decimal expansion. For instance, the instruction “ $a = 4.0\text{f0}$ ” defines a floating point number  $a$  of type **Float32**; it is equivalent to writing “ $a = \text{Float32}(4.0)$ ”.

### 1.2.1 Denormalized floating point numbers

We can decompose the set  $\mathbf{F}(p, E_{\min}, E_{\max})$  in two disjoint parts:

$$\begin{aligned} \mathbf{F}(p, E_{\min}, E_{\max}) = & \left\{ (-1)^s 2^E (1.b_1 b_2 \dots b_{p-1})_2 : \right. \\ & \left. s \in \{0, 1\}, b_i \in \{0, 1\} \text{ and } E_{\min} \leq E \leq E_{\max} \right\} \\ & \cup \left\{ (-1)^s 2^{E_{\min}} (0.b_1 b_2 \dots b_{p-1})_2 : s \in \{0, 1\}, b_i \in \{0, 1\} \right\}. \end{aligned}$$

The numbers in the second set are called *subnormal* or *denormalized*.

### 1.2.2 Relative error and machine epsilon

Let  $x$  be a nonzero real number and  $\hat{x}$  be an approximation. We define the absolute and relative errors of the approximation as follows.

**Definition 1.1** (Absolute and relative error). The absolute error is given by  $|x - \hat{x}|$ , whereas the relative error is

$$\frac{|x - \hat{x}|}{|x|}$$

The following result establishes a link between the machine  $\varepsilon_M$  and the relative error between a real number and the closest member of a floating point format.

**Proposition 1.1.** Let  $x_{\min}$  and  $x_{\max}$  denote the smallest and largest non-denormalized positive numbers in a format  $F = \mathbf{F}(p, E_{\min}, E_{\max})$ . If  $x \in [-x_{\max}, -x_{\min}] \cup [x_{\min}, x_{\max}]$ , then

$$\min_{\hat{x} \in F} \frac{|x - \hat{x}|}{|x|} \leq \frac{1}{2} 2^{-(p-1)} = \frac{1}{2} \varepsilon_M. \quad (1.3)$$

*Proof.* For simplicity, we assume that  $x > 0$ . Let  $n = \lfloor \log_2(x) \rfloor$  and  $y := 2^{-n}x$ . Since  $y \in [1, 2)$ , it admits a binary representation of the form  $(1.b_1 b_2 \dots)_2 \neq (1.\bar{1})_2$ . Thus  $x = 2^n (1.b_1 b_2 \dots)_2$ , and from the assumption that  $x_{\min} \leq x \leq x_{\max}$  we deduce that  $E_{\min} \leq n \leq E_{\max}$ . We now define the number  $x_- \in F$  by truncating the binary expansion of  $x$  as follows:

$$x_- = 2^n (1.b_1 \dots b_{p-1})_2.$$

The distance between  $x_-$  and its successor in  $F$ , which we denote by  $x_+$ , is given by  $2^{n-p+1}$ . Consequently, it holds that

$$(x_+ - x) + (x - x_-) = x_+ - x_- = 2^{n-p+1}.$$

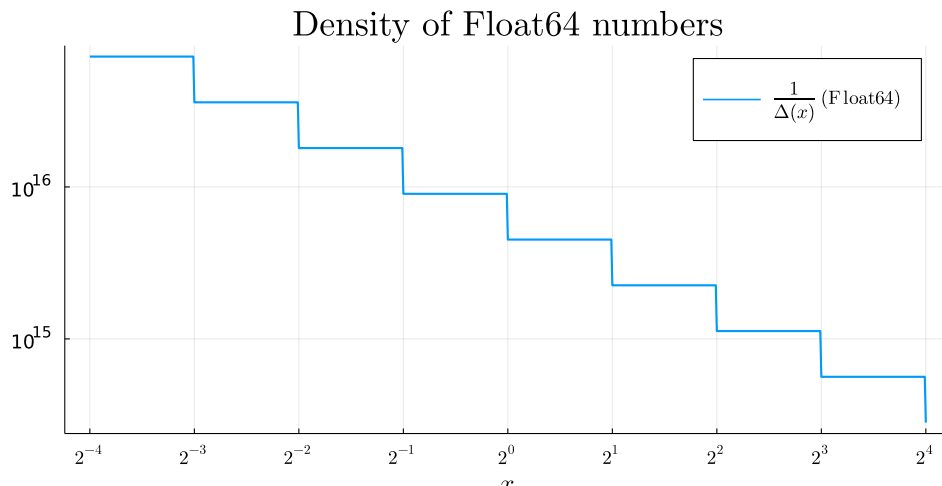


Figure 1.1: Density of the double-precision floating point numbers, measured here as  $1/\Delta(x)$  where, for  $x \in \mathbf{F}_{64}$ ,  $\Delta(x)$  denotes the distance between  $x$  and its successor in  $\mathbf{F}_{64}$ .

Since both summands on the left-hand side are positive, this implies that either  $x_+ - x$  or  $x - x_-$  is bounded from above by  $\frac{1}{2}2^{n-p+1} \leq \frac{1}{2}2^{-p+1}x$ , which concludes the proof.  $\square$

The machine epsilon, which was defined as  $\varepsilon_M = 2^{-(p-1)}$ , coincides with the maximum relative spacing between a non-denormalized floating point number  $x$  and its successor in the floating point format, defined as the smallest number in the format that is strictly larger than  $x$ .

Figure 1.1 depicts the density of double-precision floating point numbers, i.e. the number of  $\mathbf{F}_{64}$  members per unit on the real line. The figure shows that the density decreases as the absolute value of  $x$  increases. We also notice that the density is piecewise constant with discontinuities at powers of 2. Figure 1.2 illustrates the relative spacing between successive floating point numbers. Although the absolute spacing increases with the absolute value of  $x$ , the relative spacing oscillates between  $\frac{1}{2}\varepsilon_M$  and  $\varepsilon_M$ .

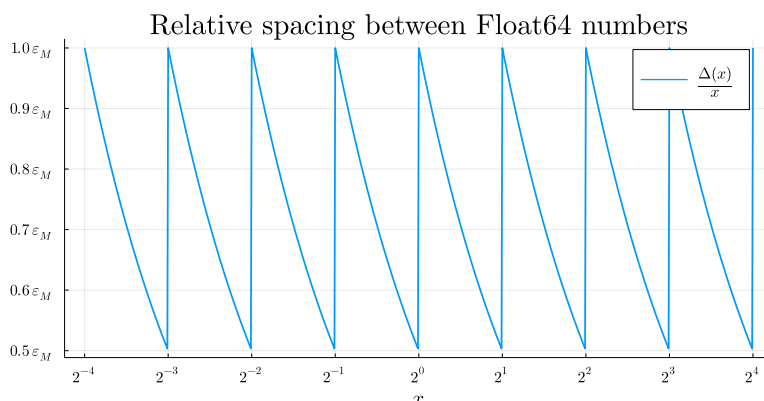


Figure 1.2: Relative spacing between successive double-precision floating point numbers in the “normal range”. The relative spacing oscillates between  $\frac{1}{2}\varepsilon_M$  and  $\varepsilon_M$ .

The picture of the relative spacing between successive floating point numbers looks quite different for denormalized numbers. This is illustrated in Figure 1.3, which shows that the relative spacing increases beyond the machine epsilon in the denormalized range. Fortunately, in the usual  $\mathbf{F}_{32}$  and  $\mathbf{F}_{64}$  formats, the transition between denormalized and non-denormalized

numbers occurs at such a small value that it rarely needs worrying about.

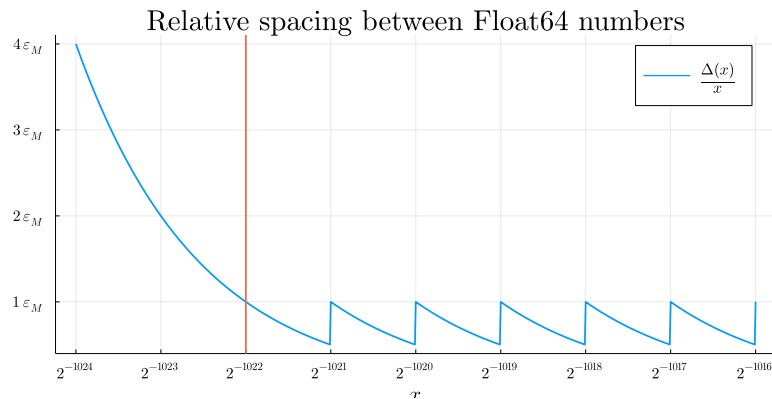


Figure 1.3: Relative spacing between successive double-precision floating point numbers, over a range which includes denormalized number. The vertical red line indicates the transition from denormalized to non-denormalized numbers.

*Example 1.3.* In Julia, the machine epsilon can be obtained using the `eps` function. For example, the instruction `eps(Float16)` returns  $\varepsilon_M$  for the half-precision format.

### 1.3 Arithmetic operations between floating point formats

Now that we have presented the set of values representable on a computer, we attempt in this section to understand precisely how arithmetic operations between floating point formats are performed. The key mechanism governing arithmetic operations on a computer is that of *rounding*, the action of approximating a real number regarded as infinitely precise by a number in a floating point format  $\mathbf{F}(p, E_{\min}, E_{\max})$ . The IEEE 754 standard stipulates that the default mechanism for rounding a real number  $x$ , called *round to nearest*, should behave as follows:

- **Standard case:** The number  $x$  is rounded to the *nearest representable number*, if this number is unique.
- **Edge case:** When there are two equally near representable numbers in the floating point format, the one with the least significant bit equal to zero is delivered.
- **Infinities:** If the real number  $x$  is larger than the largest representable number in the format, that is larger than or equal to  $x_{\max} = 2^{E_{\max}}(2 - \varepsilon)$ , then there are two cases,
  - If  $x < 2^{E_{\max}}(2 - 2^{-p})$ , then  $x_{\max}$  is delivered;
  - Otherwise, the special value `Inf` is delivered.

In other words,  $x_{\max}$  is delivered if it would be delivered by following the rules of the first two bullet points in a different floating point format with the same precision but a larger exponent  $E_{\max}$ . A similar rule applies for large negative numbers.

When a binary arithmetic operation ( $+$ ,  $-$ ,  $\times$ ,  $/$ ) is performed on floating point numbers in format  $\mathbf{F}$ , the result delivered by the computer is obtained by rounding the exact result of

the operation according to the rules given above. In other words, the arithmetic operation is performed as if the computer first calculated an intermediate exact result, and then rounded this intermediate result in order to provide a final result in  $\mathbf{F}$ .

Mathematically, arithmetic operations between floating point numbers in a given format  $\mathbf{F}$  may be formalized by introducing the rounding operator  $\text{fl} : \mathbf{R} \rightarrow \mathbf{F}$  and by defining, for any binary operation  $\circ \in \{+, -, \times, /\}$ , the corresponding machine operation

$$\hat{\circ} : \mathbf{F} \times \mathbf{F} \rightarrow \mathbf{F}; (x, y) \mapsto \text{fl}(x \circ y).$$

We defined this operator for arguments in the same floating point format  $\mathbf{F}$ . If the arguments of a binary arithmetic operation are of different types, the format of the end result, known as the *destination format*, depends on that of the arguments: as a rule of thumb, it is given by the most precise among the formats of the arguments. In addition, recall that a floating point literal whose format is not explicitly specified is rounded to double-precision format and so, for example, the addition  $0.1 + 0.1$  produces the result  $\text{fl}_{64}(\text{fl}_{64}(0.1) + \text{fl}_{64}(0.1))$ , where  $\text{fl}_{64}$  is the rounding operator to the double-precision format.

*Example 1.4.* Using the `typeof` function, we check that the floating point literal `1.0` is indeed interpreted as a double-precision number:

```
julia> a = 1.0; typeof(a)
Float64
```

When two numbers in different floating point formats are passed to a binary operation, the result is in the more precise format.

```
julia> typeof(Float16(1) + Float32(1))
Float32
```

```
julia> typeof(Float32(1) + Float64(1))
Float64
```

If a mathematical expression contains several binary arithmetic operations to be performed in succession, the result of each intermediate calculation is stored in a floating point format dictated by the formats of its argument, and this floating point number is employed in the next binary operation. A consequence of this mechanism is that the machine operands  $\hat{+}$  and  $\hat{*}$  are generally *not associative*. For example, in general

$$(x \hat{+} y) \hat{+} z \neq x \hat{+} (y \hat{+} z)$$

*Example 1.5.* Let  $x = 1$  and  $y = 3 \times 2^{-13}$ . Both of these numbers belong to  $\mathbf{F}_{16}$  and, denoting



by  $\hat{+}$  machine addition in  $\mathbf{F}_{16}$ , we have

$$(x \hat{+} y) \hat{+} y = 1 \quad (1.4)$$

but

$$x \hat{+} (y \hat{+} y) = 1 + 2^{-10}. \quad (1.5)$$

To explain this somewhat surprising result, we begin by writing the normalized representations of  $x$  and  $y$  in the  $\mathbf{F}_{16}$  format:

$$\begin{aligned} x &= (-1)^0 \times 2^0 \times (1.0000000000)_2 \\ y &= (-1)^0 \times 2^{-12} \times (1.1000000000)_2. \end{aligned}$$

The exact result of the addition  $x + y$  is given by  $r = 1 + 3 \times 2^{-13}$ , which in binary notation is

$$r = (1.\underbrace{0000000000}_{11 \text{ zeros}}11)_2.$$

Since the length of the significand in the half-precision ( $\mathbf{F}_{16}$ ) format is only  $p = 11$ , this number is not part of  $\mathbf{F}_{16}$ . The result of the machine addition  $\hat{+}$  is therefore obtained by rounding  $r$  to the nearest member of  $\mathbf{F}_{16}$ , which is 1. This reasoning can then be repeated in order to conclude that, indeed,

$$(x \hat{+} y) \hat{+} y = x \hat{+} y = 1.$$

In order to explain the result of (1.5), note that the exact result of the addition  $y + y$  is  $r = 3 \times 2^{-12}$ , which belongs to the floating point format, so it also holds that  $y \hat{+} y = 3 \times 2^{-12}$ . Therefore,

$$x \hat{+} (y \hat{+} y) = 1 \hat{+} 3 \times 2^{-12} = \text{fl}_{16}(1 + 3 \times 2^{-12}).$$

The argument of the  $\mathbf{F}_{16}$  rounding operator does not belong to  $\mathbf{F}_{16}$ , since its binary representation is given by

$$(1.\underbrace{0000000000}_{10 \text{ zeros}}11)_2.$$

This time the nearest member of  $\mathbf{F}_{16}$  is given by  $1 + 2^{-10}$ .

When a numerical computation unexpectedly returns **Inf** or **-Inf**, we say that an *overflow error* occurred. Similarly, *underflow* is said to occur when a number is smaller than the smallest representable number in a floating point format.

## 1.4 Encoding of floating point numbers

Once a number format is specified through parameters  $(p, E_{\min}, E_{\max})$ , the choice of encoding, i.e. the machine representation of numbers in this format, has no bearing on the magnitude and propagation of round-off errors. Studying encoding is, therefore, not essential for our purposes in this course, but we opted to cover the topic anyway in the hope that it will help the students

build intuition on floating point numbers. We focus mainly on the single precision format, but the following discussion applies *mutatis mutandis* to the double and half-precision formats. The material in this section is for information purposes only.

We already mentioned in [Remark 1.2](#) that a number in a floating point format may have several representations. On a computer, however, a floating point number is always stored in the same manner (except for the number 0, see [Remark 1.4](#)). The values of the exponent and significand which are selected by the computer, in the case where there are several possible choices, are determined from the following rules:

- Either  $E > E_{\min}$  and  $b_0 = 1$ ;
- Or  $E = E_{\min}$ , in which case the leading bit may be 0.

The following result proves that these rules enable to define the exponent and significand of a number in a set of floating point numbers uniquely.

**Proposition 1.2.** *Assume that*

$$(-1)^s(2^E b_0 . b_1 \dots b_{p-1})_2 = (-1)^{\tilde{s}}(2^{\tilde{E}} \tilde{b}_0 . \tilde{b}_1 \dots \tilde{b}_{p-1})_2, \quad (1.6)$$

where the parameter sets  $(s, E, b_0, \dots, b_{p-1})$  and  $(\tilde{s}, \tilde{E}, \tilde{b}_0, \dots, \tilde{b}_{p-1})$  both satisfy the above rule. Then  $E = \tilde{E}$  and  $b_i = \tilde{b}_i$  for  $i \in \{0, \dots, p-1\}$ .

*Proof.* We show that  $E = \tilde{E}$ , after which the equality of significands follows trivially. Let us assume for contradiction that  $E > \tilde{E}$  and denote the left and right-hand sides of (1.6) by  $x$  and  $\tilde{x}$ , respectively. Then  $E > E_{\min}$ , implying that  $b_0 = 1$  and so  $2^E \leq |x| < 2^{E+1}$ . On the other hand, it holds that  $|\tilde{x}| < 2^{\tilde{E}+1}$  regardless of whether  $\tilde{E} = E_{\min}$  or not. Since  $E \geq \tilde{E} + 1$  by assumption, we deduce that  $|\tilde{x}| < 2^E \leq |x|$ , which contradicts the equality  $x = \tilde{x}$ .  $\square$

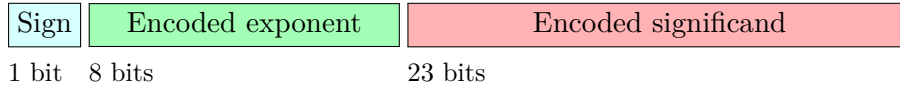
Now that we have explained how a unique set of parameters (sign, exponent, significand) can be assigned to any floating point number, we describe how these parameters are stored on the computer in practice. As their names suggest, the **Float16**, **Float32** and **Float64** formats use 16, 32 and 64 bits of memory, respectively. A naive approach for encoding these number formats would be to store the full binary representations of the sign, exponent and significand.

For the **Float32** format, this approach would require 1 bit for the sign, 8 bits to cover the 254 possible values of the exponent, and 24 bits for the significand, i.e. for storing  $b_0, \dots, b_{p-1}$ . This leads to a total number of 33 bits, which is one more than is available, and this is without the special values **NaN**, **Inf** and **-Inf**. So how are numbers in the **F<sub>32</sub>** format actually stored? To answer this question, we begin with two observations:

- If  $E > E_{\min}$ , then necessarily  $b_0 = 1$  in the unique representation of the significand. Consequently, the leading bit need not be explicitly specified in the case; it is said to be *implicit*. We will see, as a consequence, that  $p-1$  instead of  $p$  bits are in fact sufficient for the significand.
- In the **F<sub>32</sub>** format, 8 bits at minimum need to be reserved for the exponent, which enables the representation of  $2^8 = 256$  different values, but there are only 254 possible values for

the exponent. This suggests that  $256 - 254 = 2$  combinations of the 8 bits can be exploited in order to represent the special values `Inf`, `-Inf` and `NaN`.

Simplifying a little bit, we may view single precision floating point number as an array of 32 bits as illustrated below:



According to the IEEE 754 standard, the first bit is the sign  $s$ , the next 8 bits  $e_0e_1 \dots e_6e_7$  encode the exponent, and the last 23 bits  $b_1b_2 \dots b_{p-2}b_{p-1}$  encode the significand. Let us emphasize that when we say “encode the exponent” here, we just mean that the bits contain information from which the exponent can be uniquely determined, but we have not yet described how this is achieved. Let us introduce the integer number  $e = (e_0e_1 \dots e_6e_7)_2$ ; that is to say,  $0 \leq e \leq 2^8 - 1$  is the integer number whose binary representation is given by  $e_0e_1 \dots e_6e_7$ . One may determine the exponent and significand of a floating point number from the following rules.

- **Denormalized numbers:** If  $e = 0$ , then the implicit leading bit  $b_0$  is zero, the fraction is  $b_1b_2 \dots b_{p-2}b_{p-1}$ , and the exponent is  $E = E_{\min}$ . In other words, using the notation of Section 1.2, we have  $x = (-1)^s 2^{E_{\min}} (0.b_1b_2 \dots b_{p-2}b_{p-1})_2$ . In particular, if  $b_1b_2 \dots b_{p-2}b_{p-1} = 00 \dots 00$ , then it holds that  $x = 0$ .
- **Non-denormalized numbers:** If  $0 < e < 255$ , then the implicit leading bit  $b_0$  of the significand is 1 and the fraction is given by  $b_1b_2 \dots b_{p-2}b_{p-1}$ . The exponent is given by

$$E = e - \text{bias} = E_{\min} + e - 1.$$

where the exponent bias for the single and double precision formats are given in Table 1.2. In this case  $x = (-1)^s 2^{e - \text{bias}} 1.b_1b_2 \dots b_{p-2}b_{p-1}$ . Notice that  $E = E_{\min}$  if  $e = 1$ , as in the case of denormalized numbers.

- **Infinities:** If  $e = 255$  and  $b_1b_2 \dots b_{p-2}b_{p-1} = 00 \dots 00$ , then  $x = \text{Inf}$  if  $s = 0$  and  $-\text{Inf}$  otherwise.
- **Not a Number:** If  $e = 255$  and  $b_1b_2 \dots b_{p-2}b_{p-1} \neq 00 \dots 00$ , then  $x = \text{NaN}$ . Notice that the special value `NaN` can be encoded in many different manners. These extra degrees of freedom were reserved for passing information on the reason for the occurrence of `NaN`, which is usually an indication that something has gone wrong in the calculation.

	Half precision	Single precision	Double precision
Exponent bias ( $-E_{\min} + 1$ )	15	127	1023
Exponent encoding (bits)	5	8	11
Significand encoding (bits)	10	23	52

Table 1.2: Encoding parameters for floating point formats

*Remark 1.3* (Encoding efficiency). With 32 bits, at most  $2^{32}$  different numbers could in principle be represented. In practice, as we saw in [Exercise 1.10](#), the **Float32** format enables to represent

$$(E_{\max} - E_{\min})2^p + 2^{p+1} - 1 = 253 \times 2^{23} + 2^{25} - 1 = 2^{32} - 2^{24} - 1 \approx 99.6\% \times 2^{32},$$

different real numbers, indicating a very good encoding efficiency.

*Remark 1.4* (Nonuniqueness of the floating point representation of 0.0). The sign  $s$  is clearly unique for any number in a floating point format, except for 0.0, which could in principle be represented as

$$(-1)^0 2^{E_{\min}} (0.00 \dots 00)_2 \quad \text{or} \quad (-1)^1 2^{E_{\min}} (0.00 \dots 00)_2.$$

In practice, both representations of 0.0 are available on most machines, and these behave slightly differently. For example  $1/(0.0) = \text{Inf}$  but  $1/(-0.0) = -\text{Inf}$ .

## 1.5 Integer formats

The machine representation of integer formats is much simpler than that of floating point numbers. In this short section, we give a few orders of magnitude for common integer formats and briefly discuss overflow issues. Programming languages typically provide integer formats based on 16, 32 and 64 bits. In Julia, these correspond to the types **Int16**, **Int32** and **Int64**, the latter being the default for integer literals.

The most common encoding for integer numbers, which is used in Julia, is known as *two's complement*: a number encoded with  $p$  bits given by  $b_{p-1}b_{p-2} \dots b_0$  corresponds to

$$x = -b_{p-1}2^{p-1} + \sum_{i=0}^{p-2} b_i 2^i.$$

This encoding enables to represent uniquely all the integers from  $N_{\min} = -2^{p-1}$  to  $N_{\max} = 2^{p-1} - 1$ . In contrast with floating point formats, integer formats do not provide special values like **Inf** and **NaN**. The number delivered by the machine when a calculation exceeds the maximum representable value in the format, called the *overflow behavior*, generally depends on the programming language.

Since the overflow behavior of integer formats is not universal across programming languages, a detailed discussion is of little interest. We only mention that Julia uses a *wraparound* behavior, where  $N_{\max} + 1$  silently returns  $N_{\min}$  and, similarly,  $-N_{\min} - 1$  gives  $N_{\max}$ ; the numbers loop back. This can lead to unexpected results, such as  $2^{64}$  evaluating to 0.

## 1.6 Exercises

⚙ **Exercise 1.1.** Show that if a number  $x \in \mathbf{R}$  admits a finite representation (1.1) in base  $\beta$ , then it also admits an infinite representation in the same base. **Hint:** You may have learned before that  $(0.\overline{9})_{10} = 1$ .

⚙ **Exercise 1.2.** How many digits does it take to represent all the integers from 0 to  $10^{10} - 1$  in decimal and binary formats? What about the hexadecimal format?

⚙ **Exercise 1.3.** Find the decimal representation of  $(0.000\overline{1100})_2$ .

⚙ **Exercise 1.4.** Find the binary representation of  $(0.1)_{10}$ .

□ **Exercise 1.5.** Implement *Algorithm 1* on a computer and verify that it works. Your function should take two arguments: an array of integers  $[a_1, \dots, a_n]$  containing the digits after the decimal point and the  $m$  number of bits to return. The bits should be returned as an array  $[b_1, \dots, b_m]$ .

□ **Exercise 1.6.** As mentioned above, *Algorithm 1* works only for decimal numbers of the specific form  $x = (0.a_1 \dots a_n)_{10}$ . Find and implement a similar algorithm for integer numbers. More precisely, write a function that takes an integer  $n$  as argument and returns an array containing the bits of the binary expansion  $(b_m \dots b_0)_2$  of  $n$ , from the least significant  $b_0$  to the most significant  $b_m$ . That is to say, your code should return  $[b_0, b_1, \dots]$ .

```
function to_binary(n)
    # Your code comes here
end

# Check that it works
number = 123456789
bits = to_binary(number)
pows2 = 2 .^ range(0, length(bits) - 1)
@assert sum(bits*pows2) == number
```

⚙ **Exercise 1.7.** Show that the successor of 1 in  $\mathbf{F}_{64}$  is  $1 + \varepsilon_{64}$ , where  $\varepsilon_{64}$  is the machine epsilon for the double-precision format.

⚙ **Exercise 1.8.** Write down the values of the smallest and largest, in absolute value, positive real numbers representable in the  $\mathbf{F}_{32}$  and  $\mathbf{F}_{64}$  formats.

⚙ **Exercise 1.9** (Relative error and machine epsilon). Prove that the inequality (1.3) is sharp. That is to say, find  $x \in \mathbf{R}$  such that the inequality is an equality.

⚙ **Exercise 1.10** (Cardinality of the set of floating point numbers). Show that, if  $E_{\max} \geq E_{\min}$ , then  $\mathbf{F}(p, E_{\min}, E_{\max})$  contains exactly

$$(E_{\max} - E_{\min})2^p + 2^{p+1} - 1$$

distinct real numbers. (In particular, the special values `Inf`, `-Inf` and `NaN` are not counted.) **Hint:** Count first the numbers with  $E > E_{\min}$  and then those with  $E = E_{\min}$ .

⚙️ **Exercise 1.11.** Calculate the machine epsilon  $\varepsilon_{16}$  for the  $\mathbf{F}_{16}$  format. Write the results of the arithmetic operations  $1 \hat{+} \varepsilon_{16}$  and  $1 \hat{-} \varepsilon_{16}$  in the form

$$2^E(1.b_1 \dots b_{p-1})_2.$$

⚙️ **Exercise 1.12.** Let  $\varepsilon_{16}$  be the machine epsilon for the  $\mathbf{F}_{16}$  format, and define  $y = \frac{4}{3}\varepsilon_{16}$ . What is the relative error between  $\Delta = (1 + y) - 1$ , and the machine approximation  $\hat{\Delta} = (1 \hat{+} y) \hat{-} 1$ ?

□ **Exercise 1.13** (Numerical differentiation). Let  $f(x) = \exp(x)$ . By definition, the derivative of  $f$  at 0 is given by

$$f'(0) = \lim_{\delta \rightarrow 0} \left( \frac{f(\delta) - f(0)}{\delta} \right).$$

The expression within brackets on the right-hand side may be used with a small but nonzero  $\delta$  as an approximation for  $f'(0)$ . Implement this approach using double-precision numbers and the same values for  $\delta$  as in the table below. Explain the results you obtain.

$\delta$	$\frac{\varepsilon_{64}}{4}$	$\frac{\varepsilon_{64}}{2}$	$\varepsilon_{64}$
Approximation of $f'(0)$	0	2	1

□ **Exercise 1.14** (Avoiding overflow). Write a code to calculate the weighted average

$$S := \frac{\sum_{j=0}^J w_j j}{\sum_{j=0}^J w_j}, \quad w_j = \exp(j), \quad J = 1000.$$

You may need to first rewrite  $S$  differently.

□ **Exercise 1.15.** Plot the function  $x \mapsto \log(e^{e^x} - 1)$  over the interval  $[0, 10]$ .

□ **Exercise 1.16** (Calculating the sample variance). Assume that  $(x_n)_{1 \leq n \leq N}$ , with  $N = 10^6$ , are independent random variables distributed according to the uniform distribution  $\mathcal{U}(L, L + 1)$ . That is, each  $x_n$  takes a random value uniformly distributed between  $L$  and  $L + 1$  where  $L = 10^9$ . In Julia, these samples can be generated with the following lines of code:

```
N, L = 10^6, 10^9
x = L .+ rand(N)
```

It is well known that the variance of  $x_n \sim \mathcal{U}(L, L + 1)$  is given by  $\sigma^2 = \frac{1}{12}$ . Numerically, the variance can be estimated from the sample variance:

$$s^2 = \frac{1}{N-1} \left( \left( \sum_{n=1}^N x_n^2 \right) - N\bar{x}^2 \right), \quad \bar{x} = \frac{1}{N} \sum_{n=1}^N x_n. \quad (1.7)$$

Write a computer code to calculate  $s^2$  with the best possible accuracy. Can you find a formula that enables better accuracy than (1.7)?

*Remark 1.5.* In order to estimate the true value of  $s^2$  for your samples, you can use the **BigFloat** format, to which the array  $x$  can be converted by using the instruction

```
x = BigFloat.(x)
```

□ **Exercise 1.17.** Euler proved that

$$\frac{\pi^2}{6} = \lim_{N \rightarrow \infty} \sum_{n=1}^N \frac{1}{n^2}.$$

Using the default **Float64** format, estimate the error obtained when the series on the right-hand side is truncated after  $10^{10}$  terms. Can you rearrange the sum for best accuracy?

⚙️ **Exercise 1.18.** Let  $x$  and  $y$  be positive real numbers in the interval  $[2^{-10}, 2^{10}]$  (so that we do not need to worry about denormalized numbers, assuming we are working in single or double precision), and let us define the machine addition operator  $\hat{+}$  for arguments in real numbers as

$$\hat{+} : \mathbf{R} \times \mathbf{R} \rightarrow \mathbf{R}; (x, y) \mapsto \text{fl}(\text{fl}(x) + \text{fl}(y)).$$

Prove the following bound on the relative error between the sum  $x + y$  and its machine approximation  $x \hat{+} y$ :

$$\frac{|(x + y) - (x \hat{+} y)|}{|x + y|} \leq \frac{\varepsilon_M}{2} \left(2 + \frac{\varepsilon_M}{2}\right).$$

*Hint:* decompose the numerator as

$$(x + y) - (x \hat{+} y) = (x - \text{fl}(x)) + (y - \text{fl}(y)) + (\text{fl}(x) + \text{fl}(y) - (x \hat{+} y)),$$

and then use *Proposition 1.1*.

⚙️ **Exercise 1.19.** Is `Float32(0.1) * Float32(10) == 1` equal to **true** or **false** given the default rounding rule defined by the IEEE standard? Explain.

*Solution.* By default, real numbers are rounded to the nearest floating point number. This can be checked in Julia with the command `rounding(Float32)`, which prints the default rounding mode. The exact binary representation of the real number  $x = 0.1$  is

$$\begin{aligned} x &= (0.000\overline{1100})_2 \\ &= 2^{-4} \times \underbrace{(1.10011001100110011001100\overline{1100})}_{24 \text{ bits}}_2 \end{aligned}$$

The first task is to determine the member of  $\mathbf{F}_{32}$  that is nearest  $x$ . We have

$$\begin{aligned} x^- &:= \max\{y : y \in \mathbf{F}_{32} \text{ and } y \leq x\} = 2^{-4} \times (1.10011001100110011001100)_2 \\ x^+ &:= \min\{y : y \in \mathbf{F}_{32} \text{ and } y \geq x\} = 2^{-4} \times (1.10011001100110011001101)_2. \end{aligned}$$

Since the number  $(0.\overline{1100})_2$  is closer to 1 than to 0, the number  $x$  is closer to  $x^+$  than to  $x^-$ . Therefore, the number obtained when writing `Float32(0.1)` is  $x^+$ . To conclude the exercise, we need to calculate  $\text{fl}(10 \times x^+)$ , and to this end we first write the exact binary representation of the real number

$10 \times x^+ = (1010)_2 \times x^+$ . We have

$$\begin{aligned} (1010)_2 \times x^+ &= (1000)_2 \times x^+ + (10)_2 \times x^+ = 2^{-4} \times (1100.11001100110011001101)_2 \\ &\quad + 2^{-4} \times (11.0011001100110011001101)_2 \\ &= 2^{-4} \times \underbrace{(10000.0000000000000000000001)_2}_{24 \text{ bits}}. \end{aligned}$$

This can be checked in Julia by writing `bitstring(Float32(0.1) * Float64(10.0))`. Clearly, when rounding to the nearest  $\mathbf{F}_{32}$  number, the number  $2^{-4}(10000)_2 = 1$  is obtained; the boolean expression in the question is thus **true**.  $\triangle$

*Remark 1.6.* It should not be inferred from Exercise 1.19 that `Float32(1/i) * i` is always exact in floating point arithmetic. For example `Float32(1/41) * 41` does not evaluate to 1, and neither do `Float16(1/11) * 11` and `Float64(1/49) * 49`.

⚙️ **Exercise 1.20.** Given that the default rounding rule specified by the IEEE 754 standard is “round to nearest, tie to even”, does `Float16(0.1) + Float16(0.2) == Float16(0.3)` evaluate to **true** or **false**. Explain.

⚙️ **Exercise 1.21.** Explain why `Float32(sqrt(2))^2 - 2` is not zero in Julia.

*Solution.* The exact binary representation of  $x := \sqrt{2}$ , found using the `bitstring` in the `Float64` format, is given by

$$x = \underbrace{(1.01101010000010011110011001100\dots)_2}_{24 \text{ bits}}.$$

The first task is to determine the member of  $\mathbf{F}_{32}$  that is nearest  $x$ . We have

$$\begin{aligned} x^- &:= \max\{x : x \in \mathbf{F}_{32} \text{ and } x \leq \sqrt{2}\} = \underbrace{(1.01101010000010011110011)_2}_{24 \text{ bits}} \\ x^+ &:= \min\{x : x \in \mathbf{F}_{32} \text{ and } x \geq \sqrt{2}\} = \underbrace{(1.01101010000010011110100)_2}_{24 \text{ bits}}, \end{aligned}$$

and we calculate

$$\begin{aligned} x - x^- &= 2^{-24}(0.01100\dots)_2, \\ x^+ - x &= 2^{-21}(1 - (0.11001100\dots)_2) \geq 2^{-21}(1 - (0.11001101)_2) = 2^{-21}(0.00110011)_2. \end{aligned}$$

We deduce that  $x - x^- < x^+ - x$ , and so  $\text{fl}(x) = x^-$ . To conclude the exercise, we need to show that  $\text{fl}((x^-)^2)$  is not equal to 2. The exact binary expansion of  $(x^-)^2$  is

$$(x^-)^2 = \underbrace{(1.11111111111111111111111111)_2}_{24 \text{ bits}} 01101100111111010101001_2.$$

The member of  $\mathbf{F}_{32}$  nearest this number is

$$(1.111111111111111111111111)_2 = 2 - 2^{-23},$$



which is precisely the result returned by Julia.  $\triangle$

**Exercise 1.22** (Numerical differentiation). Let  $f(x) = \exp(x)$  and let  $d(\delta)$  be the approximation of  $f'(x)$  obtained from the following piece of code:

```
f, x = exp, 1
d(δ) = (f(x+δ) - f(x))/δ
```

Plot for fixed  $x = 1$  the error  $\text{abs}(d(\delta) - \text{exp}'(x))$  as a function of  $\delta$  in logarithmic scale, and explain the result.

*Solution.* We assume that  $\delta \in \mathbf{F}_{64}$  for simplicity. This is not a restrictive assumption as  $\delta$  can only take floating point values in computer programs. The proof is rather technical, and so it is given for information purposes only. We rewrite  $d(\delta)$  mathematically as

$$d(\delta) = \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) \widehat{/} \delta,$$

where  $\widehat{f}(x) = \text{fl}(f(x))$ . We wish to bound  $|f'(x) - d(\delta)|$ . By the triangle inequality,

$$|f'(x) - d(\delta)| \leq \left| f'(x) - \frac{f(x+\delta) - f(x)}{\delta} \right| + \left| \frac{f(x+\delta) - f(x)}{\delta} - \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) \widehat{/} \delta \right|. \quad (1.8)$$

**Bounding the first term in (1.8).** By Taylor's theorem, there exists  $\xi \in [x, x + \delta]$  such that

$$f(x + \delta) = f(x) + \delta f'(x) + \frac{\delta^2}{2} f''(\xi).$$

Therefore, the first term in (1.8) satisfies

$$\left| f'(x) - \frac{f(x+\delta) - f(x)}{\delta} \right| = \frac{\delta}{2} |f''(\xi)| = \frac{\delta}{2} |f''(x)| + \mathcal{O}(\delta^2).$$

**Bounding the second term in (1.8) – the roundoff error.** This is more tedious but not difficult; the main ingredient is the triangle inequality. Specifically, we will use the bound

$$\begin{aligned} & \left| \frac{f(x+\delta) - f(x)}{\delta} - \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) \widehat{/} \delta \right| \\ & \leq \delta^{-1} |f(x+\delta) - f(x \widehat{+} \delta)| + \delta^{-1} |f(x \widehat{+} \delta) - \widehat{f}(x \widehat{+} \delta)| + \delta^{-1} |\widehat{f}(x) - f(x)| \\ & \quad + \delta^{-1} \left| \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) - \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) \right| \\ & \quad + \left| \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) / \delta - \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) \widehat{/} \delta \right|. \end{aligned} \quad (1.9)$$

Note that without absolute values, both sides are indeed equal. The first three terms on the right-hand side are together a bound from above for

$$\left| \frac{f(x+\delta) - f(x)}{\delta} - \frac{\widehat{f}(x \widehat{+} \delta) - \widehat{f}(x)}{\delta} \right|, \quad (1.10)$$

while the two other terms account for the roundoff errors associated with the machine subtraction and division operators, respectively. We will show that the dominant terms in (1.9) are the first three; the latter two are negligible in comparison. Employing [Proposition 1.1](#) with  $x = f(a)$ , with  $x = a \pm b$  and

with  $x = a/b$ , we deduce that the following inequalities are satisfied:

$$\left| \widehat{f}(a) - f(a) \right| \leq \varepsilon |f(a)|, \quad |(a \pm b) - (a \widehat{\pm} b)| \leq \varepsilon |a \pm b|, \quad |a/b - a \widehat{\nearrow} b| \leq \varepsilon |a/b|. \quad (1.11)$$

The first two inequalities are valid for all  $(a, b) \in \mathbf{F}_{64} \times \mathbf{F}_{64}$ , while the third inequality is valid for all  $(a, b) \in \mathbf{F}_{64} \times \mathbf{F}_{64} \setminus \{0\}$ . The first inequality in (1.11) can be employed in order to bound the second and third term on the right-hand side of (1.9):

$$\begin{aligned} \delta^{-1} \left| f(x \widehat{+} \delta) - \widehat{f}(x \widehat{+} \delta) \right| + \delta^{-1} \left| \widehat{f}(x) - f(x) \right| &\leq \delta^{-1} \varepsilon \left( |f(x + \delta)| + |f(x)| \right) \\ &= 2\delta^{-1} \varepsilon |f(x)| + \mathcal{O}(\varepsilon). \end{aligned} \quad (1.12)$$

Using Taylor's theorem and then the second inequality in (1.11) (with +), we then bound the first term on right-hand side of (1.9) as

$$\begin{aligned} \delta^{-1} \left| f(x + \delta) - f(x \widehat{+} \delta) \right| &= \delta^{-1} \left| f'(\xi) \left( (x + \delta) - (x \widehat{+} \delta) \right) \right| \\ &\leq \delta^{-1} \varepsilon |f'(\xi)| |x + \delta| = \delta^{-1} \varepsilon |f'(x)x| + \mathcal{O}(\varepsilon). \end{aligned} \quad (1.13)$$

Combining (1.12) and (1.13), we obtain that the expression in (1.10) scales as  $\delta^{-1} \varepsilon$ :

$$\left| \frac{f(x + \delta) - f(x)}{\delta} - \frac{\widehat{f}(x \widehat{+} \delta) - \widehat{f}(x)}{\delta} \right| = \mathcal{O}(\delta^{-1} \varepsilon). \quad (1.14)$$

Next, using the second inequality in (1.11) (with -) and then (1.14), we bound the fourth term in (1.9) as follows:

$$\begin{aligned} \delta^{-1} \left| \left( \widehat{f}(x \widehat{+} \delta) - \widehat{f}(x) \right) - \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) \right| \\ \leq \delta^{-1} \varepsilon \left| \widehat{f}(x \widehat{+} \delta) - \widehat{f}(x) \right| = \delta^{-1} \varepsilon \underbrace{|f(x + \delta) - f(x)|}_{=\mathcal{O}(\delta)} + \mathcal{O}(\delta^{-1} \varepsilon^2) = \mathcal{O}(\varepsilon + \delta^{-1} \varepsilon^2). \end{aligned} \quad (1.15)$$

Notice that this term is smaller than the first three in (1.9) when  $\delta \ll 1$ . We deduce by a triangle inequality from (1.14) and (1.15) that

$$\delta^{-1} \left| \left( f(x + \delta) - f(x) \right) - \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) \right| = \mathcal{O}(\delta^{-1} \varepsilon).$$

Finally, using the third inequality in (1.11) together with this equation, we bound the fifth term on the right-hand side of (1.9):

$$\begin{aligned} \left| \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) / \delta - \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) \widehat{\nearrow} \delta \right| \\ \leq \varepsilon \left| \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right| / \delta = \frac{\varepsilon}{\delta} \underbrace{|f(x + \delta) - f(x)|}_{=\mathcal{O}(\delta)} + \mathcal{O}(\delta^{-1} \varepsilon^2) = \mathcal{O}(\varepsilon + \delta^{-1} \varepsilon^2). \end{aligned}$$

This term is also negligible in front of the other dominant contributions to the roundoff error given in (1.12) and (1.13). Going back to (1.9), we conclude that

$$\left| \frac{f(x + \delta) - f(x)}{\delta} - \left( \widehat{f}(x \widehat{+} \delta) \widehat{-} \widehat{f}(x) \right) \widehat{\nearrow} \delta \right| \leq \delta^{-1} \varepsilon \left( 2|f(x)| + |f'(x)x| \right) + \mathcal{O}(\varepsilon + \delta^{-1} \varepsilon^2).$$

**Concluding the proof.** Going back to (1.8), we conclude that

$$|f'(x) - d(\delta)| \leq \frac{\delta}{2}|f''(x)| + \frac{\varepsilon}{\delta}(2|f(x)| + |f'(x)x|),$$

up to higher order terms. For fixed  $x$ , the right-hand side is minimized when

$$\delta = \sqrt{\varepsilon} \sqrt{\frac{2|f(x)| + |f'(x)x|}{|f''(x)|}},$$

which is a well-known formula for the optimal step size in numerical differentiation. The error as a function of  $\delta$  for  $x = 1$  is depicted in Figure 1.4.  $\triangle$

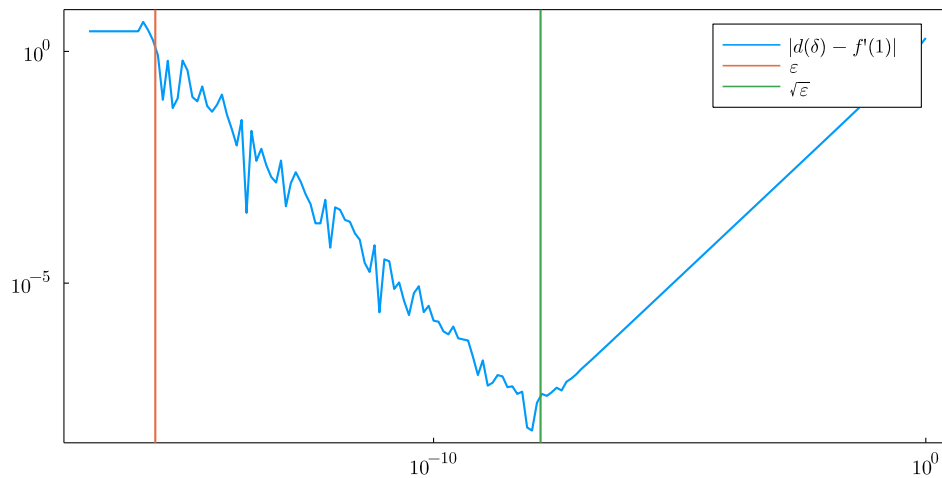


Figure 1.4: Solution to Exercise 1.22.

⚙️ **Exercise 1.23.** Explain why `exp(log(Float16(7))) == 7` is **false**.

*Solution.* We begin by finding the binary representation of  $\log(7)$ . In the  $\mathbf{F}_{64}$  format, the sign and exponent are encoded over 1 and 11 bits respectively, and so the fraction is given by the last 52 bits, which can be obtained from the command

```
julia> bitstring(log(7))[13:end]
"1111001000100111001010101110001100100101101001010111"
```

Therefore, since `exponent(log(7))` returns 0, we have

$$\log(7) = (\underbrace{1.1111001000}_{11 \text{ bits}} 10011100101010111000110010010110100101011\dots)_2 \quad (1.16)$$

The number returned by the command `log(Float16(7))` is given by  $\text{fl}_{16}(\log(7))$ , where  $\text{fl}_{16}$  denotes the half-precision rounding operator. Rounding the right-hand side of (1.16) to 11 bits, we obtain

$$\text{fl}_{16}(\log(7)) = (1.1111001001)_2 = 1.9462890625.$$

The number returned by the code `exp(log(Float16(7)))` is

$$\text{fl}_{16}\left(\exp\left(\text{fl}_{16}(\log(7))\right)\right) = \text{fl}_{16}(\exp(1.9462890625)).$$

The rounding operator appears twice on the left-hand side, because the computer rounds *after every operation*. To explain the result of the rounding operation, we begin by calculating the binary expansion of  $\exp(1.9462890625)$ .

```
julia> bitstring(exp(1.9462890625))[13:end]
"1100000000101011011101110000111000100001100010001110"
```

```
julia> exponent(exp(1.9462890625))
2
```

Therefore,

$$\exp(1.9462890625) = 2^2 \underbrace{(1.1100000000)}_{11 \text{ bits}} 1010110111011100001110001000011000100011\dots)_2,$$

and, rounding to 11 bits, we finally obtain

$$\begin{aligned} \text{fl}_{16}(\exp(1.9462890625)) &= 2^2(1.1100000001)_2 = 4 \left( 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{2^{10}} \right) \\ &= 7 + \frac{1}{2^8} = 7.00390625, \end{aligned}$$

which is different from 7. △

⚙️ **Exercise 1.24.** Determine the encoding of the following **Float32** numbers:

- $x_1 = 2.0^{E_{\min}}$
- $x_2 = -2.0^{E_{\min}-p-1} = -2.0^{-149}$
- $x_3 = 2.0^{E_{\max}}(2 - 2^{-p+1})$

Check your results using the Julia function `bitstring`.

⚙️ **Exercise 1.25** (Summary). True or false?

1. Let  $(\bullet)_2$  denote binary representation. It holds that  $(0.1111)_2 + (0.0001)_2 = 1$ .
2. It holds that  $(1000)_2 \times (0.001)_2 = 1$ .
3. It holds that
 
$$(0.\bar{1})_3 = \frac{1}{2}.$$
4. In base 16, all the natural numbers from 1 to 200 can be represented using 2 digits.
5. In Julia, `Float64(.1) == Float32(.1)` evaluates to **true**.
6. The spacing (in absolute value) between successive double-precision (**Float64**) floating point numbers is constant.
7. It holds that  $(0.\overline{10101})_2 = (1.2345)_{10}$ .
8. Machine addition  $\hat{+}$  is an associative operation. More precisely, given any three double-precision floating point numbers  $x$ ,  $y$  and  $z$ , the following equality holds:

$$(x \hat{+} y) \hat{+} z = x \hat{+} (y \hat{+} z).$$

9. The machine epsilon is the smallest strictly positive number that can be represented in a floating point format.
10. Let  $\varepsilon$  denote the machine epsilon for the double-precision format. Let also  $\hat{+}$  and  $\hat{/}$  denote respectively the machine addition and the machine division operators for the double-precision format. It holds that  $1 \hat{+} (\varepsilon \hat{/} 64) = 1$  and that  $\varepsilon \hat{/} 64 \neq 0$ .
11. Assume that  $x \in \mathbf{R}$  belongs to the double-precision floating point format (that is, assume that  $x \in \mathbf{F}_{64}$ ). Then  $-x \in \mathbf{F}_{64}$ .

## 1.7 Discussion and bibliography

This chapter is mostly based on the original 1985 IEEE 754 standard [4] and the reference book [10]. A significant revision to the 1985 IEEE standard was published in 2008 [5], adding for example specifications for the half precision and quad precision formats, and a minor revision was published in 2019 [6]. The original IEEE standard and its revisions constitute the authoritative guide on floating point formats. It was intended to be widely disseminated and is written very clearly and concisely, but is not available for free online. Another excellent source for learning about floating point numbers and round-off errors is D. Goldberg’s paper “*What every computer scientist should know about floating-point arithmetic*” [3], freely available online.