

# 1 Representation of real numbers on a computer

## 1.1 Floating-point arithmetic (IEEE754: Double precision)

Computers have a limited amount of memory. Therefore, it is not possible to handle an infinite number of digits, so it is necessary to give up at some point. Furthermore, the computer can only handle binary numbers of **0** and **1** inside, so we need to discuss how to fit them into a finite number of digits. Right now, the CPUs of mainstream computers are based on the **IEEE754** standard. The IEEE754 standard includes a technical standard for **floating-point**.

$\pi$  is an irrational number, so it cannot be held by a computer. However, in the IEEE754 standard, even 1.1 cannot be represented exactly in decimal. The error is already there before the calculation is done!

Let's take a look at a simple example.

Listing 1 roundoff.c

```
1 #include <stdio.h>
2
3 int main(){
4     double i=1.1;
5     printf("%.20f\n",i);
6 }
```

If you don't have any programming environment (if it's too much trouble), use **CES** (<https://www.ces-alpha.org/en/>) or **Ganjin** (<https://ganjin.online/>).

Enter the following command and try to execute it.

```
$ gcc roundoff.c
$ ./a.out
1.100000000000000008882
```

It is certainly not exactly 1.1, and you can see that there is an error. Let's check IEEE754 to see why such an error occurred.

First of all, IEEE754 has two formats: *single* and *double*. *single* is 32 bits long, and *double* is 64 bits long.

$$(-1)^s \times 2^e \times m$$

Sign bit  $s$  : 1 bit (0 or 1)

Exponent  $e$  :  $e_{\min} \leq e \leq e_{\max}$

Significand  $m$  : It is expressed in the normalized form,  $d_0, d_1, d_2, \dots, d_{p-1}$ . However,  $d_i$  can be 0 or 1.

The values of  $e_{\max}$ ,  $e_{\min}$ , and  $p$  vary depending on *single* and *double*, and are defined as shown in Table 1.

Table 1 Exponent and significand

Parameter	single	double
$p$	24 bit	53 bit
$e_{\max}$	127	1023

Figure 1 shows the representation for the *Single* case. In the state called normalized significand,  $d_0$  is always set to 1, and the numbers are stored starting from  $d_1$ . It is called the “hidden” or “implicit” bit.

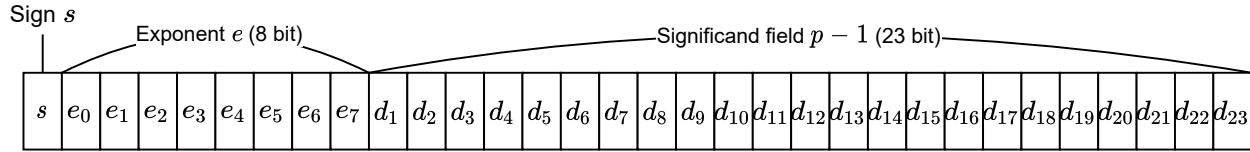


Fig.1 An example of a layout for 32-bit (*Single*) floating-point

If we convert 1.1 to binary, we get a recurring decimal like  $(1.000110011001\dots)_2$ . This is the cause of the error.

The precision of the calculation (significant digits) in *double* is  $2^{-52} \simeq 2.22 \times 10^{-16} \simeq 10^{-15.65}$ , which is approximately 16 digits.

#### Exercise

Convert the decimal number 7.25 to a floating-point number. However,  $p$  (including  $d_0$ ) must be 6 bits, and  $e$  must be 3 bits.

$$\begin{aligned} 7.25 &= (111.01)_2 \\ &= (-1)^0 \times 2^2 \times (1.1101)_2 \end{aligned}$$

Therefore, the answer is

$$0 \ 010 \ 11010$$

#### Example of a large error (quadratic equations)

The solution to the quadratic equation  $ax^2 + bx + c = 0$  is

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

Transforming this equation, we get

$$\begin{aligned} \frac{-b + \sqrt{b^2 - 4ac}}{2a} &= \frac{b^2 - b^2 + 4ac}{2a(-b - \sqrt{b^2 - 4ac})} \\ &= \frac{4ac}{2a(-b - \sqrt{b^2 - 4ac})} \\ &= \frac{2c}{-b - \sqrt{b^2 - 4ac}} \end{aligned} \tag{1}$$

When calculated with  $a = 1$ ,  $b = 10^{15}$ ,  $c = 10^{14}$ .

#### Listing 2 quadratic.c

```
1 #include <stdio.h>
2 #include <math.h>
3
4 int main(){
5     double a=1;
6     double b=1000000000000000;
7     double c=100000000000000;
```

```

8
9     double x1=(-b+sqrt(b*b-4*a*c))/(2*a);
10    printf("%.16f\n",x1);
11
12     double x2=(2*c)/(-b-sqrt(b*b-4*a*c)); // Equation (1)
13
14    printf("%.16f\n",x2);
15 }

```

```

$ gcc quadratic.c -lm
$ ./a.out
-0.125000000000000000
-0.100000000000000000

```

This result shows that

$$\frac{-b + \sqrt{b^2 - 4ac}}{2a} = -0.125000000000000000$$

$$\frac{2c}{-b - \sqrt{b^2 - 4ac}} = -0.100000000000000000$$

but it is strange that the answers are different.

## 2 Interval arithmetic

### 2.1 Operation between floating-point numbers

Let  $\mathbb{F}$  be the set of floating-point numbers defined by IEEE754. Rounding, such as truncation and rounding up, in IEEE754 is defined only for the results of the five operations  $+$ ,  $-$ ,  $*$ ,  $/$  and  $\sqrt{\cdot}$ . Note that rounding of input, output, sin, cos, etc. is not defined.

Let  $\circ \in \{+, -, \times, /\}$ . In this case, the following rounding modes exist.

Round to nearest :  $\tilde{\circ} : \mathbb{F} \times \mathbb{F} \rightarrow \mathbb{F}$ , and round to the nearest floating-point number.

Round upward :  $\bar{\circ} : \mathbb{F} \times \mathbb{F} \rightarrow \mathbb{F}$ ,  $\inf\{x \in \mathbb{F} \mid x \geq a \circ b\}$

Round downward :  $\underline{\circ} : \mathbb{F} \times \mathbb{F} \rightarrow \mathbb{F}$ ,  $\sup\{x \in \mathbb{F} \mid x \leq a \circ b\}$

Consider floating-point operations. That is, if  $a, b \in \mathbb{F}$ , then the four arithmetic operations can be defined as follows.

$$a + b \in [a_{\pm}b, a_{\mp}b] = \{c \in \mathbb{R} \mid a_{\pm}b \leq c \leq a_{\mp}b\}$$

$$a - b \in [a_{\pm}b, a_{\mp}b]$$

$$a \cdot b \in [a_{\pm}b, a_{\mp}b]$$

$$a/b \in [a_{\pm}b, a_{\mp}b]$$

In C compilers (such as gcc), `fesetround` function can be used to change the rounding mode by using `fenv.h`.

```
//You can change the rounding mode like this!
```

```
#include <fenv.h>
```

```
fesetround(FE_UPWARD); //Change to upward rounding mode
```

```
fesetround(FE_DOWNWARD); //Change to downward rounding mode
fesetround(FE_TONEAREST); //Change to nearest rounding mode
```

I mentioned earlier that it is difficult to accurately represent 1.1 on a computer, but it is now possible to calculate the correct interval.

What is the correct interval? For example,  $\pi \simeq 3.14$  is just an approximation, but if we say  $\pi \in [3.14, 3.15]$ , this is the correct interval (correct information). Let's check this in the actual code.

```
#include <stdio.h>
#include <math.h>
#include <fenv.h>

int main(){
    fesetround(FE_UPWARD);
    double a_u=11.0/10.0;

    fesetround(FE_DOWNWARD);
    double a_l=11.0/10.0;

    /*
    Do not use 1.1
    Recall that the rounding mode can be applied only to certain operations.
    */

}
```

In this case, we should have  $11/10 \in [a_l, a_u]$ . (However, we cannot control the rounding mode inside the `printf` function.)

So, if we want to use interval arithmetic to calculate  $0.1+1.1$ , how do we do that? First of all, as in the case of  $11/10$ , we can take the interval  $1/10 \in [b_l, b_u]$ . Therefore, by calculating as in Figure 2, we can get the correct interval,  $a + b \in [c_l, c_u]$ , can't we?

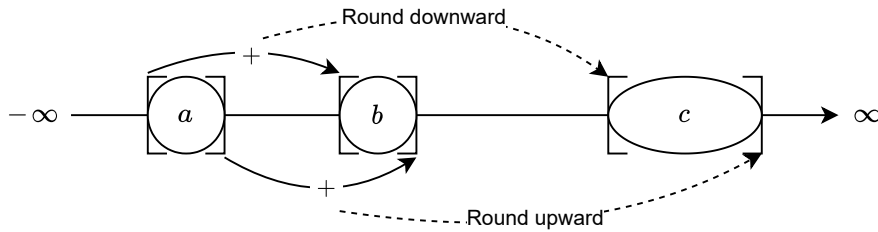


Fig.2  $[a_l, a_u] + [b_l, b_u]$

Thus, the four arithmetic operations on the interval can be summarized as follows. ( $a_l, a_u, b_l, b_u \in \mathbb{F}$ )

Addition :  $\{a + b \mid a \in [a_l, a_u], b \in [b_l, b_u]\} = [a_l, a_u] + [b_l, b_u] = [a_l \pm b_l, a_u \mp b_u]$

Subtraction :  $[a_l, a_u] - [b_l, b_u] = [a_l \mp b_u, a_u \mp b_l]$

Multiplication :  $[a_l, a_u] \times [b_l, b_u] = [\min\{a_l \cdot b_l, a_u \cdot b_l, a_l \cdot b_u, a_u \cdot b_u\}, \max\{a_l \cdot b_l, a_u \cdot b_l, a_l \cdot b_u, a_u \cdot b_u\}]$

Division :  $[a_l, a_u] / [b_l, b_u] = [\min\{a_l / b_l, a_u / b_l, a_l / b_u, a_u / b_u\}, \max\{a_l / b_l, a_u / b_l, a_l / b_u, a_u / b_u\}]$  (but  $0 \notin [b_l, b_u]$ )

## Exercise

Modify the code in `addition.c` (Listing 3) to calculate  $0.1 + 1.1$  using interval arithmetic.

Listing 3 `addition.c`

```
1 #include <stdio.h>
2 #include <math.h>
3 #include <fenv.h>
4
5 int main(){
6     double a=1.0/10.0; //0.1
7     double b=11.0/10.0; //1.1
8
9     double c=a+b;
10
11     printf("%lf\n",c);
12 }
```

(How to compile)

```
$ gcc addition.c -lm
```

(execution)

```
$ ./a.out
```

## Interval of solutions to a quadratic equation

This is just an example. When  $a = 1$ ,  $b = 10^{15}$ ,  $c = 15^{14}$ , the result is as follows.

$$\frac{-b + \sqrt{b^2 - 4ac}}{2a} \in [-0.125000, -0.062500]$$

The actual answer is  $x = -0.1$ , so we can say that we have calculated the correct interval.

Listing 4 `int_quadratic.c`

```
1 #include <stdio.h>
2 #include <math.h>
3 #include <fenv.h>
4
5 int main(){
6     double a=1;
7     double b=1000000000000000.0;
8     double c=100000000000000.0;
9
10     fesetround(FE_TONEAREST);
11     double int_sqrt=b*b-4*a*c;
12
13     fesetround(FE_UPWARD);
14     double sqrt_u=sqrt(int_sqrt);
15     fesetround(FE_DOWNWARD);
16     double sqrt_l=sqrt(int_sqrt);
17
18     fesetround(FE_UPWARD);
19     double denom_u= -b + sqrt_u;
20     fesetround(FE_DOWNWARD);
```

```

21     double denom_l= -b + sqrt_l;
22
23     fesetround(FE_TONEAREST);
24     double numerator=2*a;
25     fesetround(FE_UPWARD);
26     double ans_u[4];
27     ans_u[0]=denom_l/numerator;
28     ans_u[1]=denom_u/numerator;
29     fesetround(FE_DOWNWARD);
30     double ans_l[4];
31     ans_l[0]=denom_l/numerator;
32     ans_l[1]=denom_u/numerator;
33
34     double min,max;
35
36     max = ans_u[0];
37     for(int i=0; i<2; i++){
38         if(max < ans_u[i]){
39             max=ans_u[i];
40         }
41     }
42
43     min = ans_l[0];
44     for(int i=0; i<2; i++){
45         if(min > ans_l[i]){
46             min=ans_l[i];
47         }
48     }
49
50     printf("ans_u in [%lf, %lf]\n", min, max);
51
52 }

```

---

Similarly, we can find the interval of the solution for  $\frac{2c}{-b - \sqrt{b^2 - 4ac}}$ .

## 2.2 About matrix operations

Next, we will discuss matrix calculations.

The  $n$ -dimensional vectors  $a \in \mathbb{F}^n$  and  $b \in \mathbb{F}^n$  of floating-point numbers can be computed as follows

$$a \cdot b \subset [a \cdot_+ b, a \cdot_- b]$$

Note that  $a \cdot_+ b$  means that all inner product operations are computed with downward rounding, and  $a \cdot_- b$  means that all inner product operations are computed with upward rounding.