

Floating-point arithmetic

Academic year 2020-2021

System

```
> restart:  
Digits := 10:
```

1.7. Conversion between binary representation and floating-point

Exercise: Convert the 32-bit word "11000100111000000001100000000000" into its single-precision floating-point counterpart. Subsequently, convert the number 10/3 to its single-precision binary representation. Can you represent this number exactly?

Single-precision

For the binary representation of a single-precision floating-point number a 32-bit word is used.

- bit 1: Sign of the number.
- bits 2-9: Exponent of the number (don't forget the bias!).
- bits 10-32: Mantissa or fractional part of the number.

A normalized floating-point is then expressed as

$$(-1)^s * (1+f) * 2^e$$

where f is the mantissa, e the exponent and 2 the basis (since we work in binary).

For double precision a 64-bit word is used.

- bit 1: Sign of the number.
- bits 2-12: Exponent of the number (don't forget the bias!).
- bits 13-64: Mantissa or fractional part of the number.

Binary to floating-point

```
[> float_bin := "11000100111000000001100000000000";  
    float_bin := "11000100111000000001100000000000" (2.2.1)
```

```
[> sign_bin := parse(float_bin[1]);  
    exp_bin := parse(float_bin[2..9]);  
    mant_bin := parse(float_bin[10..32]);  
    sign_bin := 1  
    exp_bin := 10001001  
    mant_bin := 110000000011000000000000 (2.2.2)
```

```
[> sign_float := (-1)^sign_bin;  
    sign_float := -1 (2.2.3)
```

```
[> bias := 127;  
    exp_float := convert(exp_bin,decimal,2);  
    exp_float := exp_float - bias;  
    bias := 127  
    exp_float := 137  
    exp_float := 10 (2.2.4)
```

```
[> mant_float := convert(mant_bin,decimal,2)*2^(-23);  
    mant_float := add((1/2)^i*parse(float_bin[9+i]),i=1..23);  
    mant_float :=  $\frac{3075}{4096}$   
    mant_float :=  $\frac{3075}{4096}$  (2.2.5)
```

```
[> float_number := sign_float*(1+mant_float)*2^(exp_float);  
    float_number := evalf(float_number);  
    float_number :=  $-\frac{7171}{4}$   
    float_number := -1792.750000 (2.2.6)
```

Real number to floating-point

```
[> number := 10/3;  
    number :=  $\frac{10}{3}$  (2.3.1)
```

The sign is positive -- The bit is 0.

```
[> sign_bin := 0; (2.3.2)
```

```
sign_bin := 0 (2.3.2)
```

We have $2 < 10/3 < 4$ -- Hence the exponent has to be 1. However, we can't forget to add the bias to really obtain the stored 8-bit word.

```
> exp_number := 1;
   exp_number := exp_number + bias;
   exp_bin := convert(exp_number,binary);
           exp_number := 1
           exp_number := 128
           exp_bin := 10000000 (2.3.3)
```

If we would work exactly, we would find that

$$(1+f) * 2 = 10/3 \Rightarrow f = 2/3 = 1/2 + 1/8 + 1/32 + 1/128 + \dots$$

However, we have to truncate this sequence somewhere since we only have a finite number of bits. Should we afterwards round up or not?

```
> mant_float1 := add(1/2*(1/4)^i,i=0..11);
   mant_bin1 := convert(mant_float1*2^23,binary);
   mant_bin2 := convert(mant_float1*2^23+1,binary);
           mant_float1 := 5592405/8388608
           mant_bin1 := 10101010101010101010101
           mant_bin2 := 10101010101010101010110 (2.3.4)
```

```
> float_bin1 := cat(convert(sign_bin,string),convert(exp_bin,
string),convert(mant_bin1,string));
float_bin2 := cat(convert(sign_bin,string),convert(exp_bin,
string),convert(mant_bin2,string));
           float_bin1 := "010000000101010101010101010101"
           float_bin2 := "010000000101010101010101010110" (2.3.5)
```

```
> convert_bin_float := proc(float_bin)

   description "convert floating number in binary
representation to decimal representation";
   local float_bin_str, sign_bin, exp_bin, mant_bin,
sign_float, bias, exp_float, mant_float, float_number;

   sign_bin := parse(float_bin[1]);
   exp_bin := parse(float_bin[2..9]);
   mant_bin := parse(float_bin[10..32]);

   sign_float := (-1)^sign_bin;
   bias := 127;
   exp_float := convert(exp_bin,decimal,2);
   exp_float := exp_float - bias;
   mant_float := convert(mant_bin,decimal,2)*2^(-23);
   mant_float := add((1/2)^i*parse(float_bin[9+i]),i=1..23);

   float_number := sign_float*(1+mant_float)*2^(exp_float);
   float_number := evalf(float_number);
```

```
L end proc:
```

```
> float_d10 := evalf(10/3);  
float1 := convert_bin_float(float_bin1);  
float2 := convert_bin_float(float_bin2);  
float_d10 := 3.3333333333  
float1 := 3.333333254  
float2 := 3.333333492
```

(2.3.6)

If \tilde{x} is the computed solution and x^* is the exact solution, then the absolute error is given by

$e = x^* - \tilde{x}$ and the relative error is given by $\frac{\|x^* - \tilde{x}\|}{\|x^*\|}$.

In most cases only the relative error is meaningful:

$A = 0.01, A' = 0.015 \Rightarrow$ absolute error = 0.005, relative error = 0.5 = 50%

```
> abs(float1-float_d10)/abs(float_d10);  
abs(float2-float_d10)/abs(float_d10);  
2.370000000 10-8  
4.770000000 10-8
```

(2.3.7)

1.7. Basis and machine epsilon

Exercise: What is the basis of the hardware floating-point arithmetic in Maple? What is its value for the machine epsilon?

What is the basis for regular floating-point arithmetic in Maple? What is its value for the machine epsilon? Do you find the same results as before?

Hardware floats

To find the number base of floating-point computations we can look at its behaviour. We will use Malcolm's algorithm to determine the base.

```
> one := evalhf(1):  
zero := evalhf(0):  
A := one:  
B := one:  
while evalhf(B-one) = zero do  
    A := evalhf(A + A);  
    B := evalhf(A + one);  
    B := evalhf(B - A);  
  
end do:  
base := 0:  
C := one:  
while base = 0 do
```

```

      C := evalhf(C + C);
      base := evalhf(A + C);
      base := floor(evalhf(base - A));
end do:
print(`The number base is` = base);
      The number base is = 2

```

(3.1.1)

The machine epsilon is *the distance from 1 to the next larger floating-point number*.

```

> eps := one:
while evalhf(one + eps - one) <> zero do
  eps := evalhf(eps/base);
end do:
eps := evalhf(eps*base):
print(`The machine epsilon is` = eps);
print(`The machine epsilon in base 2 is` = log[base](eps));
      The machine epsilon is = 2.22044604925031308 10-16
      The machine epsilon in base 2 is = -52.000000000

```

(3.1.2)

▼ Floating-point arithmetic in Maple

```

> one := evalf(1):
zero := evalf(0):
A := one:
B := one:
while evalhf(B-one) = zero do
  A := evalf(A + A);
  B := evalf(A + one);
  B := evalf(B - A);

end do:
base := 0:
C := one:
while base = 0 do
  C := evalf(C + C);
  base := evalf(A + C);
  base := floor(evalf(base - A));

end do:
print(`The number base is` = base);
      The number base is = 10

```

(3.2.1)

```

> eps := one:
while evalf(one + eps - one) <> zero do
  eps := evalf(eps/base);
end do:
eps := evalf(eps*base):
print(`The machine epsilon is` = eps);
print(`The machine epsilon in base 10 is` = log[base](eps));
print(`The number of digits Maple uses is` = Digits);
      The machine epsilon is = 1.0000000000 10-9
      The machine epsilon in base 10 is = -9.0000000000
      The number of digits Maple uses is = 10

```

(3.2.2)

1.1. The golden ratio and roots of quadratic polynomials

Exercise: Compute the roots of the polynomial $0.001 x^2 + 1000000 x + 0.01$. Do you obtain the results you expected? How can you solve this problem?

We can obtain the roots of a quadratic polynomial with the discriminant formulas.

$$\begin{aligned} > \mathbf{s := solve(a*x^2+b*x+c=0, x);} \\ s := \frac{-b + \sqrt{-4ac + b^2}}{2a}, -\frac{b + \sqrt{-4ac + b^2}}{2a} \end{aligned} \quad (4.1)$$

$$\begin{aligned} > \mathbf{a := 10^{(-3)};} \\ \mathbf{b := 10^6;} \\ \mathbf{c := 10^{(-2)};} \\ \mathbf{r := solve(a*x^2+b*x+c=0, x);} \\ \mathbf{r1 := r[1];} \\ \mathbf{r2 := r[2];} \\ a := \frac{1}{1000} \\ b := 1000000 \\ c := \frac{1}{100} \\ r := -500000000 + \sqrt{249999999999999990}, -500000000 - \sqrt{249999999999999990} \end{aligned} \quad (4.2)$$

$$\begin{aligned} > \mathbf{Digits := 50;} \\ \mathbf{r1_d50 := evalf(r1);} \\ \mathbf{r2_d50 := evalf(r2);} \\ Digits := 50 \\ r1_d50 := -1.00000000000000001000000000000000 10^{-8} \\ r2_d50 := -9.99999999999999989999999999999990000000000000000 10^8 \end{aligned} \quad (4.3)$$

$$\begin{aligned} > \mathbf{r1_hf := evalhf((-b+sqrt(b^2-4*a*c))/(2*a));} \\ \mathbf{r2_hf := evalhf(-(b+sqrt(b^2-4*a*c))/(2*a));} \\ r1_hf := 0. \\ r2_hf := -1.000000000 10^9 \end{aligned} \quad (4.4)$$

$$\begin{aligned} > \mathbf{evalf(abs(r1_hf-r1)/abs(r1));} \\ \mathbf{evalf(abs(r2_hf-r2)/abs(r2));} \\ 1.00 \\ 1.0000000000000000200000000000000000300000000000000 10^{-17} \end{aligned} \quad (4.5)$$

We can use the product rule: the product of the two roots is the ratio of the constant term and the leading coefficient.

$$> \mathbf{expand(s[1]*s[2]);}$$

```
c/a;
10
10
```

(4.6)

```
> r1_hf_ver2 := evalhf(c/a/r2_hf);
r1_hf_ver2 := -1.000000000000000002 10-8
```

(4.7)

```
> evalf(abs(r1_hf_ver2-r1)/abs(r1));
9.9999999999999999000000000000000000000000000000000 10-18
```

(4.8)

1.2. Fibonacci numbers, sequences and series

Exercise: What is the limit of the series $\sum_{n=1}^{\infty} \frac{1}{n^2}$? What is the result you find using floating-point arithmetic? Why do we find this result and how can we obtain a higher accuracy?

```
> Digits := 5:
S := sum(1/(i^2), i=1..infinity);
evalf(S);
```

$$S := \frac{\pi^2}{6}$$

1.6450 (5.1)

We can't compute the sum of an infinite amount of numbers, hence we have to truncate the sum somewhere.

```
> n := 1:
SUM := evalf(0):
while evalf(SUM + 1/(n^2)) > SUM do
  SUM := evalf(SUM + 1/(n^2));
  n := n+1;
end do:
N := n:
print(`The sum is` = SUM);
abs(SUM-S)/abs(S);
100*abs(SUM-S)/abs(S);
print(`The number of terms is` = n);
```

The sum is = 1.6390
0.0036475
0.36475
The number of terms is = 142 (5.2)

Catastrophic cancellation occurred! We tried to add a very small number to a (relatively) larger one.

```
> evalf(1/(N^2));
evalf(SUM + 1/(N^2) - SUM);
evalf(SUM - SUM + 1/(N^2));
```

$$\begin{array}{r}
 0.000049593 \\
 0. \\
 0.000049593
 \end{array}
 \tag{5.3}$$

We can solve this problem by computing the sum backwards: first adding all the small numbers together and working our way up.

```

> n := 12000:
  SUM := evalf(0):
  while n > 0 do
    SUM := evalf(SUM + 1/(n^2));
    n := n-1;
  end do:
  print(`The sum is` = SUM);
  abs(SUM-S)/abs(S);
  100*abs(SUM-S)/abs(S);

```

The sum is = 1.6449

$$\begin{array}{r}
 0.000060792 \\
 0.0060792
 \end{array}
 \tag{5.4}$$