

Scientific Computing TU Berlin Winter 2021/22 © Jürgen Fuhrmann Notebook 08

```

• begin
•     using PlutoUI
•     using PlutoVista
•     using BenchmarkTools
• end

```

Number representation

Integer numbers
 Floating point numbers
 Decimal representation of $x \in \mathbb{R}$
 Scientific notation
 IEEE754 standard
 Prettyprinting bit representations
 Julia floating point types
 Floating point limits
 BigFloat
 Density of numbers

Number representation

Besides of the concrete names of Julia library functions everything in this chapter is valid for most modern programming languages and computer systems.

All data in computers are stored as sequences of bits. In Julia, for concrete number types, the `bitstring` function returns this information as a sequence of 0 and 1. The `sizeof` function returns the number of bytes in the binary representation.

Integer numbers

```
T_int = Int8
```

```
• T_int=Int8
```

```
i = 1
```

```
• i=T_int(1)
```

```
1
```

```
• sizeof(i)
```

```
"00000001"
```

```
• bitstring(i)
```

Positive integer numbers are represented by their representation in the binary system. For negative numbers n , the binary representation of their "two's complement" $2^N - |n|$ (where N is the number of available bits) is stored. Correspondingly, the sign of an integer number is stored in the first bit.

`typemin` and `typemax` return the smallest and largest numbers which can be represented in number type.

```
(-128, 127)
```

```
• typemin(T_int), typemax(T_int)
```

Unless the possible range of the representation $(-2^{N-1}, 2^{N-1})$ is exceeded, addition, multiplication and subtraction of integers are exact. If it is exceeded, operation results wrap around into the opposite sign region.

10

• 3+7

-119

• `typemax(T_int)+T_int(10)`

3//10

• 1//10 + 2//10

Floating point numbers

0.30000000000000004

• 0.1+0.2

But this should be 0.3. What is happening ???

Decimal representation of $x \in \mathbb{R}$

Usually we write real numbers as decimal fractions and cut the representation off if the number of digits is too large or infinite.

Any real number $x \in \mathbb{R}$ can be expressed via the **representation formula**

$$x = \pm \sum_{i=0}^{\infty} d_i \beta^{-i} \beta^e$$

with **base** $\beta \in \mathbb{N}, \beta \geq 2$, **significand** (or **mantissa**) digits $d_i \in \mathbb{N}, 0 \leq d_i < \beta$ and **exponent** $e \in \mathbb{Z}$.

This representation is infinite for periodic decimal numbers and irrational numbers.

Scientific notation

The scientific notation of real numbers is derived from the representation formula in the case of $\beta = 10$. Let e.g. $x = 6.022 \cdot 10^{23} = 6.022\text{e}23$. Then

- $\beta = 10$
- $d = (6, 0, 2, 2, 0 \dots)$
- $e = 23$

This representation is not unique, e.g. $x_1 = 0.6022 \cdot 10^{24} = 0.6022\text{e}24 = x$ with

- $\beta = 10$
- $d = (0, 6, 0, 2, 2, 0 \dots)$
- $e = 24$

IEEE754 standard

This is the actual standard format for storing floating point numbers. It was developed in the 1980ies. It is derived as well from the representation formula

- $\beta = 2$, therefore $d_i \in \{0, 1\}$
- Parameters:
 - t is the **significand** (mantissa) length
 - k : **exponent size** – number of bits for exponent
- Truncation to fixed size: $x = \pm \sum_{i=0}^{t-1} d_i 2^{-i} 2^e$
- Normalization: assume $d_0 = 1 \Rightarrow$ save one bit for the storage of the significand. This saved bit is also called "hidden bit" This requires a normalization step after operations which adjusts significand and exponent of the result.
- Exponent range: $L := -2^k + 1 \leq e \leq 2^k - 1 =: U$. Actually, $e + 2^{k-1} - 1$ is stored
- Extra bit for sign
- \Rightarrow storage size: $(t - 1) + k + 1$
- The storage sequence is: Sign bit, exponent, mantissa.
- Standardized for most modern languages
- Hardware support usually for 64bit and 32bit

| precision | Julia | C/C++ | k | t | bits |
|-----------|---------|-------------|----|-----|------|
| quadruple | n/a | long double | 16 | 113 | 128 |
| double | Float64 | double | 11 | 53 | 64 |
| single | Float32 | float | 8 | 24 | 32 |
| half | Float16 | n/a | 5 | 11 | 16 |

- See also the [Julia Documentation on floating point numbers](#), [0.30000000000000004.com](#), [wikipedia](#) and the links therein.

Storage layout for a normalized Float32 number ($d_0 = 1$):

- bit 1: sign, $0 \rightarrow +$, $1 \rightarrow -$
- bit 2...9: $k = 8$ exponent bits
 - the value $e + 2^{k-1} - 1 = e + 127$ is stored \Rightarrow no need for sign bit in exponent
- bit 10...32: $23 = t - 1$ mantissa bits $d_1 \dots d_{23}$
- $d_0 = 1$ not stored \equiv "hidden bit"

Prettyprinting bit representations

Julia allows to obtain the significand and the exponent of a floating point number

```
x0 = 100.0
```

```
• x0=100.0
```

```
(1.5625, 6)
```

```
• significand(x0), exponent(x0)
```

Calculate k – the length of the exponent – from the maximum representable floating point number by taking the base-2 logarithm of its exponent:

```
• exponent_length(T::Type{<:AbstractFloat})=Int(log2(exponent(floatmax(T))+1)+1);
```

t – the size of the significand – is calculated from the overall size of the representation minus the size of the exponent and the size of the sign bit + 1 for the "hidden bit".

```
• significand_length(T::Type{<:AbstractFloat})=8*sizeof(T)-exponent_length(T)-1+1;
```

This allows to define a more readable variant of the bitstring representation for floats.

The sign bit is the first bit in the representation:

```
• signbit(x::AbstractFloat)=bitstring(x)[1:1];
```

```
exponent_bits (generic function with 1 method)
  • exponent_bits(x::AbstractFloat)=bitstring(x)[2:exponent_length(typeof(x))+1]
```

```
• significand_bits(x::AbstractFloat)=bitstring(x)
  [exponent_length(typeof(x))+2:8*sizeof(x)];
```

```
• floatbits(x::AbstractFloat)=signbit(x)*"-"*exponent_bits(x)*"-"*significand_bits(x)
);
```

```
T = Float64
• T=Float64
```

- $x = T(0.1)$

- Numbers which are exactly represented in the decimal system may not be exactly represented in the binary system! Such numbers are always rounded to a finite (in the binary system) approximation.

```
x_per = 0.30000000000000000004
```

- $x_{\text{per}} = T(0.1) + T(0.2)$

```
"0_01111111101_00110011001100110011001100110011001100110011001100110100"
```

Finite size of representation \Rightarrow there are minimal and maximal possible numbers which can be represented

Smallest positive denormalized number: $d_i = 0, i = 0 \dots t-2, d_{t-1} = 1 \Rightarrow x_{min} = 2^{1-t}2^L$

[illegible]

Smallest positive normalized number: $d_0 = 1, d_i = 0, i = 1 \dots t-1 \Rightarrow x_{min} = 2^L$

[illegible]

Largest positive normalized number: $d_i = 1, 0 \dots t-1 \Rightarrow x_{max} = 2(1 - 2^{1-t})2^U$

[illegible]

4/7

```
• typemax(T),floatbits(typemax(T)),prevfloat(typemax(T)),  
  floatbits(prevfloat(typemax(T)))
```

- There cannot be more than 2^{t+k} floating point numbers \Rightarrow almost all real numbers have to be approximated
- Let x be an exact value and \tilde{x} be its approximation. Then $|\frac{\tilde{x}-x}{x}| < \epsilon$ is the best accuracy estimate we can get, where
 - $\epsilon = 2^{1-t}$ (truncation)
 - $\epsilon = \frac{1}{2}2^{1-t}$ (rounding)
- Also: ϵ is the smallest representable number such that $1 + \epsilon > 1$.
- Relative errors show up in particular when
 - subtracting two close numbers
 - adding smaller numbers to larger ones

E.g. Addition

- The smallest number one can add to 1 can have at most t bit shifts of normalized mantissa until mantissa becomes 0, so its value must be 2^{-t} .

- Smallest floating point number ϵ such that $1 + \epsilon > 1$ in floating point arithmetic
- In exact math it is true that from $1 + \epsilon = 1$ it follows that $0 + \epsilon = 0$ and vice versa. In floating point computations this is not true

- $\epsilon = \text{eps}(T)$

- `floatbits(ϵ)`

```
• one(T)+ $\epsilon/2$ , floatbits(one(T)+ $\epsilon/2$ ), floatbits(one(T))
```

- `one(T)+ ϵ , floatbits(one(T)+ ϵ)`

- `nextfloat(one(T))-one(T), floatbits(nextfloat(one(T))-one(T))`

- the `BigFloat` type wraps the GNU multiprecision library and allows for floating point operations with high accuracy
- one can control its precision
- However it cannot be exact: periodic dual representations must be cut off nevertheless...
- Here, we just use the default precision

- `eps(BigFloat)`

- `abig=big"0.1"; bbig=big"0.2"`

- `abig+bbig`

- **a64=0.1; b64=0.2;**

- @benchmark \$a64+\$b64

- `@benchmark $abig+$bbig`

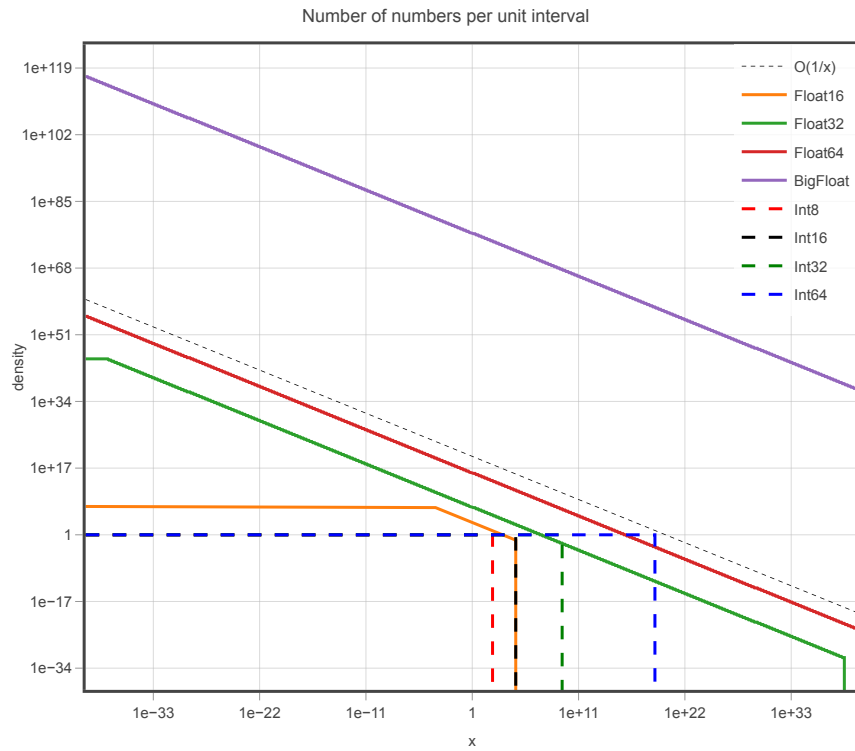
How dense are numbers on the real axis?

```

1 function dens(x, ::Type{T}; sample_size=1000) where T<:AbstractFloat
2     xleft=T(x)
3     xright=T(x)
4     for i=1:sample_size
5         xleft=prevfloat(xleft)
6         xright=nextfloat(xright)
7     end
8     return Float64(prevfloat(2.0*sample_size/(xright-xleft)))
9 end;

```

```
• dens(x, ::Type{T}) where T<:Integer = x < typemax(T) ? 1 : 0 ;
```



```

• let
•     X=10.0.^collect(-40:0.1:40)
•     p=plot(resolution=(600,500),
•           title="Number of numbers per unit interval",
•           xscale=:log,yscale=:log,xlabel="x",ylabel="density",legend=:rt)
•     plot!(p,X,map(x->1.0e20/x,X),
•           label="O(1/x)",linewidth=0.5,color=:black,linestyle=:dot)
•     plot!(p,X,dens.(X,Float16),label="Float16")
•     plot!(p,X,dens.(X,Float32),label="Float32")
•     plot!(p,X,dens.(X,Float64),label="Float64")
•     plot!(p,X,dens.(X,BigFloat),label="BigFloat")
•     plot!(p,X,dens.(X,Int8),label="Int8",linestyle=:dash,color=:red)
•     plot!(p,X,dens.(X,Int16),label="Int16",linestyle=:dash,color=:black)
•     plot!(p,X,dens.(X,Int32),label="Int32",linestyle=:dash,color=:green)
•     plot!(p,X,dens.(X,Int64),label="Int64",linestyle=:dash,color=:blue)
• end

```