

- 
- ° Floating point representation of numbers

## Representing numbers.

The real number system.  $x = \pm (d_0.d_1d_2d_3\ldots)_\beta \times \beta^E$  or base  $\beta$  rep.

Ex:  $x = 3.1415926538\ldots$

$$= + \left( \frac{3}{10^0} + \frac{1}{10^1} + \frac{4}{10^2} + \ldots \right) \times 10^{+0}$$

$$= + (3.14159265\ldots)_{10} \times 10^{+0}$$

$$= + S \times 10^{+0}, \text{ with } 1 \leq S < 10$$

The floating-point number system.

A floating-point number  $x$  has the form

$$x = \pm S_\beta \times \beta^e = \pm \left( \frac{d_0}{\beta^0} + \frac{d_1}{\beta^1} + \ldots + \frac{d_{p-1}}{\beta^{p-1}} \right) \times \beta^e, \text{ with } e_{\min} \leq e \leq e_{\max}.$$

$\underbrace{\hspace{10em}}_{p \text{ digits}}$

The set  $F$  of floating-point numbers ~~is~~ is defined by

base  $\beta$

floating-point precision  $p$

exponent range  $e_{\min} \leq e \leq e_{\max}$

Binary system:  $\beta = 2$

$$x = \pm (1.b_1b_2\ldots b_{p-1})_2 \times 2^e, \text{ where digits } \{b_i\} \text{ satisfy } b_i \in \{0, 1\}.$$

normalized rep

For uniqueness of representation

$$m = 1 + \frac{b_1}{2} + \frac{b_2}{2^2} + \ldots + \frac{b_{p-1}}{2^{p-1}} \text{ is called the mantissa.}$$

The spacing of numbers in  $\mathbb{F}$ :  
Consider the numbers on either side of 1.

1.8-4

$$1 = (1.00 \dots 0) \times 2^0$$

the next number in  $\mathbb{F}(p, m)$  is

$$(1.0 \dots 0 \underset{\substack{\uparrow \\ 1/2^{p-1}}}{1})_2 \times 2^0 = 1 + \frac{1}{2^{p-1}} \quad \text{let } \epsilon = \frac{1}{2^{p-1}}$$

the next number is

$$(1.00 \dots 0 \underset{\substack{\uparrow \\ 1/2^{p-2}}}{01})_2 \times 2^0 = 1 + \frac{1}{2^{p-2}} + \frac{1}{2^{p-1}} = 1 + 2\epsilon$$

so the numbers following 1 are

$$1, 1+\epsilon, 1+2\epsilon, \dots, 1+k\epsilon.$$

Notice:  $1+k\epsilon = (1.1 \dots 1)_2 \times 2^0 = 2 - \frac{1}{2^{p-1}} = 2 - \epsilon$

Notice:  $1+k\epsilon = (1.1 \dots 1)_2 \times 2^0 = 2 - \frac{1}{2^{p-1}}$   
 $k\epsilon = 1 - \frac{1}{2^{p-1}} = \frac{1}{2^{p-1}} (2^{p-1} - 1) \Rightarrow k = 2^{p-1} - 1$

The numbers preceding 1 are:

$$(1, 1+\epsilon, 1+2\epsilon, \dots, 1+k\epsilon) \times 2^{-1} = 1 - (2 - \frac{1}{2^{p-1}}) \cdot \frac{1}{2} = 1 - (1 - \frac{1}{2^p}) = \frac{\epsilon}{2}$$

In particular,  $1 - (1+k\epsilon) \times \frac{1}{2} = 1 - (2 - \frac{1}{2^{p-1}}) \cdot \frac{1}{2} = 1 - (1 - \frac{1}{2^p}) = \frac{\epsilon}{2}$

The spacing between numbers

$$(1, 1+\epsilon, \dots, 1+k\epsilon) \times 2^e \quad \text{is} \quad 2^e \times \epsilon = \frac{2^e}{2^{p-1}}$$

So, some problems arise when small numbers are added to big ones.

double precision:  $\epsilon_p = \frac{1}{2^{52}}$ , called "machine  $\epsilon$ "

~~1.5~~  
1.5-2

$\Rightarrow$  max command: ~~44~~ EPS

$\epsilon_{\max} = 10^{23}$ , so the spacing of numbers from

$1 \times 2^{55}, \dots, \cancel{X \times (2 - \frac{1}{2^{52}})} \times 10^{23}$  is greater than one

$\Rightarrow$  If you need to ~~represent~~ ~~float~~ integers (and only integers) more ~~precisely~~ exactly, you should use the integer class.

$\Rightarrow$  There are also extended arithmetic packages available.

$\Rightarrow$  common uses: cryptography algorithms, very large primes + factorizations.

## Rounding

The absolute error between  $x$  and ~~the~~  $m_f = fl(x)$ , its floating-pt representation is  $|x - m_f|$

The relative error is  $\frac{|x - m_f|}{|x|}$ .

Using the "round to nearest" rule:  $\frac{|x - m_f|}{|x|} \leq \frac{\epsilon_p}{2}$ .

Other "rounding standards" include

"round up"

"round down"

~~NaN~~

Notice:  $\infty$  is not represented in  $\mathbb{F}(p, M)$ :

it is added to the set of floating-pt numbers.

Other additional numbers:

$\pm \infty$ : overflow:  $x > z_{max}$

NaN: not a number: IF  $|x| < z_{min}$ ,  $fl(x) = 0$   
and any ratio of the form  $0/0$  is a NaN.

Ex:

$$fl \left[ \frac{(z^{54} - 1) - z^{54}}{(z^{54} - 0.5) - z^{54}} \right] = \text{NaN} \quad \text{exact: } z$$

What if you replace  $z^{54}$  with  $z^{53}$ ?

Ex:

$$fl \left[ \frac{(z^{54} - z) - z^{54}}{(z^{54} - 1) - z^{54}} \right] = -\text{Inf} \quad \text{exact: } z$$

Ex:

$$fl \left[ \frac{(z^{54} - 1) - z^{54}}{(z^{54} - z) - z^{54}} \right] = 0 \quad \text{exact: } \frac{1}{z}$$

— If you are not aware of these issues, they can cause problems...

Ex:  $\sin(k\pi) = 0$  for all integers  $k$ .

BUT:  $\sin(z^{53})$  in Matlab  $(-0.8489\dots)$

— Try and predict what Matlab gives for:

$$(z^{54} - z) - z^{54}$$

$$(z^{54} + z) - z^{54}$$

Why??