## Representation of Numbers:

### Integer Numbers:

N bits: $2^N$ integers

$2^{N-1}$

$$I = (-1)^S (\alpha_n 2^n + \alpha_{n-1} 2^{n-1} + \ldots + \alpha_0 2^0) \quad \Leftarrow$$

$$S = \begin{cases} 0 \to + \\ 1 \to - \end{cases} \quad n = N-2 \quad \alpha_i = \begin{cases} 0 \\ 1 \end{cases}$$

single precision: 32 bits

$$\Rightarrow 2^{31} \sim 2 \times 10^9$$

a) the answer respects the range

b) division is interpreted as the integer part

## Real Numbers

$$r = (-1)^S \underbrace{(m_1 2^{-1} + m_2 2^{-2} + \ldots + m_{23} 2^{-23})}_{23 \text{ bits}} 2^{\overset{\longrightarrow}{\alpha - 127}} \quad \leftarrow \text{"bias"}$$

1 bit

8 bits $\to \alpha \in [0, 255]$

$\Downarrow$ $\qquad\qquad$ $\Downarrow$

$2^{-23} \sim 10^{-7}$ $\qquad$ $2^{-127} \to 2^{128}$

real numbers: $10^{-44} \underset{\sim}{\leq} r \underset{\sim}{\leq} 10^{+38}$

double precision: $10^{-12}$

Machine Precision: $(\varepsilon)$

for 32-bits: $\varepsilon = 10^{-7}$   for 64 bits: $\varepsilon = 10^{-16}$

smallest #:

32 bits

$\# > 2^{128}$   OVERFLOW

$\# < 2^{-127}$   UNDERFLOW

Roundoff Error:

$1 + \varepsilon + \varepsilon + \varepsilon \ldots = 2$ but 1

$\underbrace{\qquad}_{10^{-7} \text{ times}}$

$\propto \sqrt{N} \, \varepsilon_m$     $\propto N$

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

Truncation Error:

$$f(x) = \sum_{n=0}^{\infty} a_n x^n \quad \Rightarrow \quad f_N(x) = \sum_{n=0}^{N} a_n x^n$$

derivative

$$f'(x) = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$

as $h \to 0$    $x+h \sim x$ in compute