

Floating point

$$\pm 2^w \cdot (1+M)$$

w - wykładnik

M - mantyza $(0,1)$

Wykładnik

Zapis biased na n bitach: $w + 2^{n-1}$

dla 8 bitów (float):

$$w \in [-126, 127] \quad w + 127 \text{ w pamięci}$$

Specjalne wartości

0 - wykładnik i mantyza to same 0 (± 0)

∞ - wykładnik to same 1, mantyza same 0 ($\pm \infty$)

NaN - wykładnik same 1, niezerowa mantyza

Denormalized numbers

Wykładnik same 0, dowolna mantyza: $0.x \dots x \cdot 2^{-126}$

Błędy

bezwzględny $|x - x^*| \quad \pm \epsilon$

względny $\left| \frac{x - x^*}{x} \right| \quad \pm 1\%$

Catastrophic cancellation - odejmowanie bliskich liczb

- alternatywny sposób obliczania wyniku
- sprzężenie

$$x^2 = 2 \Rightarrow x^* = 1.4$$

$$\text{forward error: } |\sqrt{2} - 1.4|$$

$$\text{backward error: } |1.4^2 - 2|$$

Uwarunkowanie problemu

małe zaburzenie wejścia \rightarrow małe zaburzenie wyniku

złe uwarunkowany problem: pierwiastki wielomianu

Stabilność jest cechą algorytmu