

# BŁĘDY OBLICZEŃ NUMERYCZNYCH

- błędy zaokrągleń  
skończona liczba cyfr (bitów) w reprezentacji numerycznej
- błędy obcięcia  
rozwinienia w szeregi i procesy iteracyjne -  
w praktyce muszą być skończone
- błędy metody  
przybliżone rozwiązania problemu

poza tym:

- \* błędy modelu
- \* błędy danych wejściowych
- \* pomyłki

## Zaokrąglenia

1. nieskończone reprezentacje dziesiętne
2. nieskończone reprezentacje binarne  
skończonych reprezentacji dziesiętnych, np.:

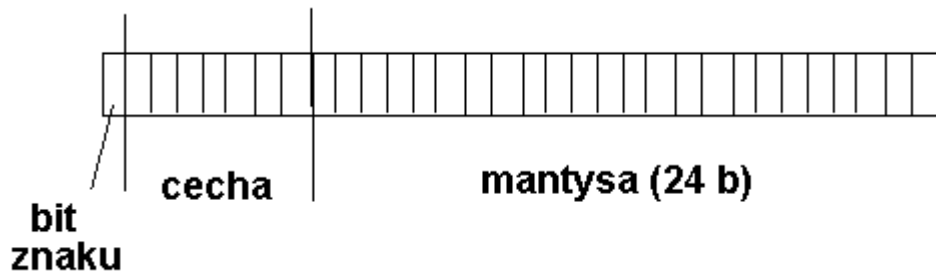
$$(0.1)_{10} = (0.000110011001100110011001100 \dots)_2$$
$$0.09999999 \dots$$

$$(0.000110011)_2 \rightarrow 0.099609375$$

przykład:

# zmiennoprzecinkowa reprezentacja w IBM 370 (najlepsza maszyna obliczeniowa w latach 70-tych)

## 32-bitowe słowo



podstawa systemu liczenia - 16

rozpatrzmy liczbę, której reprezentacja ma postać

**0 100010 101100110000010000000000**

wykładnik podstawy:  $[(100010)_2 = 66] - 64$

(cecha bez znaku – odejmowanie 64 zapewnia zakres wartości cechy:

od -64 do +63

wartość (uwzgl. postać znormalizowaną)

$$L = m \cdot p^C$$

$$[(1/2)^1 + (1/2)^3 + (1/2)^4 + (1/2)^7 + (1/2)^8 + (1/2)^{14}] \cdot 16^{(66-64)} = 179.015625$$

minimalnie większa wartość różniąca się 1 bitem

**0 100010 101100110000010000000001**

**179.0156097412109375**

minimalnie mniejsza wartość różniąca się 1 bitem:

**0 100010 101100110000001111111111**

**179.0156402587890625**

**porównajmy:**

**179.0156097412109375**

**179.015625**

**179.0156402587890625**

**1) oznacza to, że reprezentacja**

**0 1000010 101100110000010000000000**

**odpowiada**

**nieskończonej „ilości” liczb rzeczywistych z połowy  
powyższego przedziału**

**2) maksymalnie 7 znaczących cyfr 10-tnych**

**przy 24 znaczących bitach reprezent. maszynowej**

**3) największa możliwa liczba**

**0 1111111 11111111111111111111111111111111  $\approx 16^{63} \approx 10^{76}$**

**najmniejsza**

**0 0000000 0001000000000000000000000000 =  $16^{-65} \approx 10^{-78}$**

**(najmniejsza możliwa mantysa znormalizowana =  $1/16 = 16^{-1}$  ,**

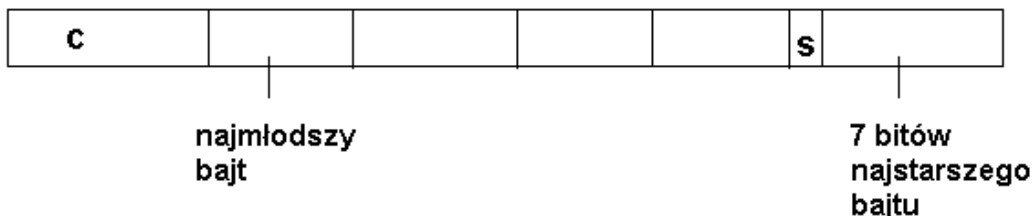
**najmniejsza możliwa cecha)**

**rozważmy bardziej przyziemny przykład:**

**reprezentacja REAL w PASCALU (PC – słowo jednobajtowe)**



w PAO, reprezentacja wygląda inaczej gdyż z pamięci kolejne BAJTY sączytywane są sekwencyjnie (od lewej) i „ustawiane” w rejestrze „poprawnie” – od prawej ... zatem:



program REAL\_BIN – reprezentacja REAL w PAO (..dla +/- 0.5)

program MAX\_REAL – znajdowanie bitu znaku

## utrata dokładności

- operacja dodawania wymaga wyrównania cech

(wykonajmy dodawanie na cyfrach, a nie na bitach i założmy że mamy tylko 7 dostępnych cyfr na mantysę)

$$0.1460071 \times 10^3 \quad \text{--->} \quad 0.0000015 \times 10^8$$

$$+ \quad 0.2389034 \times 10^8$$

a próbując dodać np.  $0.1460071 \times 10^{-1}$  do  $0.2389034 \times 10^8$  w wyniku otrzymamy  $0.2389034 \times 10^8$

a co jeśli do  $0.2389034 \times 10^8$  mielibyśmy dodać 100 milionów takich  $0.1460071 \times 10^{-1}$  liczb?

dodając je sukcesywnie, wciąż mielibyśmy w wyniku  $0.2389034 \times 10^8$  ale dodając w pierw do siebie 100 milionów ( $10^8$ ) razy  $0.1460071 \times 10^{-1}$  dostaniemy  $0.1460071 \times 10^7$

a następnie dodając to do  $0.2389034 \times 10^8$  dostaniemy  $0.2535041 \times 10^8$

- odejmowanie liczb bliskich  $(a - b) * c$

$$a = 0.2356073 \times 10^6$$

$$b = 0.2356102 \times 10^6$$

$$- \quad 0.2900000 \times 10^1 \text{ (bo zawsze postać znormalizowana)}$$

$$c = 0.1112050 \times 10^1$$

wynik

$$0.3224945 \times 10^1$$

| śmieci |

mówimy, że  $p^*$  jest zmiennoprzecinkowym przybliżeniem do  $p$  z dokładnością do  $t$  cyfr znaczących jeśli  $t$  jest największą nieujemną liczbą całkowitą taką, że

$$\frac{|p - p^*|}{|p|} < 5 \times 10^{-t}$$

### Przykład

rozwiązanie równania

$$\frac{d^2 y(x)}{dx^2} = c^2 y(x),$$

$C$  - rzeczywiste

np. równanie Schrödingera dla cząstki wewnątrz bariery

$$y(x) = A e^{-Cx} + B e^{Cx}$$

żądamy rozwiązania, które  $y(0)=1$ ,  $y'(0)=-C$

$$\Rightarrow y(x) = e^{-Cx}$$

jeśli chcemy numerycznie znaleźć to rozwiązanie np. w przedziale  $[0,P]$ ,

czyli dla wielu arg.  $x$ , to musimy nadać wartość szukanego rozwiązania w  $x=0$

a następnie „znajdować” numeryczne  $y(x)$  w kolejnych „węzłach” dyskretnej siatki zmiennej  $x$  w przedziale  $[0,P]$ ; czyli -

zadajemy  $1$  oraz  $C$  w  $x=0$  i numerycznie szukamy wartości  $Y$  dla  $x>0$

**rozwiązując numerycznie na dyskretnej siatce węzłów:**

$1$  (jedyńka) oraz  $C$

są określone (zaokrąglone) z numeryczną dokładnością  $\varepsilon$

$$\Rightarrow y_{\text{numeryczne}} \approx e^{-Cx} + \varepsilon e^{Cx}$$

już w  $X=0$  będzie zawierało domieszkę rozwiązania

$e^{Cx}$  (z wagą rzędu  $\varepsilon$  - rzędu dokładności numerycznej )

dla dużych i dodatnich  $X$  "eksplodujące" rozwiązanie  $e^{Cx}$  zdominuje

numeryczne rozwiązanie  $y_{\text{num}}$

### błędy obcięcia

**pewne wielkości wyliczane są zawsze z nieskończonych sum, szeregów czy iteracji**

**np:**

$$\sin(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots$$

przykład unikania utraty dokładności

$$f(x) = x - \sin(x)$$

dla  $x \approx 0$   $x \approx \sin(x)$  (bliskie wartości)

$$\implies f(x) = \frac{x^3}{3!} - \frac{x^5}{5!} + \frac{x^7}{7!} - \dots$$

.....

numeryczne sumowanie szeregu  $1/(n^3)$

- w przód i w tył

...program... suma\_szeregu

**METODY NUMERYCZNE =**

sposoby przybliżonego rozwiązywania zagadnień

- dyskretyzacja
- aproksymacja
- propagacja błędów i numeryczna niestabilność

np. przykład (prawo rozpadu promieniotwórczego)



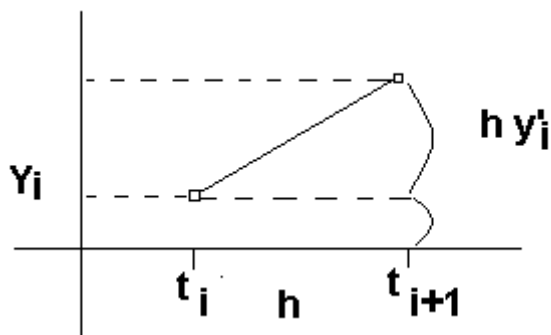
rozwiązanie równania  $\frac{dY(t)}{dt} = -\lambda Y(t)$

metodą Eulera (poznamy później)

rozwiązaniem ścisłym jest  $y(t) = \exp(-\lambda t)$

(dla  $\lambda > 0$  wartości  $Y(t)$  powinny monotonicznie maleć w czasie)

na siatce punktów  $\{t_i\}$ ,  $t_{i+1} - t_i = h$



$$Y_{i+1} = Y_i + h Y_i' = (1 - \lambda h) Y_i$$

metoda jest niestabilna

jeśli

$$h > 2 / \lambda$$

...program EULER\_GR ....

problemy źle uwarunkowane

(przy zadanej dokładności numerycznej)

np.:

układ równań

$$6.025 x + 109.774 y = 4.12$$

$$14.054 x + 256.125 y = 10.50$$

ma rozwiązanie:

$$x = -250.15 \quad y = 13.77 ;$$

jeśli dysponujemy dokładnością do 2 miejsc po kropce i wyniki pośrednich operacji zaokrąglamy do 2 miejsc po kropce, to rozwiązaniem jest

$$x = -44.33 \quad y = 2.47$$

powód:

równania "prawie" zależne (niezależne na granicy dokładności)

