

## Laboratorium 02

### Temat: Błędy metod numerycznych - źródła i rodzaje.

Celem ćwiczenia jest zapoznanie się ze źródłami błędów metod numerycznych oraz ich rodzajami. Podstawowy problem dotyczący dokładności obliczeń z wykorzystaniem metod numerycznych wynika z konieczności zapisu liczb w postaci ich reprezentacji komputerowej. Reprezentacja komputerowa liczb zależy przede wszystkim od ich typu:

- liczby całkowite, np. 1, 3, itp: są reprezentowane w postaci zapisu binarnego na 32 lub 64 bitach,
- liczby zmiennoprzecinkowe, np. 1.10, 3.14: są reprezentowane z wykorzystaniem trzech liczb binarnych zapisanych na 64 bitach,
- liczby zespolone: składają się z dwóch liczb zmiennoprzecinkowych.

Liczby w postaci cyfrowej są kodowane z wykorzystaniem dwóch napięć, tj. niskiego (np. 0V) i wysokiego (np. 5V). Zatem do dyspozycji mamy tylko dwie liczby tj. 0 i 1, które określamy nazwą bitu (z ang. bit to **binary digit**). Ze względu na to, że ludzie posługują się systemem dziesiętnym (10 cyfr), a w komputerze posługujemy się systemem dwójkowym (2 liczby), istnieje konieczność zamiany liczb pomiędzy systemami. Zamiana liczb z systemu dwójkowego na dziesiętny jest stosunkowo prosta i wymaga jedynie sumowania kolejnych potęg liczby 2 z odpowiednim współczynnikiem odpowiadającym wartości poszczególnych bitów

$$\sum_{i=0}^n b_i 2^i$$

gdzie współczynniki  $b_i$  są kolejnymi liczbami zapisu dwójkowego, np.

1	0	0	1
---	---	---	---

$$1001 \Rightarrow 1 \cdot 2^0 + 0 \cdot 2^1 + 0 \cdot 2^2 + 1 \cdot 2^3 = 1 + 0 + 0 + 8 = 9$$

W tabeli przedstawiono wartości liczb dla systemu dziesiętnego oraz ich odpowiedniki w systemie dwójkowym.

Binary	1	10	11	100	101	110	111	1000	1001
Decimal	1	2	3	4	5	6	7	8	9

O ile zapis liczb całkowitych jest stosunkowo prosty to zapis liczb zmiennoprzecinkowych sprawia pewne problemy i wymaga zastosowania znaku rozdzielającego części przed i po przecinku:

$$b_n \dots b_0 . b_{-1} b_{-2} \dots$$

gdzie  $m < n$ . Ponieważ liczby w komputerze są przechowywane zazwyczaj w komórkach o określonej wielkości, to istnieją pewne ograniczenia na ich wielkość, jak i dokładność obliczeń. Informację o właściwościach liczb dostępnych w języku Python w wersji 2.7 można uzyskać korzystając z modułu „sys” oraz odpowiednich metod:

- liczby całkowite (int): sys.maxint

```
>>> import sys
>>> print sys.maxint
9223372036854775807
```

- liczby zmiennoprzecinkowe (float): sys.float\_info

```
>>> import sys
>>> print sys.float_info
sys.float_info(max=1.7976931348623157e+308, max_exp=1024, max_10_exp=308, min=
2.2250738585072014e-308, min_exp=-1021, min_10_exp=-307, dig=15, mant_dig=53,
epsilon=2.220446049250313e-16, radix=2, rounds=1)
```

w przypadku Pythona w wersji 3.6:

- liczby całkowite (int): sys.int\_info

```
>>> import sys
>>> print(sys.int_info)
sys.int_info(bits_per_digit=30, sizeof_digit=4)
```

- liczby zmiennoprzecinkowe (float): sys.float\_info

```
>>> import sys
>>> print(sys.float_info)
sys.float_info(max=1.7976931348623157e+308, max_exp=1024, max_10_exp=308, min=
2.2250738585072014e-308, min_exp=-1021, min_10_exp=-307, dig=15, mant_dig=53,
epsilon=2.220446049250313e-16, radix=2, rounds=1)
```

W przypadku liczb całkowitych ich maksymalną liczbę całkowitą można wyznaczyć z zależności:

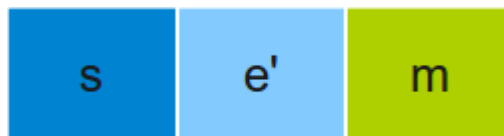
$$\text{liczba maksymalna} = 2^{\text{liczba bitów}} - 1$$

$$2^{63} - 1 = 9223372036854775807$$

gdzie liczba bitów wynosi 64.

Zapis liczb zmiennoprzecinkowych wg standardu IEEE-754 wymaga trzech pól, tzn.:

- bitu określającego znak liczby (sign)
- pola bitów przechowujących wykładnik (exponent)
- pola bitów przechowujących mantysę (mantissa)



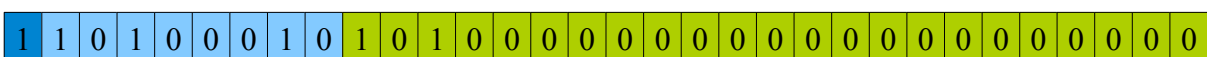
Zamiana zawartości pól na zmiennoprzecinkową liczbę binarną odbywa się wg formuły:

$$\text{liczba} = (-1)^s \times (1.m) \times 2^{e'-127}$$

Standard IEEE-754 przewiduje między innymi możliwość zapisu liczb zmiennoprzecinkowych z wykorzystaniem:

- 32 bitów [s=>1 bit, e'=>8 bitów, m=>23 bity]
- 64 bitów [s=>1 bit, e'=>11 bitów, m=>52 bity]

np.





gdzie funkcja `Decimal()` z biblioteki `decimal` zwraca dokładną wartość liczby typu `float` przechowywaną w pamięci.

### Ćwiczenia:

1. Napisz program do zamiany liczb z systemu dwójkowego na system dziesiętny, przyjmując założenia, że korzystamy jedynie z komórki pamięci o wielkości 1 bajta tj. 8 bitów.
2. Sprawdź i wyjaśnij zachowanie języka Python w wersji 2.7 oraz 3.6 dla przypadku wykonania następujących operacji:
  - dodanie 1 do wartości maksymalnej typu „int”,
  - dwukrotne dodanie wartości maksymalnej typu „float”.
3. Wykonaj wykres zależności błędu tzw. machine epsilon w zależności od ilości bitów przeznaczonych na zapis mantysy, przyjmując założenie, że liczba bitów dla wykładnika wynosi 8, a liczba bitów dla mantysy jest zmienna i może przyjmować wartości od 8 do 52.



4. Zaproponuj / wyszukaj algorytm, który pozwoli uniknąć błędu obliczeń związanego z tzw. machine epsilon, np. w przypadku banków, itp.
5. Oszacuje ilość liczb typu `float` pomiędzy w zakresie 0-1.

### Uwagi:

- funkcja `type(zmienna)` – zwraca/wyświetla typ zmiennej, np.

*Program:*

```
a=1
print type(a) lub print(type(a))
```

*Wynik:*

```
<type 'int'>
```

- zmiana skali liniowej na logarytmiczną

*Program:*

```
...
plt.yscale("log")
```

- suma zbioru liczb bez udziału błędu, np. => `math.fsum`

```
1 import math
2 suma=0
3 sumb=0
4 zbior=[]
5 for i in range(0,1000000):
6     suma=suma+0.1
7     zbior+=0.1
8 sumb=math.fsum(zbior)
9 print(suma)
10 print(sumb)
```

Shell

```
>>> %Run lab02c.py
100000.00000133288
100000.0
```

```
1 import math
2 suma=0
3 sumb=0
4 zbior=[]
5 for i in range(0,1000000):
6     suma=suma+0.1
7     zbior+=0.1
8 sumb=math.fsum(zbior)
9 print(suma)
10 print(sumb)
```

Shell

```
>>> %Run lab02c.py
100000.00000133288
100000.0
```