

Uvod

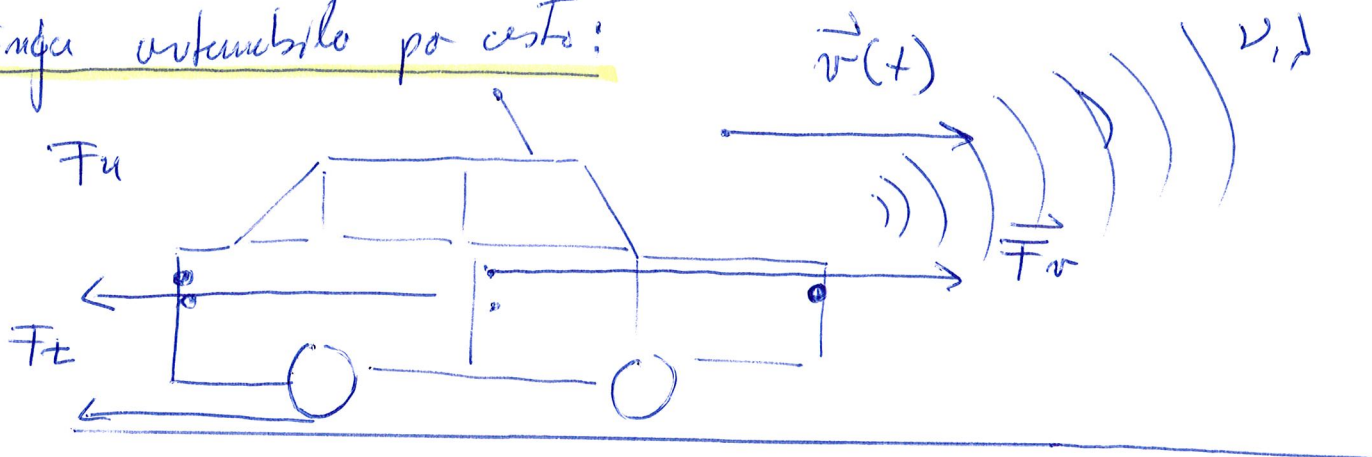
14. 2. 2023



MATEMATIČNO FIZIKALNI SEMINAR

Cilji predmeta: Razumevanje in uporaba numeričnih metod v fizikalnih izračunih.

Vožnja avtomobila po cesti:



Šofer pozna $v(t)$: Od tega lahko izračunamo prevoženo pot in prepušča z integracijo in odvajanjem hitrosti po času:

$$\vec{a}(t) = \frac{d\vec{v}}{dt} = \dot{\vec{v}} \quad \text{---}$$

To navedeno znamo vedno izračunati, če le poznamo funkcije odvisnost.

$$\vec{s}(t) = \int_0^t \vec{v}(t') dt' \quad \text{---}$$

Teh integralov pa ne znamo vedno izračunati analitično.

Tedaj se poslužimo numerične integracije.

Hiterost homogeno le v določeni točki: \Rightarrow Amplitudno
rodimo miso možno, saj se poznamo prave funkcije
odvisnosti \Rightarrow Za izračun pospeškov in hitrosti
se zateemo numerična odčitavanja in integracije

Poznamo sile, ki delujejo na objektu:

$$x: \sum_i F_i = m \cdot a = F_v - F_u - F_t$$

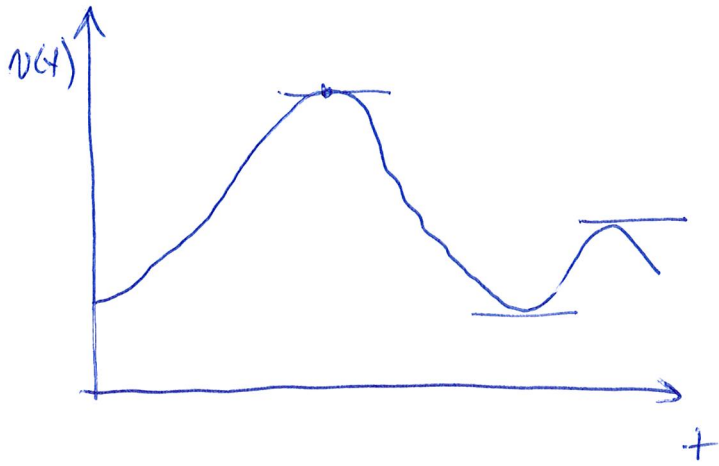
$$m \cdot a = F_v(t) - c \beta z \cdot \frac{S \cdot v^2}{2} - mg \cdot k_t$$

$$m \ddot{x}(t) = F_v(t) - \frac{c \beta z S}{2} \dot{x}(t)^2 - mg k_t$$

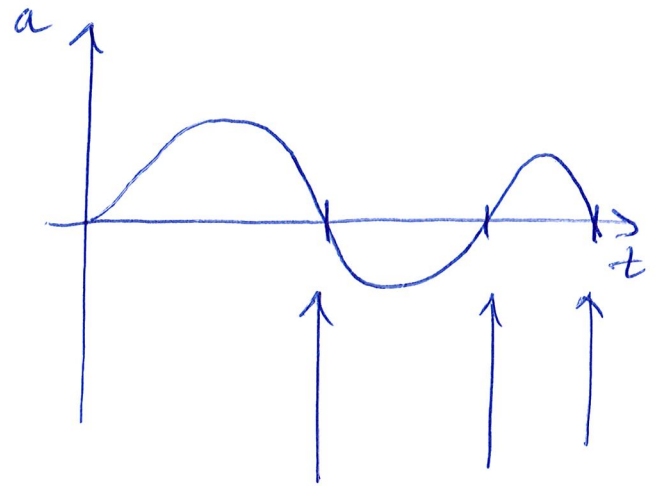
Zapletene diferencialne enačbe se z mano rosti
analitično \Rightarrow Numerična reševanja diferencialnih
enačb.

↓
Reševanje parcialnih diferencialnih
enačb. (npr. Difuzijska enačba)
↓
Reševanje diferencialnih enačb v res
dimenzijah.

Kolej auto zadne zawrato \Rightarrow klasy model popycha



\Rightarrow



Ishkanije ekstremov
funkcij

Numerično iskanje
model funkcij

Automobil med različnimi oddaljenostmi s frekenco ν_0 .
 Mi zaradi Dopplerjevega efekta slišimo zvočni
 pri drugodnih frekvencah, glede na to ali se
 manj čisto približuje ali oddaljuje.

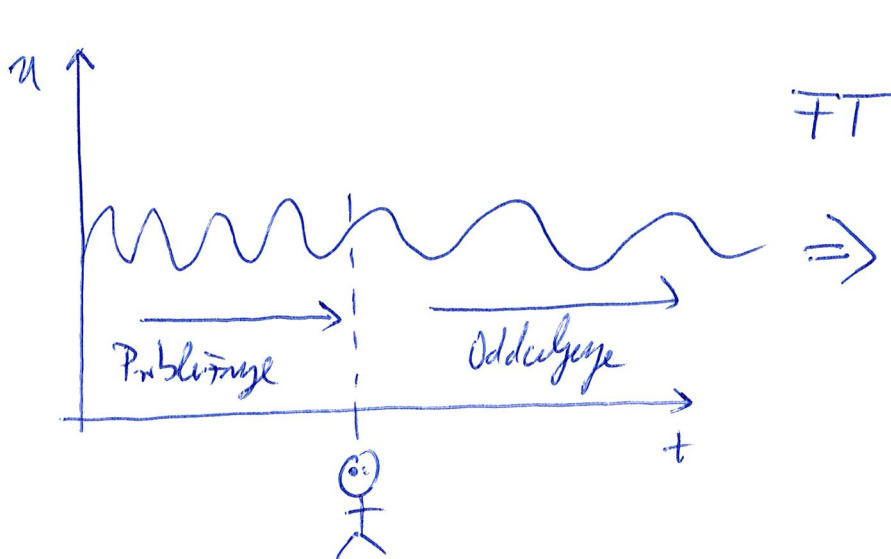
Zvočni:
$$u(t) = u_0 \cdot \sin(2\pi \nu t)$$

$$\nu = \nu_0 \cdot \left(1 \pm \frac{v}{c}\right)^{-1}$$

Če se mi približuje.

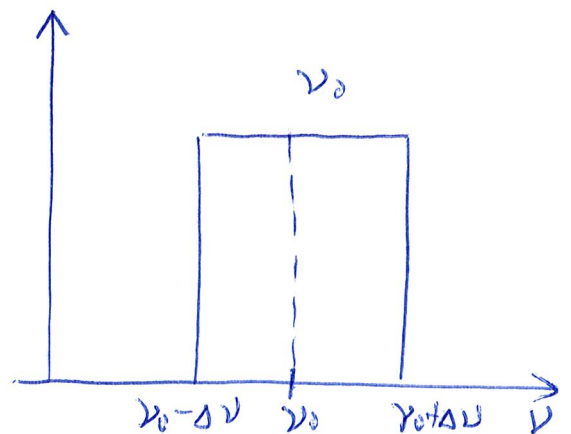
↓
 Približuje
 ali oddaljuje.

Kako analizirati zvočni avtomobila \Rightarrow Fourierova analiza



Osnovna odzivnost
 signala

FT \Rightarrow



Spekter signala

Daljšanje nuj hitrost.

Znova imamo apračati določeno urovo
(pametno raba!) \Rightarrow Numerične Fourierove
transformacije in lube Fourierove transformacije

* Dodaten motivacijski zglede \longrightarrow

Kdaj uporabljamo numerične metode?

— Ko določa matrico se postavlja:
 \longrightarrow kot bazo modula $\int_0^1 e^{x^2} dx$

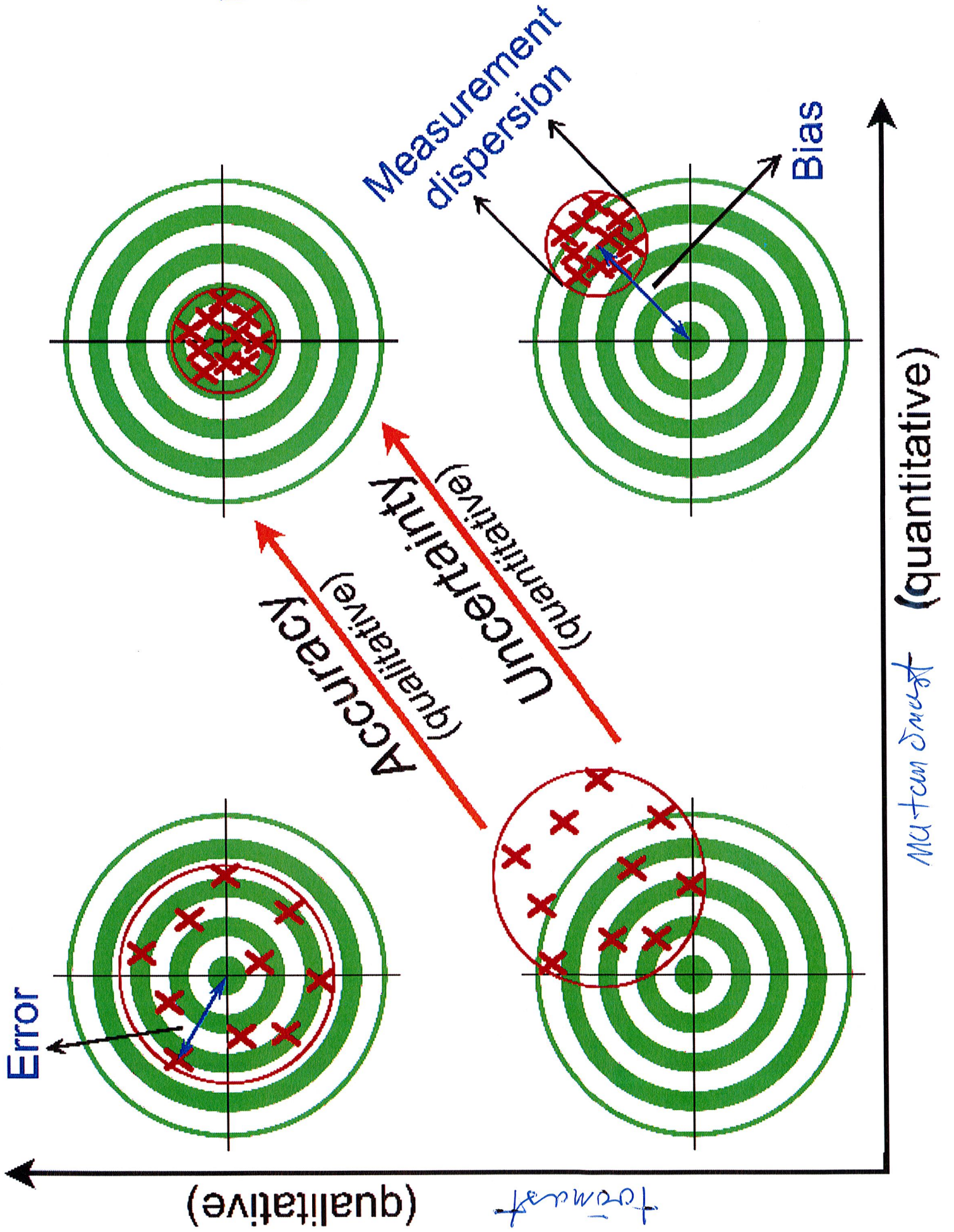
— kadar je to udobnejše kot računati manj
zahtevno od iskanja analitičnih rešitev

Od dobrih numeričnih metod pričakujemo:

- EKONOMIČNOST: časovno (številne operacije)
prostorsko (paralelna izvedba)
- UNIVERZALNOST: da odpravi le pri redkih izjemih
- NUMERIČNA STABILNOST: dim naj bo numerične
napake.

Prav potrditveno se o tem ne bomo ukvarjali.
Uporabljamo bomo že izdelane algoritme.

Tocnost vs. natančnost



Od rezultata, ko jih dobimo z numerično analizo (boljše modelske napovedi boljši analiza podatkov), pričakujemo, da so karkoli točnejši, nato pa tudi analitičnejši nastanejo.

ZAPIS ŠTEVIL V BIMARNI BAZI

V računalniku so vsi števila predstavljena z zaporedjem bitov 0,1. Števila predstavljamo

z navedeno dolžino 32,64 bitov (navede se bitovi ki niso krajši, lahko pa so tudi daljši)

"char" = 8 bit
"int" = 32 bit
"long" = 64 bit

Cela števila: (tip "long")



↑
Predznak
(sign)

Največje število: $2^{63} - 1$ (za vsa mesta)

Pohablje Matematičar, kaj je to
so eno število!

Najmanjše število: -2^{63}



Cela števila niso problemi, so pa omejena.

Zapis steru:

$$X = (-1)^s \sum_{i=0}^{63} b_i \cdot 2^i$$

Zglad: Pretworu steru
w deseti5ku system:

0000000...0101010
↑
s
64 bitow
6 5 4 3 2 1 0

$$X = (-1)^0 \cdot [2^5 + 2^3 + 2] = 32 + 8 + 2 = \underline{\underline{42}}$$

Realna števila

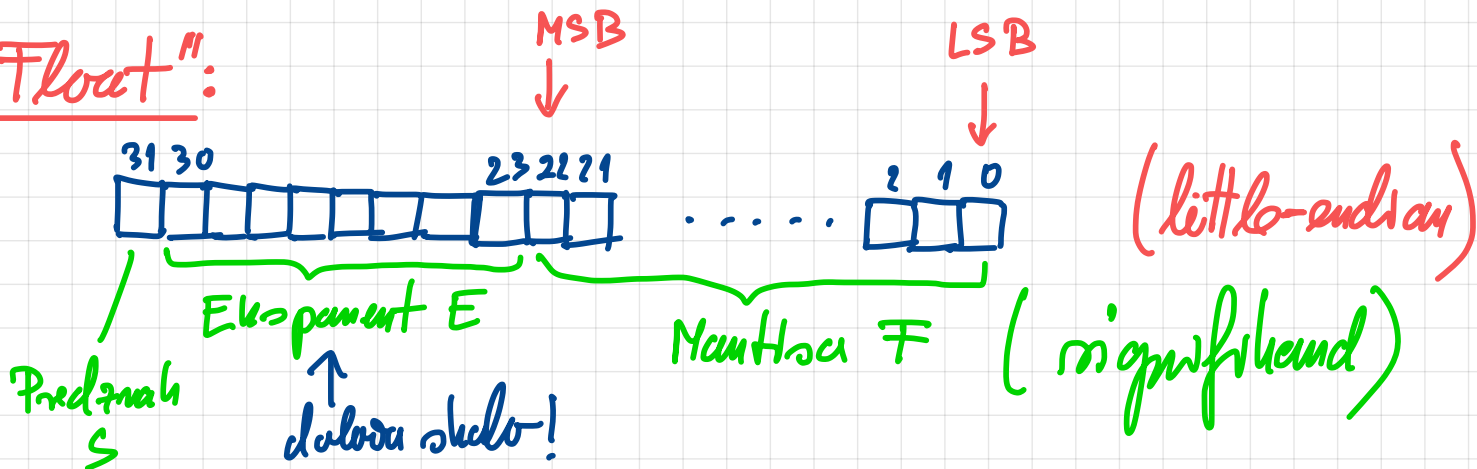
Določimo sorodnosti med
naturnostjo in širjenjem števil.

Za predstavitve realnih števil uporabljamo aritmetiko v plavajoči decimalni vejici. V računalnikih jih predstavimo z mzi dolžini 32 oziroma 64 bitov. Po tem upoštevamo standard IEEE 754.

- Števila v plavajoči vejici v enojni natančnosti (v C/C++ "float")

- Števila v plavajoči vejici v dvojni natančnosti (v C/C++ "double")

"Float":



Števila so predstavljena na naslednji način:

$0 < E < 255$: "normalizirane vrednosti"

$$f_p(x) = (-1)^s \cdot 2^{E-127} \cdot \left[1 + \sum_{i=1}^{23} b_{23-i} 2^{-i} \right]$$
$$= (-1)^s \cdot 2^{E-127} \cdot \left(1 + \frac{b_{22}}{2^1} + \frac{b_{21}}{2^2} + \dots + \frac{b_0}{2^{23}} \right)$$

$E = 0$: "denormalizirane vrednosti"

$$f_p(x) = (-1)^s \cdot 2^{-126} \cdot \left[0 + \sum_{i=1}^{23} b_{23-i} 2^{-i} \right]$$
$$= (-1)^s \cdot 2^{-126} \cdot \left(\frac{b_{22}}{2^1} + \frac{b_{21}}{2^2} + \dots + \frac{b_0}{2^{23}} \right)$$

Zgled: (Martin Horvat)

0 10000000 000000000000000000000000 = +1 * 2**(128-127) * 1.0 = 2
0 10000001 101000000000000000000000 = +1 * 2**(129-127) * 1.101 = 6.5
1 10000001 101000000000000000000000 = -1 * 2**(129-127) * 1.101 = -6.5

0 00000001 000000000000000000000000 = +1 * 2**(1-127) * 1.0 = 2**(-126)
0 00000000 100000000000000000000000 = +1 * 2**(-126) * 0.1 = 2**(-127)
0 00000000 000000000000000000000001 = +1 * 2**(-126) *
0.0000000000000000000000000001 =
2**(-149) (Smallest positive value)



Web calculator

Najmanjšo število predstavljen
v enajmih mestih dvajseti.

The screenshot shows a web browser window with the URL <https://www.h-schmidt.net/FloatConverter/IEEE754.html>. The page title is "IEEE-754 Floating Point Converter". The main content area displays the conversion of the decimal value "0.1" into IEEE 754 floating point format. The input field "You entered" contains "0.1". The "Value actually stored in float" is shown as "0.100000001490116119384765625". The "Error due to conversion" is "1.490116119384765625E-9". The "Binary Representation" is "00111101110011001100110011001101". The "Hexadecimal Representation" is "0x3dccccd". The interface also shows the IEEE 754 Converter (JavaScript), V0.22 layout with fields for Sign, Exponent, and Mantissa, and a grid of checkboxes for the binary digits.

Pravilna številca:

$$E = 255 \quad \text{in} \quad F \neq 0 \Rightarrow fl(x) = \text{NaN} \\ \text{"Not a number"}$$

$$E = 255 \quad \text{in} \quad F = 0 \quad \text{in} \quad S = 1 \Rightarrow fl(x) = -\infty$$

$$E = 255 \quad \text{in} \quad F = 0 \quad \text{in} \quad S = 0 \Rightarrow fl(x) = +\infty$$

Tetova rodinaya s plavajoca vesica je ta, da so bitne operacije obravnavane z napako.

Natančnost aritmetike interpretiramo kot napredje število, za katero velja:

$$fl(1 + \epsilon) = 1$$

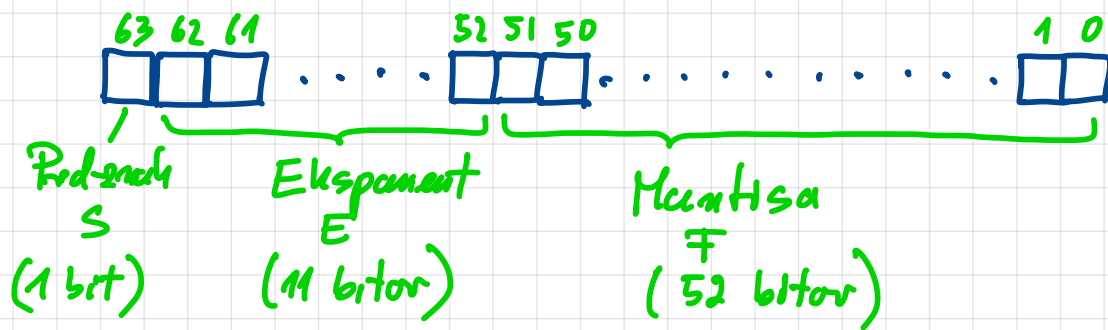
Tu se miha ne pozna, ker je manjše od daljše mantise.

Odnosno je od zaokroževanja. Po IEEE standardu zaokrožujemo k najbližjemu predstavljenemu rezultatu:

$$\epsilon \approx \frac{\epsilon_M}{2} \quad \text{Skrajna natančnost.}$$

$$\epsilon_M = 2^{-23} \approx 1.2 \cdot 10^{-7}$$

"Double":



Števila so predstavljena na naslednji način:

$0 < E < 2047$: "denormalizirane vrednosti"

$$f_p(x) = (-1)^s \cdot 2^{E-1023} \cdot \left[1 + \sum_{i=1}^{52} b_{52-i} 2^{-i} \right]$$
$$= (-1)^s \cdot 2^{E-1023} \cdot \left(1 + \frac{b_{51}}{2^1} + \frac{b_{50}}{2^2} + \dots + \frac{b_0}{2^{52}} \right)$$

$$= (-1)^s 2^{E-1023} \cdot \underbrace{1. b_{51} b_{50} b_{49} \dots b_0}_{\text{binarna baza}}$$

$E = 0$: "denormalizirane vrednosti"

$$f_p(x) = (-1)^s \cdot 2^{-1022} \cdot \left[0 + \sum_{i=1}^{52} b_{52-i} 2^{-i} \right]$$
$$= (-1)^s \cdot 2^{-1022} \cdot \underbrace{0. b_{51} b_{50} b_{49} \dots b_0}_{\text{binarna baza}}$$

Pravilna številca:

$$E = 2047 \text{ in } F \neq 0 \Rightarrow fl(x) = \text{NaN} \\ \text{"Not a number"}$$

$$E = 2047 \text{ in } F = 0 \text{ in } S = 1 \Rightarrow fl(x) = -\infty$$

$$E = 2047 \text{ in } F = 0 \text{ in } S = 0 \Rightarrow fl(x) = +\infty$$

Tetova rodinaya s plavayosa vesica je ta, da so bitume operayosa obravnavane z napako.

Natančno aritmetike interpretiramo kot napredje štelo, za kalibra velja:

$$fl(1 + \epsilon) = 1$$

Ta se miha ne potma, ker je manjše od daljše mantise.

Odnosno je od zaokroževanja. Po IEEE standardu zaokrožujemo k najbližjemu predstavljenu rezultatu:

$$\epsilon \approx \frac{\epsilon_M}{2}$$

Štefma natančnost.

$$\epsilon_M = 2^{-52} \approx 2.22 \cdot 10^{-16}$$

Najveći exponent: $E_{max} = 1023$

Najmanji exponent: $E_{min} = -1022$

Najmanje predstavljiva številka: $4.94 \cdot 10^{-324}$

Najveći normalna številka: $1.80 \cdot 10^{308}$

Za več informacij glej: Širca, RMF, str 2, 595

Takeaway message:

- Z realnim številom v enojni natančnosti, ie 32 bit ("float") lahko sodenamo na 7 mest natančno!
- Z realnim številom v dvojni natančnosti, ie 64 bit ("double") lahko sodenamo na 16 mest natančno
- Paziti, ko odštevamo, primerjamo številke! Upoštevaj ϵ natančnost!

Zgled PlavajocaVejica.py

```
1 #!/usr/bin/env python3
2
3 import numpy as np
4
5
6 x=np.finfo(np.float64);
7 print("Lastnosti stevil tipa: ", x.dtype)
8 print("- Najmanjse stevilo je", x.min)
9 print("- Najvecje stevilo je", x.max,'\n')
10 print("- Najmanjsi eksponent v dvojski bazi je ", x.minexp)
11 print("- Najvecji eksponent v dvojski bazi je ", x.maxexp,'\n')
12 z=1
13 print("- Najblizji vecje stvilo stevila",z," je: ", '%.20f'%np.nextafter(z,z+1))
14 print("- Najblizji manjse stvilo stevila",z," je: ", '%.20f'%np.nextafter(z,z-1))
15
16
17
18 # The smallest number in double: 1/2^53
19 dx=1./pow(2,52);
20 print("Najmanjsa predstavljiva razlika dx je ", dx);
21 print('1+dx = ', '%.20f'%np.double(1+dx));
22
23 dx=1./pow(2,53);
24 print("Ze premajhna razlika dx je ", dx);
25 print('1+dx = ', '%.20f'%np.double(1+dx));
26
27 print('Zadnja decimalka v stevilu 2.0000000000000001 ', '%.20f'%np.double(1+dx));
28
29
30 def shownum(num):
31     y=np.double(num)
32     print('Stevilo', num,'se izpise kot:')
33     print('%s'%y)
34     return 0
35
36 # sprememba je premajhna, da bi jo lahko shranili v double
37 shownum('2.0000000000000001')
38
39 # majhne vrednosti (na robu) pa se shranijo narobe.
40 shownum('2.0000000000000001')
41
42 # Primer nauka, kadar imamo dve skoraj enaki veliki stevili moramo paziti, kako ju
43 # odstevamo in gledamo razliko.
44
```

```
In [187]: runfile('/Users/miham/Desktop/Sola/MFSEM-2022/Uvod/PlavajocaVejica/
PlavajocaVejica.py', wdir='/Users/miham/Desktop/Sola/MFSEM-2022/Uvod/PlavajocaVejica')
Lastnosti stevil tipa: float64
- Najmanjse stevilo je -1.7976931348623157e+308
- Najvecje stevilo je 1.7976931348623157e+308

- Najmanjsi eksponent v dvojski bazi je -1022
- Najvecji eksponent v dvojski bazi je 1024

- Najblizji vecje stvilo stevila 1 je: 1.00000000000000022204
- Najblizji manjse stvilo stevila 1 je: 0.99999999999999988898

Najmanjsa predstavljiva razlika dx je 2.220446049250313e-16
1+dx = 1.00000000000000022204

Ze premajhna razlika dx je 1.1102230246251565e-16
1+dx = 1.00000000000000000000

Zadnja decimalka v stevilu 2.0000000000000001 1.00000000000000000000

Stevilo 2.0000000000000001 se izpise kot:
2.00000000000000000000

Stevilo 2.0000000000000001 se izpise kot:
2.000000000000000088818
```