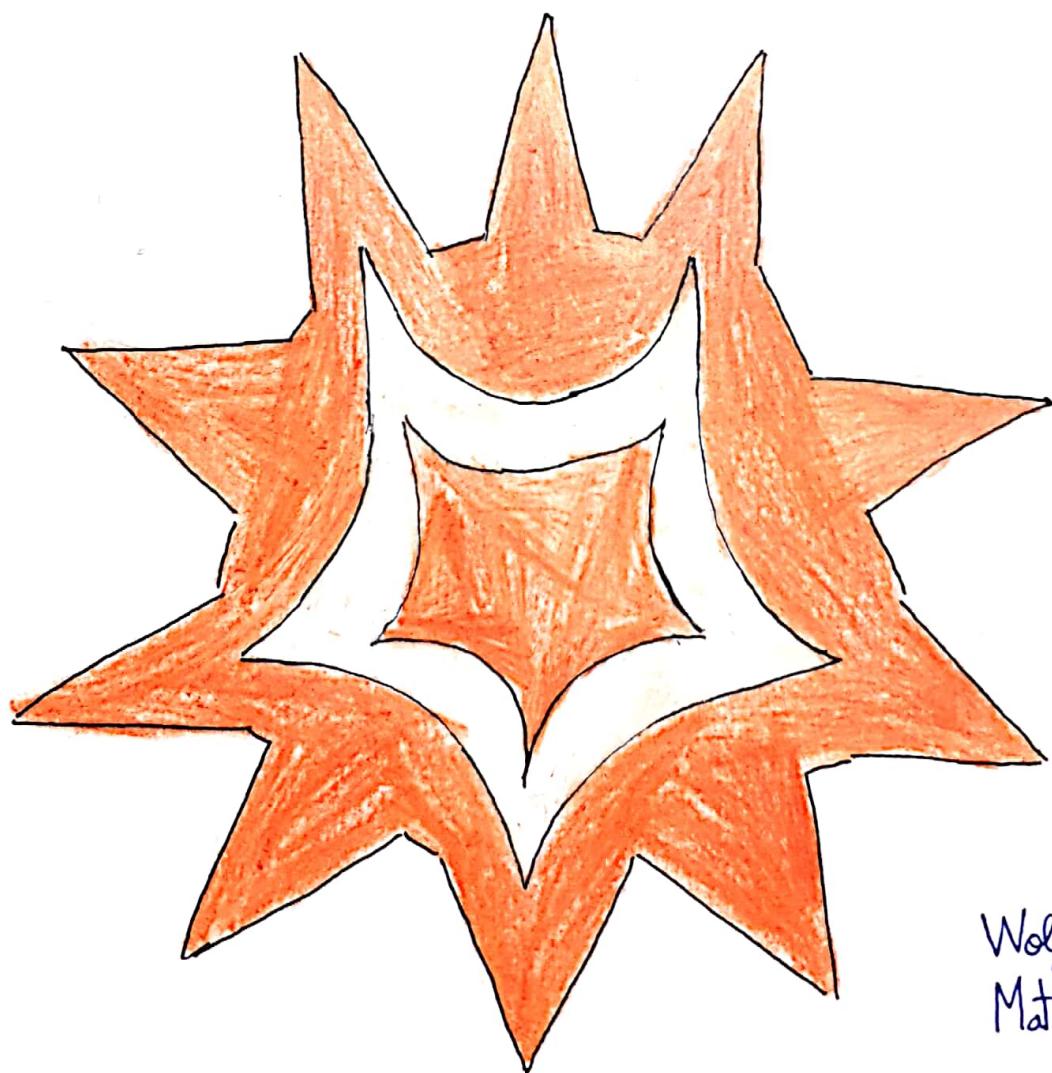


MÉTODOS



Wolfram
Mathematica

José P.Z.



Chapter 1

Radix conversion

We use decimal notation:

$$231'45_{10} = 2 \cdot 10^2 + 3 \cdot 10^1 + 1 \cdot 10^0 + 4 \cdot 10^{-1} + 5 \cdot 10^{-2}$$

→ We say that $231'45$ is in radix ten.

Suppose we want to convert it to radix eight.

$$\begin{aligned} 231'45_8 &= 2 \cdot 8^2 + 3 \cdot 8^1 + 1 \cdot 8^0 + 4 \cdot 8^{-1} + 5 \cdot 8^{-2} = \\ &= 153'578125_{10} \end{aligned}$$

Most common ones: 10, 2, 8, 16 → hexadecimal

Radix 16:

0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F

Radix 10:

0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15

base of Radix 10. (equivalents for A, ..., F)

Radix 2:

0 0 0 0 → 0

0 0 0 1 → 1

0 0 1 0 → 2

0 0 1 1 → 3

0 1 0 0 → 4

0 1 0 1 → 5

0 1 1 0 → 6

0 1 1 1 → 7

1 0 0 0 → 8

1 0 0 1 → 9

1 0 1 0 → 10 — A

1 0 1 1 → 11 — B

1 1 0 0 → 12 — C

1 1 0 1 → 13 — D

1 1 1 0 → 14 — E

1 1 1 1 → 15 — F

Error on memory 0xABDF8DE7! A classic error

that (advises of incoming hex number) references memory location
as hexadecimal number.

To binary equivalent:

$\widehat{ABDF8DE7} \rightarrow$

$= \widehat{0101010111011111100011011100111}$

To octal:

$= 25367706747$

Radix 8: (Octal)

$000 \rightarrow 0$
 $001 \rightarrow 1$
 $010 \rightarrow 2$
 $011 \rightarrow 3$
 $100 \rightarrow 4$
 $101 \rightarrow 5$
 $110 \rightarrow 6$
 $111 \rightarrow 7$

Decimal rounding:

$$2' \widehat{9}_{10} = 3_{10}$$

$$13' \widehat{247}_8 = 13' 25_8$$

$$101' \widehat{1001}_2 = 101' 1010_2$$

Briefing:

- From any radix to decimal: just substitute its value ex:

$$AB7_{16} = 10 \cdot 16^2 + 11 \cdot 16^1 + 7 \cdot 16^0 = 2743_{10}$$

$$745_8 = 7 \cdot 8^2 + 4 \cdot 8^1 + 5 \cdot 8^0 = 485_{10}$$

$$1001 = 2^3 + 2^0 = 9_{10}$$

- From hex to binary or vice versa, just substitute each hexadecimal character for its corresponding 4 bits (starting from left). (Solve)

$$AB7_{16} = \dots 0111_2 = \underbrace{1010}_{A_{16}} \underbrace{1011}_{B_{16}} \underbrace{0111}_{7_{16}} 1_2$$

$$\underbrace{11001}_{1} \underbrace{001}_{9} \underbrace{01}_{2} \underbrace{0111}_{B} \underbrace{1000}_{8_{16}} = 192B8_{16}$$

(if we have at the left and incomplete group of numbers, odd to the left until we have a group of four bits) This works because of powers of 2 (both bases are).

$$A00_{16} = 10 \cdot 16^2 = 10 \cdot (2^4)^2 = 10 \cdot 2^8$$

$$\underbrace{1010}_{1} \underbrace{0000}_{2} \underbrace{0000}_{5} = 2^1 + 2^0 = (\underbrace{2^2 + 1}_{5}) \cdot 2^0 = 10 \cdot 2^8$$

- The same goes for the other relations between radix 16, 8 and 2.

Examples:

$$317_8 = 011001111_2 = 207_{10}$$

$$\underbrace{011}_{1} \underbrace{01}_{5} = 15_8 = 8 + 5 = 13_{10}$$

$$317_8 = 011 \underbrace{00}_{C} \underbrace{1111}_{F} = CF_{16} = 207_{10}$$

(from radix 8 to 16, an intermediate step. Works similar
the other way around) (Left to right to make the groups!)

$$CF_{16} = \underbrace{11}_{3} \underbrace{00}_{1} \underbrace{1111}_{7} = \underbrace{011}_{3} \underbrace{001}_{1} \underbrace{111}_{7} = 317_8 = 207_{10}$$

- From radix ten to binary. Mark all the powers of two from 2^1 to 2^x so that $2^x \geq n$. And subtract the highest possible power of 2. Or, divide by two:

$$109 \rightarrow \overline{128} \quad \overline{64} \quad \overline{32} \quad \overline{16} \quad \overline{8} \quad \overline{4} \quad \overline{2} \quad \overline{1}$$

2⁷ 2⁶ 2⁵ 2⁴ 2³ 2² 2¹ 2⁰

$$109 = 64 + 45 = 64 + 32 + 13 = 64 + 32 + 8 + 4 + 1 \rightarrow$$

$$\Rightarrow \underline{0} \underline{1} \underline{1} \underline{0} \underline{1} \underline{1} \underline{0} \underline{1} \rightarrow 109_{10} = \underline{\underline{1101101}}_2 =$$

$$= \underbrace{0110}_6 \underbrace{1101}_{D_{16}} = \underbrace{001}_1 \underbrace{101}_5 \underbrace{101}_{58} \Rightarrow$$

$$\Rightarrow 109_{10} = 1101101_2 = 6D_{16} = 155_8$$

Or:

$$\begin{array}{r} 109 \\ 09 \end{array} \quad \begin{array}{r} 54 \\ 14 \end{array} \quad \begin{array}{r} 27 \\ 07 \end{array} \quad \begin{array}{r} 13 \\ 1 \end{array} \quad \begin{array}{r} 6 \\ 0 \end{array} \quad \begin{array}{r} 3 \\ 1 \end{array} \quad \begin{array}{r} 1 \\ 0 \end{array}$$

$\xleftarrow{\text{write in this numbers order}}$

$$\Rightarrow 109 = 1011_2$$

$$\left\{ \begin{array}{ll} 109 \% 2 & r = 1 \\ 54 \% 2 & r = 0 \\ 27 \% 2 & r = 1 \\ 13 \% 2 & r = 1 \\ 6 \% 2 & r = 0 \\ 3 \% 2 & r = 1 \\ 1 \% 2 & r = 1 \\ 0 \% 2 & r = 0 \end{array} \right. \rightarrow 109_{10} = 1101101_2$$

Why does this work? Think modularly:

$$109_{10} = 9 + 0 \cdot 10 + 1 \cdot 100 \Rightarrow$$

(know we don't mind these).

$$109_{10} \bmod(10) = 9 \rightarrow \text{units}$$

(know we don't mind these).

$$109_{10} \bmod(100) = 09 = 0 \cdot 10 + 9 \text{ units} \Rightarrow 0 \text{ tens.}$$

$$109_{10} - (09) = 100 \bmod(1000) = 1 \text{ hundreds.}$$

Ex:

$$238 \equiv 8 \bmod(10) \rightarrow \begin{cases} 238 \equiv 38 \bmod(100) \\ 230 \equiv 30 \bmod(100) \end{cases}$$

Better seen this way:

$$109 = a_0 \cdot 2^0 + a_1 \cdot 2^1 + \dots + a_7 \cdot 2^7 = \\ = a_0 + 2 \cdot (a_1 + 2 \cdot (a_2 + \dots + 2 \cdot (a_7))) \dots$$

$$\Rightarrow 109 \equiv a_0 \pmod{2} \Rightarrow a_0 = \text{resto } 109 / 2 \rightarrow 1$$

$$\left(\frac{109 - a_0}{2} \right) \equiv a_1 \pmod{2} \Rightarrow a_1 = \text{resto } \left(\frac{109 - a_0}{2} \right) / 2$$

$$\text{Now, } \overbrace{109}^{\text{dividend}} = \underbrace{2}_{\text{divisor}} \cdot \underbrace{\left(\frac{109 - a_0}{2} \right)}_{\text{coefficient}} + \underbrace{a_0}_{\text{remainder}}$$

So if:

$$109 = a_0 \cdot 2^0 + \dots + a_7 \cdot 2^7 \text{ and we define}$$

$$c_0 \text{ as } \frac{109 - a_0}{2}, \quad c_n = \frac{c_{n-1} - a_n}{2} \quad (c_{n-1} \geq c_n)$$

And thus our algorithm:

$$x = \underline{a_0 \cdot 2^0 + \dots + a_k \cdot 2^k} \Rightarrow$$

$$x \equiv a_0 \pmod{2}$$

$$\frac{x - a_0}{2} = c_0 \equiv a_1 \pmod{2}$$

$$\frac{\left(\frac{x - a_0}{2} \right) - a_1}{2} \equiv \frac{c_0 - a_1}{2} \equiv a_2 \pmod{2}$$

... We start computing $x = 2 \cdot c_0 + a_1$ then:

$$\Rightarrow c_0 = 2 \cdot c_1 + a_2 = \left(\frac{x - a_1}{2} \right) \dots$$

$$\left\{ \begin{array}{l} x = 2 \cdot c_0 + a_0 \\ c_0 = 2 \cdot c_1 + a_1 \\ \vdots \\ c_k = 2 \cdot c_{k+1} + a_{k+1} \end{array} \right.$$

Example:

$$\begin{aligned} 109 &= 2 \cdot 54 + 1 \rightarrow a_0 = 1 \\ 54 &= 2 \cdot 27 + 0 \rightarrow a_1 = 0 \\ 27 &= 2 \cdot 13 + 1 \rightarrow a_2 = 1 \\ 13 &= 2 \cdot 6 + 1 \rightarrow a_3 = 1 \quad \rightarrow \\ 6 &= 2 \cdot 3 + 0 \rightarrow a_4 = 0 \\ 3 &= 2 \cdot 1 + 1 \rightarrow a_5 = 1 \\ 1 &= 2 \cdot 0 + 1 \rightarrow a_6 = 1 \end{aligned}$$

$$\begin{aligned} \rightarrow \boxed{109}_{10} &= a_0 \cdot 2^0 + \dots + a_6 \cdot 2^6 + a_7 \cdot 2^7 = \\ &= 1 \cdot 2^0 + \dots + 1 \cdot 2^6 = \\ &= a_6 a_5 a_4 \dots a_0 (2) = \boxed{1101101_2} \end{aligned}$$

(These are my notes, now notes from class).

Our aim is to convert a b-radix number into a B-radix number.

$$u_m \dots u_2 u_1 u_0 . u_{-1} u_{-2} u_{-3} \dots u_{-n}$$

Example:

$$253_8 = 2 \cdot 8^2 + 5 \cdot 8^1 + 3 \cdot 8^0 = 2 \cdot 64 + 40 + 3 = \\ = (2 \cdot 8 + 5) \cdot 8 + 3 = 3 + 8 \cdot (5 + 8 \cdot (2)) = 171_{10}$$

What is its radix 7 equivalent?

$$(2_7 \cdot 11_7 + 5_7) \cdot 11_7 + 3_7 = (22_7 + 5_7) \cdot 11_7 + 3_7 = \\ = 30_7 \cdot 11_7 + 3_7 = 330_7 + 3_7 = 337_7$$

We check:

$$337_7 = (3 \cdot 7 + 3) \cdot 7 + 3 = 24 \cdot 7 + 3 = 168 + 3 = 171$$

Book: Donald Knuth, the art of computer programming
(.djvu format (eastern pdf *10 better compression)).

To transform from radix-10 to other radix:

$$\begin{array}{r} 171 \\ \underline{-11} \quad 21 \\ \hline 35 \end{array} \quad \begin{array}{r} 18 \\ \underline{-18} \\ \hline 0 \end{array} \quad \begin{array}{r} 18 \\ \underline{-18} \\ \hline 0 \end{array}$$

$$171_{10} = 253_8$$

$$\begin{array}{r}
 333_7 \\
 \times 3 \\
 \hline
 03 \\
 3 \\
 \hline
 10 \\
 \hline
 \end{array}
 \quad
 \begin{array}{r}
 11_7 \\
 \hline
 30_7 \\
 -05 \\
 \hline
 2 \\
 \hline
 0 \\
 \hline
 \end{array}
 \quad
 \begin{array}{r}
 11_7 \\
 \hline
 2 \\
 \hline
 0 \\
 \hline
 \end{array}
 \quad
 \Rightarrow 333_7 = 253_8$$

(Able to change from any radix to any radix with this method!)

$$\left(X_b = \sum_{i=0}^n r_i b^i \right) \rightarrow X_b = r_n r_{n-1} \dots r_3 r_2 r_1 B$$

For the decimal part; from any radix to radix ten, similar:

Example: $1.32_4 = 1 + 3 \cdot 4^{-1} + 2 \cdot 4^{-2} = 1 + 3 \cdot \frac{1}{4} + 2 \cdot \frac{1}{16} = 1.82$

$$0.132 = 1 \cdot 4^{-1} + 3 \cdot 4^{-2} + 2 \cdot 4^{-3} =$$

$$\begin{aligned}
 &= \left(\left(2 \cdot \frac{1}{4} + 3 \right) \cdot \frac{1}{4} + 1 \right) \cdot \frac{1}{4} \\
 &= \left(\left(\frac{1}{2} + \frac{6}{4} \right) \cdot \frac{1}{4} + 1 \right) \cdot \frac{1}{4} = \left(\frac{7}{8} + \frac{8}{8} \right) \cdot \frac{1}{4} = \frac{15}{32} = \\
 &= 0.46875_{10}
 \end{aligned}$$

And the other way around: (radix 10 or other, to any radix).

$$\begin{aligned}
 0.46875_{10} &\rightarrow 0.46875 \cdot 4 = 1.875 = 1 + 0.875 \rightarrow \lfloor 1.875 \rfloor = 1 \\
 0.875 \cdot 4 &= 3 + 0.5 \quad (\text{Note, what if } 23 + 0.5?)
 \end{aligned}$$

$$\begin{aligned}
 0.5 \cdot 4 &= 2 + 0 \\
 (0 \cdot 4 = 0 + 0)
 \end{aligned}$$

$$\Rightarrow 0.46875_{10} = 0.132_4$$

(Note: $x \in [0, 1] \rightarrow 0 \leq x \cdot N \leq N \rightarrow$ we can always take the outcome.)

Remember:

$$\text{Floor function} \rightarrow \lfloor 3.3 \rfloor = 3$$

$$\text{Ceiling function} \rightarrow \lceil 3.3 \rceil = 4$$

Example:

$$12112^1\overbrace{02}_3 = x$$

For periodic decimals.

$$12112_3 = 1 \cdot 81 + 2 \cdot 27 + 1 \cdot 9 + 1 \cdot 3 + 2 \cdot 3^0 = 149$$

$$0^1\overbrace{02}_3 = \frac{2}{3^2} + \frac{2}{3^4} + \frac{2}{3^6} + \dots = 2 \cdot \sum_{n=1}^{\infty} \left(\frac{1}{9}\right)^n = 2 \cdot \frac{\frac{1}{3^2}}{1 - \frac{1}{3^2}} = \frac{\frac{2}{3^2}}{1 - \frac{1}{3^2}} = \frac{1}{4}_{10}$$

Example:

$$x_4 = 12112^1\overbrace{02}_3$$

$$x_3 = 12112^1\overbrace{02}_3$$

$$100_3 \cdot x_3 = 1211202^1\overbrace{02}_3$$

$$-(100_3 - 1_3) x_3 = 12112^1\overbrace{02}_3 - 1211202^1\overbrace{02}_3 \Rightarrow$$

$$\Rightarrow (21_4 - 1_4) x_4 = \frac{1211202^1\overbrace{02}_3}{1122020_3} \Rightarrow$$

$$\Rightarrow 20_4 x_4 = 1122020_3$$

1122020₃

02
22

00
02
20
02

$$\underline{4_{10} = 11_3}$$

102001

$$\underline{\begin{array}{r} 11_3 \\ -2 \\ \hline \end{array}}$$

... (finish as exercise).

RADIX CONVERSION BRIEFING

- Integer from base N to base 10:

$$a_m a_{m-1} \dots a_0 N = a_0 \cdot N^0 + a_1 \cdot N^1 + \dots + a_m \cdot N^m = X_{10}$$

- From base 10 to N or generally, from base N₁ to N₂:

$$a_m a_{m-1} \dots a_0 N_1 = b_m' b_{m-1}' \dots b_0 N_1 \cdot (N_2)_{N_1} + r_1$$

N₂ written
in base N₁

$$b_m' b_{m-1}' \dots b_0 N_1 = c_m'' c_{m-1}'' \dots c_0 \cdot N_1 \cdot (N_2)_{N_1} + r_2$$

⋮

$$c_0 \cdot N_1 \cdot (N_2)_{N_1} + r_2 = X_{10}$$

$$y_m''' y_{m-1}''' \dots y_0 N_1 = \underbrace{z_m''' z_{m-1}''' \dots z_0 N_1}_{\text{ON}_1} \cdot (N_2)_{N_1} + r_n \leq (N_2)_{N_1}$$

$$0 = //$$

$$\rightarrow a_m a_{m-1} \dots a_0 N_1 = r_n r_{n-1} \dots r_2 r_1 N_2$$

(Note: r_i can be $0 \leq r_i \leq N_1 \leq r_i \leq N_2$, for example, when passing from radix 4 to 16 the remainder r_i can be $0 \leq r_i < 15 < 16$, but since $0 \leq r_i < N_2 \Rightarrow$ it's only one character in radix N₂)

(in this case, 15 = G)

Floating in base N to base 10:

Analogous to integer base N to base 10 but this time;

$$0.a_1a_2\dots a_m N = a_1 \cdot N^{-1} + a_2 \cdot N^{-2} + \dots + a_m \cdot N^{-m} = X_{10}$$

Floating in base N_1 to base N_2 :

$$0.a_1a_2\dots a_m N_1 \rightarrow$$

$$0.a_1a_2\dots a_m N_1 \cdot (N_2)_{N_1} = r_1 \cdot b_1 b_2 \dots b_m, \quad 0 \leq r_i < N_2$$

$$0.b_1 b_2 \dots b_m N_1 \cdot (N_2)_{N_1} = r_2 \cdot c_1 c_2 \dots c_m, \quad 0 \leq r_i < N_2$$

(* check last note on r_i)

$$0.y_1 y_2 \dots y_m N_1 \cdot (N_2)_{N_1} = r_n \cdot 0, \quad 0 \leq r_n < N_2 \Rightarrow$$

$$\Rightarrow 0.a_1a_2\dots a_m N_1 = r_1 r_2 \dots r_n N_2 //$$

Periodic base N to base 10:

$$0' b_1 b_2 \dots b_m \overbrace{a_1 a_2 \dots a_n}^n N \rightarrow \underbrace{0' b_1 b_2 \dots b_m}_\text{previous case} N + N^{-m+1} \cdot \underbrace{0' a_1 a_2 \dots a_n}_\text{what we really care} N$$

$$0' \underbrace{a_0 a_1 \dots a_{n-1}}_\text{n elements} N = \sum_{i=0}^{\infty} \frac{a_0}{N^{1+n \cdot i}} + \dots + \sum_{i=0}^{\infty} \frac{a_{n-1}}{N^{n-1+n \cdot i}}$$

$$= \frac{a_0}{N} \cdot \underbrace{\sum_{i=0}^{\infty} \left(\frac{1}{N^n}\right)^i}_{\text{formula}} + \dots + \frac{a_{n-1}}{N^{n-1}} \cdot \underbrace{\sum_{i=0}^{\infty} \left(\frac{1}{N^n}\right)^i}_{\text{formula}}$$

$$= \frac{\left(\frac{1}{N^n}\right)^{k_0=0}}{1 - \frac{1}{N^n}} \cdot \left[\frac{a_0}{N^n} + \dots + \frac{a_{n-1}}{N^{n-1}} \right] = X_{10}$$

• Periodic N_1 to N_2 (full number):

$$a_1 \dots a_m \cdot \overbrace{a'_1 \dots a'_m b_1 \dots b_n}^{N_1} = X_{N_2} \rightarrow ?$$

$$\Rightarrow (10^{n+m})_{N_2} \cdot a_1 \dots a_m \cdot \overbrace{a'_1 \dots a'_m b_1 \dots b_n}^{N_1} = a_1 \dots a_m a'_1 \dots a'_m b_1 \dots b_n \cdot \overbrace{b_1 \dots b_n}^{N_1} \Rightarrow$$

$$(10^m)_{N_2} \cdot a_1 \dots a_m \cdot \overbrace{a'_1 \dots a'_m b_1 \dots b_n}^{N_1} = a_1 \dots a_m a'_1 \dots a'_m \cdot \overbrace{b_1 \dots b_n}^{N_1} \Rightarrow$$

$$(10^{n+m} - 10^m)_{N_2} \cdot \underbrace{a_1 \dots a_m \cdot a'_1 \dots a'_m b_1 \dots b_n}_{X_{N_2}} = [a_1 \dots a_m a'_1 \dots a'_m b_1 \dots b_n - a_1 \dots a_m a'_1 \dots a'_m]_{N_1}$$

$$\Rightarrow a_1 \dots a_m \cdot \overbrace{a'_1 \dots a'_m b_1 \dots b_n}^{N_1} = \frac{[a_1 \dots a_m a'_1 \dots a'_m b_1 \dots b_n - a_1 \dots a_m a'_1 \dots a'_m]}{(10^{n+m} - 10^m)_{N_1}}$$

(Transform this very last fraction) \Rightarrow

$$X_{N_2} = \frac{([a_1 \dots a_m a'_1 \dots a'_m b_1 \dots b_n - a_1 \dots a_m a'_1 \dots a'_m])_{N_2}}{([10^{n+m} - 10^m]_{N_2})}$$

- X_{N_2} is not necessarily the irreducible fraction, to achieve this on intermediate step from $N_1 \rightarrow$ base 10 \rightarrow reduce $\rightarrow N_2$ might be necessary.
- They might ask for the decimal expansion, in which case we have to do the division in base N_2 of the respective fraction.

Note: Don't mix any of these methods.

Each one of them has its own demonstration, so using the integer methods to the floating decimals it's a big mistake!

Example:

$$x_3 = 12112.0\overline{2}_3 \xrightarrow{\text{to}} \text{radix 4?}$$

$$100_3 \cdot x_3 = 1211202.0\overline{2}_3$$

$$(100_3 - 1_3) x_3 = (21_4 - 1_4) x_4 = 20_4 x_4 = - \frac{12112.0\overline{2}_3}{1122020_3}$$

$$\begin{array}{r} 4 \\ \hline 1122020_3 \\ 02 \\ 22 \\ 000 \\ 02 \\ 20 \\ \hline 02 \end{array} \quad \begin{array}{r} 11_3 \\ \hline 102001 \\ 100 \\ 10 \\ 101 \\ 2 \\ \hline \end{array}$$

$$\begin{array}{r} 11_3 \\ \hline 2202 \\ 002 \\ \hline 002 \end{array}$$

$$1211202.0\overline{2}_3$$

$$\begin{array}{r} 12112.0\overline{2}_3 \\ \hline 1122020_3 \end{array}$$

$$\begin{array}{r} 11_3 \\ \hline 200 \\ 20 \\ 2 \\ \hline 11 \\ 11 \\ 1 \\ \hline 0 \end{array}$$

$$\Rightarrow 1122020_3 = 102222_4$$

Another method:

$$(((\underbrace{(1 \cdot 3 + 1) \cdot 3 + 2}_{10_4}) \cdot 3 + 2) \cdot 3 + 2)$$

$$\begin{array}{r} 32_3 \\ \hline 12 \\ 21 \\ \hline 222 \\ \hline \end{array} \quad \begin{array}{r} 230_3 \\ \hline 0 \\ 21 \\ 12 \\ \hline 2010 \\ \hline \end{array} \quad \begin{array}{r} 2012_3 \\ \hline 12 \\ 03 \\ 2 \\ \hline 12102 \\ \hline \end{array}$$

(Incorrect result)

Solution:

$$(((\underbrace{(1 \cdot 3 + 1) \cdot 3 + 2}_{10_4}) \cdot 3 + 2) \cdot 3 + 0) \cdot 3 + 0$$

$$\begin{array}{r} 10_4 \\ \hline 32_4 \\ \hline 230_4 \end{array}$$

$$\begin{array}{r} 2010_3 \\ \hline 3 \\ \hline 0 \end{array}$$

$$\begin{array}{r} 12030_3 \\ \hline 12 \\ 0 \\ \hline 12030 \end{array}$$

$$2010_4$$

$$\begin{array}{r} 12032_3 \\ \hline 3 \\ 12 \\ 02 \\ 12 \\ 02 \\ \hline 0222 \end{array}$$

$$\begin{array}{r} 12032_3 \\ \hline 3 \\ 12 \\ 00 \\ 12 \\ 02 \\ \hline 102222 \end{array}$$

$$12030_4 + 2_4$$

$$102222_4$$

$$\begin{array}{r}
 102222_4 \\
 \underline{-02} \\
 02 \\
 \underline{-02} \\
 02 \\
 \underline{-02} \\
 0
 \end{array}
 \quad
 \begin{array}{r}
 | 24 \\
 \hline
 21111 \\
 | 104 \\
 \hline
 104
 \end{array}
 \quad = 2111.1$$

$$\Rightarrow 12112.0\overline{2}_3 = 2111.1_4 //$$

Example:

$$2111.1_4 \longrightarrow \text{Radix } 16?$$

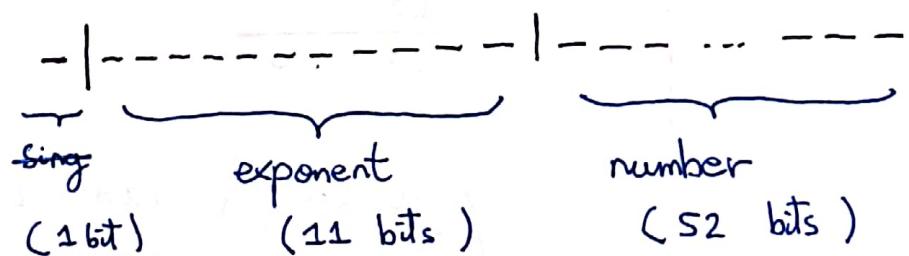
Remember: $16 = 4^2 \rightarrow \text{groups of 2}$

$$\begin{array}{r}
 21 \quad 11 \cdot 10_4 \\
 \underline{9_{10}} \quad \underline{5_{10}} \quad \underline{4_{10}}
 \end{array}
 = 0 \times 95.4 = 95.4_{16}$$

Computer Arithmetic

How do computers store numbers? The IEEE754 standard.

Storing numbers in 64 bits.



$$\begin{array}{ccc}
 \text{sign} & & || \\
 \parallel & & \parallel
 \end{array}$$

$$(-1)^s \cdot 2^{e-1023} \cdot (1 + f)$$

$(s=0 \text{ or } 1)$ $\left\{ \begin{array}{l} c; \text{exponent} \\ -1023; \text{offset} \end{array} \right.$ \downarrow mantissa

16

Example:

27. 56640625

- Since it's positive $\Rightarrow (-1)^s = 1$
- There exist: $\exists n \mid 2^n \leq x < 2^{n+1} \Rightarrow \frac{x}{2^n} = 1$ we choose that power of 2: $\underline{2^n}$ ($1023+n$, because of the offset).
- We add $\{27.56640625 / 2^n\} = 8$ (the mantissa).

So:

$$\begin{cases} s = 0 \\ c = 1027 \\ f = 0.722900390625 \end{cases}$$

Note: $0_{10} = 0_2$
 $1_{10} = 1_2 \Rightarrow$
 s is already binary, no operations needed

$$\Rightarrow 27.56640625 = (-1)^0 \cdot 2^{1027-1023} \cdot (1 + 0.722900390625)$$

Now we pass these numbers to binary. (s is already binary)

- $1027 = 1024 + 2 + 1 = 2^0 + 2^1 + 2^{10} = 10000000011_2$
- $0.722900390625 = \lfloor 2 \cdot 0.722900390625 \rfloor = 1 \quad (1.44580078125)$

$$1.44580078125 \cdot 2 = 0.8916015625 \rightarrow u_{-2} = 0$$

$$u_{-3} = \lfloor 2 \cdot 0.8916015625 \rfloor = \lfloor 1.783203125 \rfloor = 1$$

$$u_{-4} = \lfloor 2 \cdot 0.783203125 \rfloor = \lfloor 1.56640625 \rfloor = 1$$

$$u_{-5} = \lfloor 2 \cdot 0.56640625 \rfloor = \lfloor 1.1328125 \rfloor = 1$$

$$u_{-6} = \lfloor 2 \cdot 0.1328125 \rfloor = \lfloor 0.265625 \rfloor = 0$$

$$u_{-7} = \lfloor 2 \cdot 0.265625 \rfloor = \lfloor 0.53125 \rfloor = 0$$

$$u_{-8} = \lfloor 2 \cdot 0.53125 \rfloor = \lfloor 1.0625 \rfloor = 1$$

$$u_{-9} = \lfloor 2 \cdot 0.0625 \rfloor = \lfloor 0.125 \rfloor = 0$$

$$u_{-10} = \lfloor 2 \cdot 0.125 \rfloor = \lfloor 0.25 \rfloor = 0$$

$$u_{-11} = \lfloor 2 \cdot 0.25 \rfloor = \lfloor 0.5 \rfloor = 0$$

$$U_{-12} = \lfloor 2 \cdot 0.5 \rfloor = \lfloor 1 \rfloor = 1$$

$$U_{-13} = \lfloor 2 \cdot 0 \rfloor // \text{ we finish here.}$$

$$0.722900390625_{10} = 10111001000100000_2 = 8$$

So the binary representation of our number is:

$\begin{array}{r|ccccc|ccccccccc} 0 & | & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & | & 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & \dots & 0 \end{array} \Rightarrow$
the remaining
bits up to 52 spaces
are 0.

$$27.56640625_{10} =$$

$$= 01000000001110111001000100000\dots 0 \text{ in IEEE754}$$

Why the offset? $0, \dots, 2^{15} - 1 = 2047$ possible exponents. To represent exponents from $-1023, \dots, 1024$, we put the offset $(\text{off}) - 1023$, and so the stored c remains positive while the exponent can be both positive and negative.

$$\frac{2^{\text{off}}}{2} - 1 = 1023 \quad (\text{possible positive and negative exponents, 2047 total}).$$

Instead of 64, we could have other possible representations.

16 bits	\rightarrow	5 bits	for the exponent	Most important and used ones.
32 bits	\rightarrow	8 "	" "	
64 bits	\rightarrow	11 "	" "	
128 bits	\rightarrow	15 "	" "	
256 bits	\rightarrow	19 "	" "	

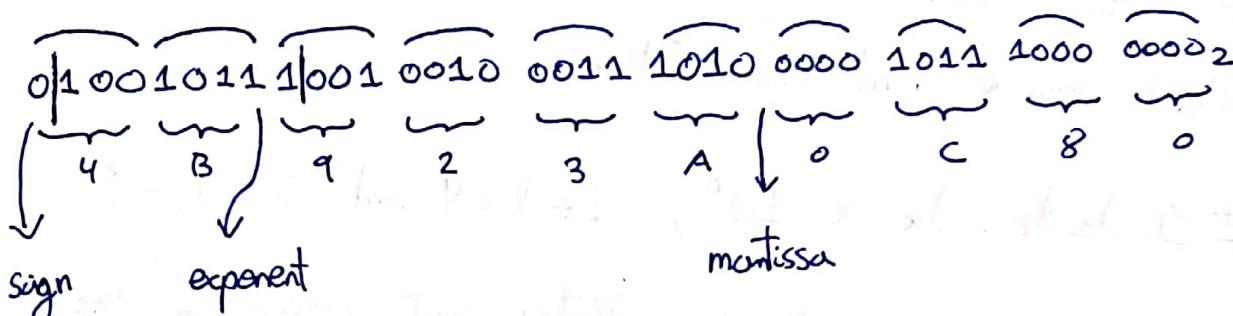
Example:

40 bits machine (1 bit sign, 8 exponent, 31 mantissa)

Convert to its floating representation: $0x4B923A0C80$

To binary (four bits by four)

0 (number is positive) (~~use 2^8 - 1 = 127~~) ($\text{offset} = \frac{2^8}{2} - 1 = 127$)



$$s = 0$$

$$c = 10010111_2 = 1 \cdot 2^0 + 1 \cdot 2^1 + 1 \cdot 2^2 + 1 \cdot 2^4 + 1 \cdot 2^7 = 151$$

$$f = 1.0010010011101000001011100000000_2 =$$

$$\begin{aligned}
 &= 2^7 + 2^8 + 2^9 + 2^{11} + 2^{17} + 2^{19} + 2^{20} + 2^{21} + 2^{25} + 2^{28} \\
 &= 1 + 2^{-3} + 2^{-6} + 2^{-9} + 2^{-10} + 2^{-11} + 2^{-12} + 2^{-19} + 2^{-20} + 2^{-23} \\
 &= 1 + 2^{-3} + 2^{-6} + 2^{-9} + 2^{-10} + 2^{-11} + 2^{-12} \quad \text{check blackboard}
 \end{aligned}$$

$$= 2^{-3} + 2^{-6} + 2^{-10} + 2^{-11} + 2^{-12} + 2^{-14} + 2^{-20} + 2^{-21} + 2^{-24}$$

$$x = (-1)^0 \cdot 2^{151-127} \cdot f \xrightarrow{\text{distribute}} \frac{f}{2^{24}}$$

$$\begin{aligned}
 x &= (2^{24} + 2^{21} + 2^{18} + 2^{14} + 2^{13} + 2^{12} + 2^{10} + 2^9 + 2^8 + 2^7 + 2^6 + 2^5 + 2^4 + 2^3 + 2^2)_2 = \\
 &= 19166233_{10}
 \end{aligned}$$

$$\Rightarrow 0x4B923A0C80 = 19166233_{10}$$

(40 bits representation)

How could we represent ∞ ? And NaN? What is the largest possible number? The smallest? Periodic numbers? Round numbers?
64 bit repr.

Wikipedia of IEEE754

Decimal Machine Numbers

Normalized decimal point representation:

$$\pm 0.d_1d_2\dots d_k \times 10^n, \quad 1 \leq d_1 \leq 9 \text{ and } 0 \leq d_i \leq 9$$

Suppose $y \in \mathbb{R}$. We denote the floating point representation, the number the machine really stores as; $fl(y)$.

Suppose the mantissa can only store k digits and y has more than k digits. There are two strategies:

- Chopping / Truncating:
Just discard the extra digits.

$$y = 0.d_1d_2\dots d_k d_{k+1} d_{k+2} \dots \Rightarrow fl(y) = 0.d_1d_2\dots d_k$$

- Rounding:
Round the number. (Important for accuracy and errors).

$$y = 0.d_1d_2\dots d_k d_{k+1} d_{k+2} \dots 10^n$$

$$fl(y) = 0.d_1d_2\dots d_k + \delta \cdot 5 \cdot 10^{n-(k+1)} \quad \text{with}$$

$$\delta = \begin{cases} 1 & \text{if } d_{k+1} \geq 5 \\ 0 & \text{if } d_{k+1} < 5 \end{cases} \quad (d_1 \in [1, 9], d_i \in [0, 9])$$

(in case $\delta = 1$, $d_i = d_i + 1$ until i_0 where $d_{i_0} + 1 \leq 9$)

We explain further the IEEE 754 standard:

- Smallest possible number that can be represented:
(0s for every bit) (No, $0 \leq 1 \times 1$, we want smallest 1×1)

$$\begin{cases} s = 0 \\ f = 0 \end{cases} \rightarrow 1 \cdot 2^{0-1023} \cdot (1+0) = \boxed{2^{-1023}}$$

- Number closest to 1 that can be represented: ($x \neq 1$)

$$(-1)^0 \cdot 2^{1023-1023} \cdot \left(1 + \frac{1}{2^{52}}\right) = \boxed{1 + 2^{-52}} \text{ and}$$

we call:

- Epsilon of the machine; $\epsilon_{ps} = 2^{-52} \approx 2.22 \cdot 10^{-16}$
- Greatest number that can be represented:
(1s for every bit) (but positive number $\Rightarrow s = 0$)

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

$$\begin{cases} s = 1 \\ c = 2047 \\ f = 2^{-1} + 2^{-2} + \dots + 2^{-52} \end{cases} \rightarrow (-1)^0 \cdot 2^{2047-1023} \cdot \left(1 + \sum_{i=1}^n 2^{-i}\right) =$$

$$\stackrel{(1)}{=} 2^{1023} \cdot 2 \cdot \left(1 - 2^{-53}\right) = \boxed{2^{1023} \cdot (2 - 2^{-52})}$$

In (1) we use:

$$\left. \begin{array}{l} S_n = r^0 + \dots + r^n \\ S_n \cdot r = r^1 + \dots + r^n + r^{n+1} \end{array} \right\} S_n (1-r) = r^0 - r^{n+1} \Rightarrow r = \frac{1}{2}$$

$$\left. \begin{array}{l} \Rightarrow S_n = \sum_{i=0}^n 2^{-i} = \frac{1 - 2^{-53}}{1 - \frac{1}{2}} = 2 \cdot \left(1 - 2^{-53}\right) = (2 - 2^{-52}) \end{array} \right\}$$

Representation of certain numbers:

One: 1

$$\left\{ \begin{array}{l} S = 0 \\ C = 1023 \\ f = 0 \end{array} \right. \rightarrow (-2)^0 \cdot 2^{\frac{1}{1023}} \cdot (1+g) = 1 ,$$

Zero: 0

$$\left\{ \begin{array}{l} S = \{ 0 \\ 1 \} \\ C = 0 \\ f = 0 \end{array} \right. \rightarrow \pm 2^{0-1023} \cdot (1+0) = \pm 2^{-1023} \equiv \pm 0 ,$$

(All 0s in the 64 bit representation).

Infinity: $\pm \infty$

$$\left\{ \begin{array}{l} S = \{ 0 \\ 1 \} \\ C = 2046 \\ f = 0 \end{array} \right. \rightarrow \pm 2^{2046-1023} \cdot (1+0) = \pm 2^{1023} \equiv \pm \infty ,$$

NaN: (Non-Arithmetic Numbers; Not a Number; wrong number)

Examples of NaN: $\sin(\infty), \cos(\infty), \%, \% \%, \dots$

$$\left\{ \begin{array}{l} S = \{ 0 \\ 1 \} \\ C = 2046 \\ f \neq 0 \end{array} \right. \rightarrow \pm 2^{2046-1023} (1+f) = \pm 2^{1023} (1+f) \equiv NaN ,$$

Denormalized numbers: (Denormal numbers)

$$\left\{ \begin{array}{l} S = \{ 0 \\ 1 \} \\ C = 0 \\ f \neq 0 \end{array} \right.$$

$$\rightarrow \pm 2^{-1023} (1+f) \rightarrow$$

Numbers greater than 0, and more important, smaller than 2^{-1022}

I.E:

$$0 | \underbrace{0000000000}_{11} 0 | \underbrace{11000 \dots 0}_{52}$$

Instead of having $(1+f)$, this numbers have only (f) so that they can be smaller; (also, they are 2^{-1022}).

$$2^{-1023} (2^{-1} + 2^{-2}) = 0.75 \cdot 2^{-1022}$$

Here is the explanation:

$$1) 01000000000100 \dots$$

$$2^{1-1023} = 2^{-1022}$$

$$2) 010000000000100 \dots$$

$$2^{0-1023} \equiv 0$$

To represent $2^{-1023} < x < 2^{-1022}$, we adopt the system representation). (which is versatile and use the following

$$3) 010000000001 * \text{not zero} *$$

$$2^{-1022} \cdot (f)$$

Note this is only if $f \neq 0$, if $f=0$, then $x = 2^{-1023} = 0$

Smallest possible denormal number:

$$\begin{cases} S = \{0 \\ 1\} \\ C = 0 \\ f = 0 \dots 1 \end{cases} \rightarrow 2^{-1022} \cdot 2^{-52} //$$

Largest possible denormal number:

$$\left\{ \begin{array}{l} S = \begin{cases} 0 \\ 1 \end{cases} \\ C = 0 \\ f = 11\dots1 \end{array} \right. \rightarrow \pm 2^{-1022} \cdot \underbrace{\left(\frac{1}{2} + \frac{1}{4} + \dots + \frac{1}{2^{52}} \right)}_{(\text{explanation page 21})} = \pm$$

$$= \pm 2^{-1022} \cdot \left(1 - \frac{1}{2^{52}} \right)$$

Rounding

- Rounding towards nearest, ties to even:

- $0.\underbrace{1101}_\text{mantissa}01 \rightarrow 0.1101$

- $0.\underbrace{110111}_\text{} \rightarrow 0.111$

- $0.\underbrace{11011}_\text{} \rightarrow \begin{cases} 0.1101 \\ 0.1110 \end{cases} \leftarrow \text{default option}$

(Both options are equally close, we choose the significant lost bit to be 0, we tie to even).

- Rounding towards 0
- Rounding towards $+\infty$
- Rounding towards $-\infty$
- Recent option:

It is used for decimal, not to binary. It consists of rounding to nearest, ties away from zero. (Example, 4 decimal places).

- $0.\overline{23474} \rightarrow 0.2374$
- $0.23476 \rightarrow 0.2378$
- $0.23475 \rightarrow 0.2378$

} if last $f \in \{0, \dots, 4\}$,
we truncate;
} if last $f \in \{5, \dots, 9\}$,
we add +1 to previous.

Remember: (generally, two options).

→ chopping or truncating

→ rounding

Now, consider a five-digit floating point representation with chopping.

$$\text{fl}(1.23578) = \text{fl}(\underline{0.123578} \cdot 10^2) = \\ = 0.12357 \cdot 10^2 = 1.2357$$

Consider the same number, but rounding instead of chopping.

$$\text{fl}(1.23578) = \text{fl}(\underline{0.123578} \cdot 10^2) = \\ = 0.12358 \cdot 10^2 = 1.2358$$

Def.:

Let p^* be an approximation of $p \in \mathbb{R}$. The **absolute error** is defined as:

$$\text{Abs. error.}: |p - p^*|$$

And the **relative error**:

$$\text{Relative. error.}: \frac{|p - p^*|}{|p|}, \quad p \neq 0$$

Def.:

p^* approximates to p to t significant digits if and only if (if and only if = iff) t is the largest nonnegative number such that (such that = s.t.).

$$\frac{|p - p^*|}{|p|} \leq 5 \cdot 10^{-t}$$

(Such that the rel. error $\leq 5 \cdot 10^{-t}$).

$$\begin{aligned} 0'01 &\leq 0'05 = 5 \cdot 10^{-2} \\ \rightarrow & 2 \text{ digits} \\ 0'06 &> 0'05 \\ 0'06 &\leq 0'5 = 5 \cdot 10^{-1} \\ \rightarrow & 1 \text{ digit} \\ 0'05 &\leq 5 \cdot 10^{-2} \rightarrow 2 \text{ digit} \end{aligned}$$

Example:

Suppose we consider k -digit chopping arithmetic. Error?

$$y = 0.d_1 d_2 \dots d_k d_{k+1} \dots \cdot 10^n \quad (d_1 \neq 0)$$

Then

$$g_p(y) = 0.d_1 d_2 \dots d_k \cdot 10^n = y^*$$

The error:

$$\frac{|y - g_p(y)|}{|y|} = \frac{|0.0\dots 0 d_{k+1} \dots \cdot 10^n|}{|y|} = \frac{|0.d_{k+1} \dots \cdot 10^{n-k}|}{|0.d_1 \dots d_k \dots \cdot 10^n|}$$

$$= \frac{|0.d_{k+1} d_{k+2} \dots|}{|0.d_1 d_2 \dots d_k d_{k+1} \dots|} \cdot 10^{-k} \leq \frac{1}{0.1} \cdot 10^{-k} = 10^{1-k} <$$

$< 5 \cdot 10^{1-k} \Rightarrow k-1$ significant digits.

Note:

$$0.d_{k+1} d_{k+2} \dots \leq 1$$

26

$$\text{Since } d_1 \neq 0 \Rightarrow 0.d_1 d_2 \dots d_k d_{k+1} \dots \geq 0.1 \Rightarrow \frac{1}{0.d_1 \dots d_k \dots} \leq \frac{1}{0.1}$$

Homework:

Compute the number of significant digits when using k-digit rounding arithmetic. (Note, it depends if $0 \leq d_{k+1} \leq 4$ or $5 \leq d_{k+1} \leq 9$).

In general, when using a computer, all the operations are rounded, chopped, etc. So truly:

$$a + b + c \xrightarrow{\text{computer}} \text{fl}(a) + \text{fl}(b) + \text{fl}(c)$$

- Because of floating point approximation and so, associativity may not hold true for lots of operation.

In math:

$$(a + b) + c = a + (b + c)$$

But in computer science:

$$(a + b) + c = \text{fl}(\text{fl}(\text{fl}(a) + \text{fl}(b)) + \text{fl}(c))$$

$$a + (b + c) = \text{fl}(\text{fl}(a) + \text{fl}(\text{fl}(b) + \text{fl}(c)))$$

And for certain cases it doesn't hold true (for most of them because of approximation).

Example next page.

Example:

$$ax^2 + bx + c = 0$$

$$a \left(x^2 + \frac{b}{a}x + \frac{c}{a} \right) = 0$$

$$a \left(\left(x + \frac{b}{2a} \right)^2 - \frac{b^2}{4a^2} + \frac{c}{a} \right) = 0$$

$$\left(x + \frac{b}{2a} \right)^2 = \frac{b^2}{4a^2} - \frac{c}{a} = \frac{b^2 - 4ac}{4a^2}$$

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

$$\begin{aligned} &= \frac{(-b \pm \sqrt{b^2 - 4ac})}{2a} \cdot \frac{(-b \mp \sqrt{b^2 - 4ac})}{(-b \mp \sqrt{b^2 - 4ac})} = \frac{(-b)^2 - (\sqrt{b^2 - 4ac})^2}{2a(-b \mp \sqrt{b^2 - 4ac})} = \\ &= \frac{4ac}{2a(-b \mp \sqrt{b^2 - 4ac})} = \frac{2c}{-b \mp \sqrt{b^2 - 4ac}} \end{aligned}$$

Example:

Let us consider: $x^2 + 62.10x + 1 = 0$ (with calculator)

$$x_1 = -0.01610723741$$



$$x_2 = -62.08389276$$

Let's approximate x_1 and x_2 by using 4-digit rounding

In general try to avoid sums of odds with different signs and very different sizes.

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2}$$

1) $\text{fl} \left(\sqrt{b^2 - 4ac} \right)$

$$62 \cdot 10^2 = 3856.41 = 0. \underbrace{385641}_{\text{fl}(62 \cdot 10^2)} \cdot 10^4 \Rightarrow$$

$\Rightarrow \text{fl}(62 \cdot 10^2) = 3856$

2)

$$\begin{aligned} \text{fl}(\sqrt{b^2 - 4ac}) &= \text{fl}(\sqrt{\text{fl}(62 \cdot 10^2) - 4}) = \\ &= \text{fl}(\sqrt{3852}) = \text{fl}(62.0644) = 62.06 \end{aligned}$$

So:

$$\begin{aligned} \bullet x_1 &= \text{fl} \left(\frac{\text{fl}(-62 \cdot 10 + 62.06)}{2} \right) = \text{fl} \left(\frac{\text{fl}(-0.04)}{2} \right) = \\ &= \text{fl}(-0.02) = -0.02 \end{aligned}$$

$$\begin{aligned} \bullet x_2 &= \text{fl} \left(\frac{\text{fl}(-62 \cdot 10 - 62.06)}{2} \right) = \text{fl} \left(\frac{\text{fl}(-124.16)}{2} \right) = \\ &= \text{fl} \left(\frac{-124.2}{2} \right) = \text{fl}(-62.1) = -62.1 \end{aligned}$$

Let's compare the relative errors.

$$\begin{cases} x_1 = -0.01610723 ; \text{ fl}(x_1) = -0.02 \\ x_2 = -62.08390 ; \text{ fl}(x_2) = -62.1 \end{cases}$$

$$\frac{|x_1 - \text{fl}(x_1)|}{|x_1|} \approx 2.4 \cdot 10^{-1} \leq 5 \cdot 10^{-1} \rightarrow 1 \text{ significant digit} \\ (\text{1 is correct}).$$

$$\frac{|x_2 - \text{fl}(x_2)|}{|x_2|} \approx 3.2 \cdot 10^{-4} \leq 5 \cdot 10^{-4} \rightarrow 4 \text{ significant digits} \\ (4 \text{ correct digits}).$$

This can be justified because we had to subtract for x_1 which is bad for precision.

Let's use the other formula for x_1 :

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2c} \Rightarrow \text{fl}(x_1) = \text{fl}\left(\frac{2}{\text{fl}(-62.1) - \underbrace{\text{fl}(-62.06)}_{\text{from before}}}\right) = \\ = \text{fl}\left(\frac{2}{\text{fl}(-124.16)}\right) = \text{fl}\left(\frac{2}{-124.2}\right) = -0.0161$$

$$\Rightarrow \frac{|x_1 - \text{fl}(x_1)|}{|x_1|} \approx 6.2 \cdot 10^{-4} \leq 5 \cdot 10^{-3} \rightarrow 3 \text{ significant digits} \\ (3 \text{ correct}).$$

So there's a difference, a noticeable one. Note that for the second case, despite having 4 numbers equal, our system identifies three equal, that's because we defined it this way.

Note; for decimal machines: $(-1)^s \cdot 10^x \cdot (1+g)$,
 but the rest remains equal.

Fixed point representation:

We reserve certain bits for the integer and other for the mantissa.

The habitual thing is a floating point representation (as we have already studied for IEEE 754).

Exercise: solution.

Let $y = 0.d_1 \dots d_k d_{k+1} \dots \cdot 10^n$ and $fl(y)$ the digit (k -digit) approximation using rounding arithmetic. Obtain f significant digits.

$$fl(y) = \begin{cases} 0.d_1 \dots d_k \cdot 10^n & \text{if } d_{k+1} < 5 \\ (0.d_1 \dots d_k + 1 \cdot 10^{-k}) \cdot 10^n & \text{if } d_{k+1} \geq 5 \end{cases} \quad \begin{array}{l} d_1 \neq 0 \text{ because} \\ \text{we normalize} \\ \text{the number.} \end{array}$$

If $d_{k+1} < 5$:

$$\left| \frac{y - fl(y)}{y} \right| = \left| \frac{0.d_1 d_2 \dots d_k d_{k+1} \dots - 0.d_1 d_2 \dots d_k}{0.d_1 d_2 \dots d_k d_{k+1} \dots} \right| \cdot \frac{10^n}{10^n} =$$

$$= \left| \frac{0.d_{k+1} \dots \cdot 10^{-k}}{0.d_1 \dots d_{k+1} \dots} \right| \stackrel{(2)}{\leq} \frac{0.5}{0.1} \cdot 10^{-k} = 5 \cdot 10^{-k} \Rightarrow \boxed{k \text{ digits of precision.}}$$

(Because $d_{k+1} < 5$ and $d_1 \neq 0$)

if $d_{k+1} \geq 5$

$$\begin{aligned} \left| \frac{y - fl(y)}{y} \right| &= \left| \frac{0.d_1 d_2 \dots d_k d_{k+1} \dots - (0.d_1 \dots d_k + 10^{-k})}{0.d_1 d_2 \dots d_k d_{k+1} \dots} \right| \cdot \frac{10^k}{10^n} \\ &= \left| \frac{0.d_{k+1} \dots \cdot 10^{-k} - 10^{-k}}{0.d_1 \dots d_k \dots} \right| = \left| \frac{0.d_{k+1} \dots - 1}{0.d_1 \dots d_k \dots} \right| \cdot 10^{-k} \leq \\ &\leq \left| \frac{0.5 - 1}{0.1} \right| \cdot 10^{-k} = 5 \cdot 10^{-k} \Rightarrow \boxed{k \text{ correct digits.}} \end{aligned}$$

Example:

$$f(x) = x^3 - 6.1x^2 + 3.2x + 1.5, \quad x = 4.71$$

Consider 3 digit rounding arithmetic: $(fl(x) = -14.263899)$

$$\begin{aligned} f(x) &= x^2(x - 6.1) + 3.2x + 1.5 = \left(\begin{array}{l} \text{They could give us or} \\ \text{we would have to} \\ \text{compute it with the} \\ \text{calculator} \end{array} \right) \\ &= x(3.2 + x(x - 6.1)) + 1.5 \\ &= ((x - 6.1)x + 3.2)x + 1.5 \quad (\rightarrow \text{Nested manner}). \end{aligned}$$

$$a = fl(x - 6.1) = fl(4.71 - 6.1) = fl(-1.39) = -1.39$$

$$b = fl(a \cdot x) = fl(-1.39 \cdot 4.71) = fl(-6.5469) = -6.55$$

$$c = fl(b + 3.2) = fl(-6.55 + 3.2) = fl(-3.35) = -3.35$$

$$d = fl(c \cdot x) = fl(-3.35 \cdot 4.71) = fl(-15.7785) = -15.8$$

$$e = fl(d + 1.5) = fl(-15.8 + 1.5) = fl(-14.3) = -14.3$$

$$\rightarrow \boxed{f(4.71) = -14.3}$$

Use calculator
for the middle
steps.

Relative error and number of significant digits:

$$\frac{|y - \text{fl}(y)|}{|y|} = \frac{|-14.263899 + 14.3|}{|-14.263899|} = \boxed{\frac{215 \cdot 10^{-3}}{5 \cdot 10^{-3}}} \leq 5 \cdot 10^{-3}$$

$\underbrace{215 \cdot 10^{-3}}_{\text{relative error}}$

\Rightarrow 3 significant digits.

This method is much better than doing the operations from left to right, we perform less operations and it's more precise. (With normal method $\rightarrow f(x) = -13$. something and 1 significant digit).

• Usual method:

$$a = \text{fl}(x) = 4.71$$

$$b = \text{fl}(x \cdot \text{fl}(x)) = \text{fl}(4.71 \cdot 4.71) = \text{fl}(22.1841) = 22.2$$

$$c = \text{fl}(x \cdot \text{fl}(x \cdot \text{fl}(x))) = \text{fl}(4.71 \cdot 22.2) = \text{fl}(104.562) = 105$$

$$d = \text{fl}(-6.1 \cdot b) = \text{fl}(-6.1 \cdot 22.2) = \text{fl}(-135.42) = -135$$

$$e = \text{fl}(c - d) = \text{fl}(105 - 135) = -30$$

$$f = \text{fl}(3.2 \cdot a) = \text{fl}(4.71 \cdot 3.2) = \text{fl}(15.072) = 15.1$$

$$g = \text{fl}(e + f) = \text{fl}(-30 + 15.1) = \text{fl}(-14.9) = -14.9$$

$$\boxed{h(x) = \text{fl}(g + 1.5) = \text{fl}(-14.9 + 1.5) = \text{fl}(-13.4) = -13.4}$$

Relative error:

(next page)

$$\left| \frac{y - fl(y)}{y} \right| = \left| \frac{-14.263899 + 13.4}{-14.263899} \right| = \boxed{6.05654 \cdot 10^{-2} \leq 5 \cdot 10^{-1}}$$

1 correct digit.

Consider 3 digit Chopping arithmetic: (Nested manner)

$$f(x) = ((x - 6.1) \cdot x + 3.2) \cdot x + 1.5$$

$$a = fl(x - 6.1) = fl(-1.39) = -1.39$$

$$b = fl(a \cdot x) = fl(-6.5469) = -6.54$$

$$c = fl(b + 3.2) = fl(-3.34) = -3.34$$

$$d = fl(c \cdot x) = fl(-15.7314) = -15.7$$

$$e = fl(d + 1.5) = fl(-14.2) = -14.2$$

$$\boxed{f(4.71) = -14.2}$$

$$\left| \frac{y - fl(y)}{y} \right| = \left| \frac{-14.263899 + 14.2}{-14.263899} \right| = \boxed{4.479 \cdot 10^{-3} \leq 5 \cdot 10^{-3}}$$

3 significant digits

Here it wasn't the case, but usually chopping is less precise than rounding, as we proved before.

Let's compute some matrix algorithms now.