# Numbers, errors and computers.

Calcoli di Processo dell' Ingegneria Chimica

Timoteo Dinelli

07th of October 2025.

Department of Chemistry, Materials and Chemical Enginering, "Giulio Natta", Politecnico di Milano.

email: timoteo.dinelli@polimi.it

# Numbers representation.

"*In computing, floating-point arithmetic (FP) is arithmetic that represents subsets of real numbers using an integer with a fixed precision, called the significand, scaled by an integer exponent of a fixed base. Numbers of this form are called floating-point numbers. For example, 12.345 is a floating-point number in base ten with five digits of precision.*" Wikipedia.

$$12.345 = \boxed{12345} \times \boxed{10^{-3}}$$

Significand (Mantissa)                          Exponent

Additional RECOMMENDED read:

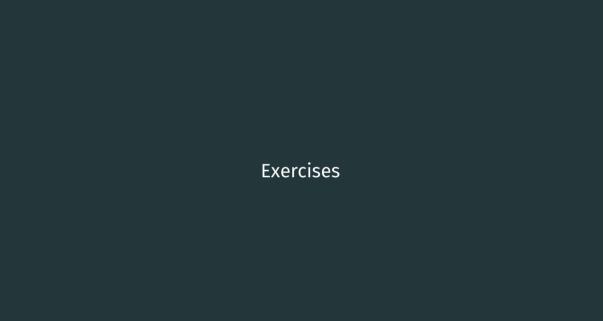▶ Floating point representations.
▶ LLM quantization.

## Elementary operations

Initialize a single precision variable in MATLAB (e.g. $x = 1e^{+25}$) using the function:
$$x = \text{single}(1.e+25)$$
Working in single precision predict and calculate from $x = 1.e+25$ and $y = 1.e+18$ the following values of z:

|  | Single Precision | Double Precision |
|---|---|---|
| $z = x * y$ | inf | 1.0000e+43 |
| $z = x/y$ | 10000000 | 10000000 |
| $z = y/x$ | 1.0000e-07 | 1.0000e-07 |
| $z = x^2$ | inf | 1.0000e+50 |
| $z = y^2$ | 1.0000e+36 | 1.0000e+36 |
| $z = 1./(x * y)$ | 0 | 1.0000e-43 |
| $z = 1./x/y$ | 9.9492e-44 | 1.0000e-43 |
| $z = y + 1e10$ | 1.0000e+18 | 1.0000e+18 |
| $z = x * y/(x * y + 1)$ | NaN | 1 |

Exercises

# Compute the Machine Epsilon (MACHEPS)

Machine epsilon (denoted as $\varepsilon_{\text{mach}}$) is the smallest positive floating-point number such that $1.0 + \varepsilon_{\text{mach}} \neq 1.0$ in floating-point arithmetic. It represents the relative precision of floating-point numbers and is a fundamental constant for a given floating-point format (e.g., approximately $2.22 \times 10^{-16}$ for double precision).

Machine epsilon is important because it quantifies the limit of precision in numerical computations and helps estimate round-off errors in calculations.

▶ The algorithm iterates by halving epsilon until $1.0 + \varepsilon/2$ is indistinguishable from 1.0 in floating-point arithmetic. At this point, we've found the machine epsilon.

# Implementation

```matlab
1 function epsilon = macheps_implementation()
2     epsilon = 1.0; % Initialize epsilon to 1
3     while (1.0 + epsilon/2) > 1.0
4         % Iterate as long as 1.0 + epsilon/2 is distinguishable from 1.0
5         epsilon = epsilon / 2; % Halve epsilon at each iteration
6     end
7     % At the end of the loop, epsilon is the machine epsilon
8     % (the smallest value where 1.0 + epsilon > 1.0)
9 end
```

Usage:

```matlab
eps_computed = macheps_implementation()
eps_builtin = eps(1.0)  % Compare with built-in function
```

# Sum of the inverse of numbers (order matters!)

Write a script which computes:

$$\sum_{n=1}^{1000000} \frac{1}{n}$$

in single and double precision. Then compare with the results obtained inverting the order of the sum, so by computing:

$$\sum_{n=1000000}^{1} \frac{1}{n}$$

# Implementation

```matlab
1 % Sum from 1 to 1000000 (forward)
2 sum_forward_single = single(0);
3 sum_forward_double = 0;
4 for n = 1:1000000
5     sum_forward_single = sum_forward_single + single(1/n);
6     sum_forward_double = sum_forward_double + 1/n;
7 end
8
9 % Sum from 1000000 to 1 (backward)
10 sum_backward_single = single(0);
11 sum_backward_double = 0;
12 for n = 1000000:-1:1
13     sum_backward_single = sum_backward_single + single(1/n);
14     sum_backward_double = sum_backward_double + 1/n;
15 end
```

```matlab
16
17 % Display results
18 fprintf('Forward sum (single):  %.10f\n', sum_forward_single);
19 fprintf('Forward sum (double):  %.15f\n', sum_forward_double);
20 fprintf('Backward sum (single): %.10f\n', sum_backward_single);
21 fprintf('Backward sum (double): %.15f\n', sum_backward_double);
22
23 fprintf('\nDifferences:\n');
24 fprintf('Single precision:  %.10e\n', ...
25     abs(sum_forward_single - sum_backward_single));
26 fprintf('Double precision:  %.15e\n', ...
27     abs(sum_forward_double - sum_backward_double));
```

# Vancouver: a nickel at a time

Analogously to what happened on the Vancouver stock market (reference), starting from a stock value of 1000, check what happens when a random variation of $\pm 1\%$ in its value is iterated for 10000 times. Try to round or truncate to two decimal places using floor, ceil, and round. Compare the results with the number obtained using full computer precision.

Please verify in the MATLAB help how the functions *rand, floor, ceil, fix, round* work.

# Solution - Part 1: Setup

```matlab
1  % Initial stock value
2  initial_value = 1000;
3  n_iterations = 10000;
4
5  % Initialize values for different rounding methods
6  value_full_precision = initial_value;
7  value_floor = initial_value;
8  value_ceil = initial_value;
9  value_round = initial_value;
10
11  % Set random seed for reproducibility (optional)
12  rng(42);
```

# Solution - Part 2: Iteration Loop

```matlab
13  % Iterate 10000 times with random ±1% variation
14  for i = 1:n_iterations
15      % Generate random variation: ±1%
16      % rand() gives [0,1], so 2*rand()-1 gives [-1,1]
17      variation = (2 * rand() - 1) * 0.01;
18
19      % Full precision (no rounding)
20      value_full_precision = value_full_precision * (1 + variation);
```

# Solution - Part 3: Rounding Methods

```matlab
21      % Floor (round down to 2 decimal places)
22      temp = value_floor * (1 + variation);
23      value_floor = floor(temp * 100) / 100;
24
25      % Ceil (round up to 2 decimal places)
26      temp = value_ceil * (1 + variation);
27      value_ceil = ceil(temp * 100) / 100;
28
29      % Round (round to nearest 2 decimal places)
30      temp = value_round * (1 + variation);
31      value_round = round(temp * 100) / 100;
32  end
```

# Solution - Part 4: Display Results

```matlab
34 % Display results
35 fprintf('Results after %d iterations:\n', n_iterations);
36 fprintf('=======================================\n');
37 fprintf('Initial value:       %.2f\n', initial_value);
38 fprintf('Full precision:      %.6f\n', value_full_precision);
39 fprintf('Floor (round down):  %.2f\n', value_floor);
40 fprintf('Ceil (round up):     %.2f\n', value_ceil);
41 fprintf('Round (to nearest):  %.2f\n', value_round);
42 fprintf('=======================================\n');
```

# Solution - Part 5: Differences Analysis

```matlab
44 % Calculate differences from full precision
45 fprintf('\nDifferences from full precision:\n');
46 fprintf('Floor:  %.6f (%.2f%%)\n', ...
47     value_full_precision - value_floor, ...
48     100*(value_full_precision - value_floor)/value_full_precision);
49 fprintf('Ceil:   %.6f (%.2f%%)\n', ...
50     value_full_precision - value_ceil, ...
51     100*(value_full_precision - value_ceil)/value_full_precision);
52 fprintf('Round:  %.6f (%.2f%%)\n', ...
53     value_full_precision - value_round, ...
54     100*(value_full_precision - value_round)/value_full_precision);
```

# Expected observations

- ▶ Floor will consistently lose value (downward bias, similar to Vancouver bug)
- ▶ Ceil will consistently gain value (upward bias)
- ▶ Round should be closer to full precision (unbiased rounding)
- ▶ The cumulative effect over 10,000 iterations will be significant, demonstrating the Vancouver Stock Exchange issue where the index lost about 50% of its value due to truncation errors!

Thank you for the attention!