

NUR A lecture 1

In this lecture:

Introduction to the course
Integer representation
Integer arithmetic and logic
Binary operations
Floating point numbers
Sources of numerical error
Efficiency vs accuracy

Numerical Recipes for Astrophysics A

Spring semester 2025

Things to know

- Lecturer: dr. Marcel van Daalen
- Teaching assistants: Maria Marinichenko & Konstantinos Kilmatis
- Course website: Brightspace
- Please do the 20-minute programming test on Brightspace!
- Outside class hours you can ask questions on the discussion boards
- Note: two lectures this week, two tutorials + one lecture next week

The screenshot shows the course navigation menu with the following items:

- Course Home
- Content
- Grades
- Course Tools ▾
- Help ▾

Below the menu is a search bar labeled "Search titles, descriptions". The main content area contains the following sections:

- General Information
- Prior Knowledge
- Assignments
- 20-minute test of student level** (This item is highlighted with a blue border and has a question mark icon next to it.)
- Discussion boards

A callout bubble on the right side of the screen says "View the quiz in the Quizzing tool." and includes a "Go to Quiz" button. There is also a circular icon with a document and a question mark inside.

Things to know

- The lectures discuss theory and specific algorithms
- After today, we start each lecture/half with a Mentimeter quiz (menti.com)
- During the tutorials you will solve problems by coding up the algorithms discussed (collaboration with each other is encouraged – but not with LLMs)
- Three larger exercises should be handed in (individually!) and count towards 50% of your grade (total)
- Participating in the tutorials means the hand-ins (and studying for the exam) are MUCH less work!
- At the end of the course a theory exam will determine the final 50% of your grade

Introduction

- Algorithms A series of steps to get from A to Z
 - E.g. Runge-Kutta
 - Sieve of Eratosthenes
 - long division
- Broad scientific libraries: GSL and SciPy
 - C T
 - Python
- Stability, accuracy and efficiency
 - "how reliable"
 - "how close"
 - "how quickly"
 - Reduce step for GPC

Topics

- L1: Numbers and arithmetic in computers
 - L2: Inter- and extrapolation
 - L3: Solving linear equations
 - L4: Function integration
 - L5: Differentiation
- L5: Random numbers and sampling
 - L6: Sorting and selection
 - L6: Root-finding/intersection
- L7: Minimization/maximization/optimization
- L8 & L9: Fitting and modelling data
 - L9: Statistical tests
 - (L10: Q&A)

Numbers and arithmetic in computers

NR Ch 1

Bits, bytes and ints

- Binary representation

$$32\text{-bit (4 bytes)}: 2^{32} \approx 4.3 \times 10^9$$

$$64\text{-bit (8 bytes)}: 2^{64} \approx 1.8 \times 10^{19}$$

Bit: 0, 1 → byte = 8 bits $\rightarrow 2^8$ values

2^7	2^6	2^5	2^4	2^3	2^2	2^1	2^0
0	0	1	0	1	1	0	1

values: $[0, 255]$ or $[-128, 127]$

- Integer addition and multiplication

$$\begin{array}{r} 00001011 \quad (11) \\ 000000110 \quad (6) \\ \hline 00010001 \quad (17) \\ \downarrow \\ 11111111 \quad (255) \\ \times 00000001 \quad (1) \\ \hline \dots 00 \end{array}$$

Mult. x by $2^y \equiv$ shifting bits of x left by y positions

$$\begin{aligned} & 00010110 \quad (22) \times 000000101 \quad (5) = \\ & 00010110 \quad (22) + 010110 \quad \underline{00} \quad (88) = 1110 \end{aligned}$$

More integer arithmetic

- Bit-wise operators

$x \& y$: AND, 1 if both bits are 1, 0 otherwise

$x | y$: OR, 1 if at least 1 of 2 bits is 1

$x \wedge y$: XOR, 1 if exactly 1 of 2 bits is 1

- Examples

$$(7 >> 2) \ll 2 = 4$$

$$(4 \wedge 13) | 10 = 11$$

Bit-shifts: $x \ll y \quad (\equiv x \cdot 2^y)$

$$x \gg y \quad (\equiv x / 2^y)$$

$\sim x$: NOT

↳ Complements of x
 $(0 \leftrightarrow 1)$

$$(3 \ll 2) \& \sim 5 = 8$$

$$x \% 2^y = x \& (2^y - 1)$$

More integer arithmetic

- Negative numbers

$$-X = \sim(X-1) = \sim X + 1 \quad (2\text{'s complement})$$

$$-X + X = 0 \rightarrow \sim X + 1 + X = 0$$

- Subtraction and division

$$Y - X = Y + (-X) = Y + \sim X + 1$$

$X/2^y$ = shift X right by y pos

in general:

$$\begin{array}{r} 00010110 \\ 01100 \\ \hline 01010 \\ 00110 \\ \hline 0100 \\ 0011 \\ \hline 0001 \end{array}$$

(22)/00000011 (3) = $2^2 + 2^1 + 2^0 = 7$

Representing floating point numbers

... $\boxed{2^2 | 2^1 | 2^0 | 2^{-1} | 2^{-2} | \dots}$...? no: inefficient \Rightarrow scientific notation

base 10: $a \times 10^b$, $a = X.xxxxx\dots$

\Rightarrow base 2: $1 \frac{\text{Fraction } F}{\text{mantissa } M} \cdot 2^E$ $E \leftarrow \text{exponent(int)}$

S-bit: $(-1)^S$

$1.\overbrace{xxxx\dots}^{\text{mantissa } M}$

F: $\boxed{2^{-1} | 2^{-2} | 2^{-3} | \dots}$

$2^{[-126, 127]}$

IEEE 754

32-bit (float, single-pre.): $\begin{cases} S: 1 \text{ bit} \\ E: 8 \\ F: 23 \end{cases}$

$1.625 \rightarrow E=0$

64-bit (double, -prec.): $\begin{cases} S: 1 \\ E: 11 \\ F: 52 \end{cases}$

$0.625 \rightarrow F = \boxed{10100\dots}$

FLOPs and reducing operations

Floating point Operation

Addition, subtraction, multiplication = 1 FLoP

Division: 4-16 FLOPs, sqrt : 6-32 FLOPs

$$x/2.0 \rightarrow x * 0.5$$

$x/9.0$ many times \rightarrow define oneNinth=1.0/9.0 \leftarrow multiply!

$$\sqrt{dx^2 + dy^2 + dz^2} < 5 \rightarrow dx^2 + dy^2 + dz^2 < 25 (= 5^2)$$

Sources of numerical error

- Overflow and underflow
 - ✓ "too large to represent" → ints: wraps around
 - ✓ "too small to represent" → floats: $\pm \infty$ (NaN , not a number, e.g. $\times/0$)

"too small to represent" → rounded to zero

$$P_1(x) = \frac{\lambda^x e^{-\lambda}}{x!} \quad (\text{single prec.})$$

$$P_5(34) \approx 1.33 \times 10^{-17} \quad (\text{correct})$$

$$P_5(35) = 0 \quad (2^{-127} \sim 10^{-39})$$

Solve: rewrite equation, order of operations,
or log space

Correct: 1.9×10^{-18}

Sources of numerical error

- Machine precision

Fractional machine accuracy: ϵ_m

$$\text{e.g. } 1.0 + 10^{-100} = 1.0$$

ϵ_m : Smallest number we can add to 1.0 and not get 1.0

ϵ_m : Single prec. $\sim 10^{-7}$, double: $\sim 10^{-16}$

e.g. 5 can be represented exactly, but 4.9 and 0.1 cannot

if $((4.9 + 0.1) == 5)$. - 

Sources of numerical error

$N = \# \text{ of operations}$

- Round-off, stability and truncation

Accumulation of machine error

unstable: error is magnified

$$\text{e.g. } (1.0002 \pm \epsilon_m) - (1.0001 \pm \epsilon_m) = 10^{-4} \pm \sim \epsilon_m$$

under our control!

$$\sum_{\text{continuous}} \rightarrow \sum_{\text{discrete}}$$

if random: frac. error $\sqrt{N} \epsilon_m$

in practice: $\sim N \epsilon_m$

frac. error $\times 10^4$!

Efficiency vs accuracy

e.g. truncation error

- At odds or hand-in-hand?

e.g. round-off

E.g. $\bar{X} = \sum_i \left(\frac{x_i}{n} \right)$ or

$$\left(\sum_i x_i \right) / n$$

in general better

- Vectorized operations/avoiding loops

Array X: $X^* = 0.2$ is better than a loop and/or divide by 5

avoid loops in Python

**Thursday:
Inter- and extrapolation**