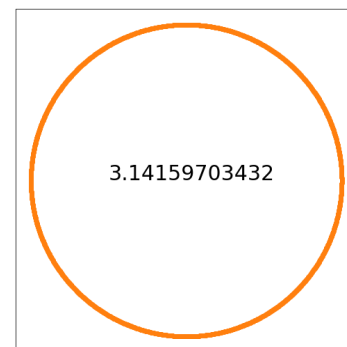
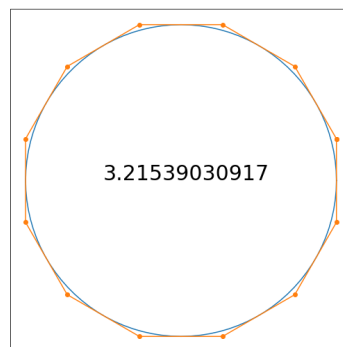
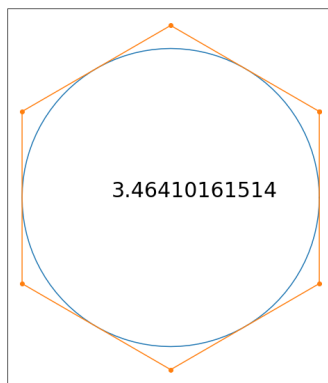


Introduction to Computational Mathematics



```
1  ## loading python libraries
2
3  # necessary to display plots inline
4  %matplotlib inline
5
6  # load the libraries
7  import matplotlib.pyplot as plt #
8  plt.rcParams.update({'font.size':
9  import numpy as np              #
10
```

Introduction

We will deal with two different types of approximation

- Mathematical approximation, or truncation errors.

$$e^x \approx 1 + x + \frac{x^2}{2} \quad \text{or} \quad f'(x) \approx \frac{f(x+h) - f(x)}{h}$$

- Approximations made by the computer, or rounding errors.

```
1 x = sqrt(2)
2 y = x*x
3 print(y)
```

2.000000000000000004

Definition. If x is an approximation of x^* ,

- the absolute error between x and x^* is $|x - x^*|$
- the relative error between x and x^* is $\frac{|x - x^*|}{|x^*|}$.

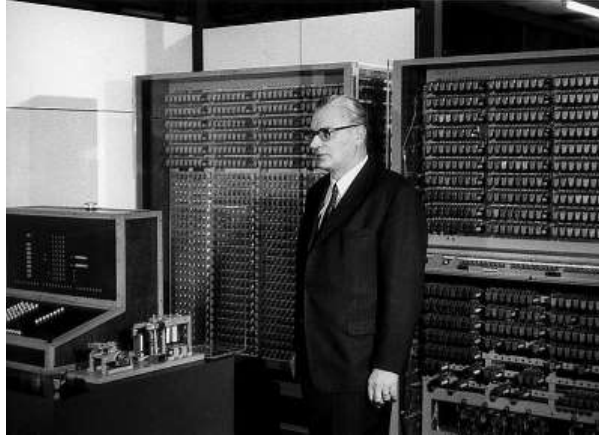
Example. Take

$$x = 11, x^* = 10 \quad \text{and} \quad y = 1001, y^* = 1000.$$

In both cases, the absolute error is 1, but one might consider that y approximates y^* more accurately than x does x^* . This is reflected in the relative error:

$$\frac{|x - x^*|}{|x^*|} = 0.1 \quad \frac{|y - y^*|}{|y^*|} = 0.001.$$

Machine representation of numbers: rounding errors



Konrad Zuse (1910-1995) and the Z3-computer (1941).

```
1 x = sqrt(2)
2 y = x**2
3 print(y)
```

2.000000000000000004

Machine representation of numbers

Example. The decimal system (or base-10 system).

$$[1000]_{10} = 1 \times 10^3$$

$$[6743.7]_{10} = 6 \times 10^3 + 7 \times 10^2 + 4 \times 10^1 + 3 \times 10^0 + 7 \times 10^{-1}$$

Example. The binary system (or base-2 system).

$$[1000]_2 = 1 \times 2^3$$

$$\lfloor 1011.1 \rfloor_2 = 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 + 1 \times 2^{-1} \quad (= \lfloor 11.5 \rfloor_{10})$$

Definition.

Normalized exponent representation in base 10.

A real number x can be written as

$$x = (-1)^{s_a} \times [0.a]_{10} \times 10^{(-1)^{s_b} [b]_{10}}$$

where the first digit of a is different from 0, and b is an integer.

Example.

$$-0.0047258 = -0.47258 \times 10^{-2} = (-1)^1 \times 0.47258 \times 10^{(-1)^1 \times 2}$$

Example.

$$\frac{10}{3} = 3.3333333 \dots = (-1)^0 \times 0.33333333 \dots \times 10^{(-1)^0 \times 1}$$

Definition.

Machine numbers.

Let m and n be given integers. A *machine number* is a real number x as above, with the length of a and b respectively lower than m and n . It can be exactly represented by the following word of size $N = n + m + 2$:

$$\begin{array}{|c|c|c|c|} \hline s_a & a & s_b & b \\ \hline \end{array}$$

Example. Take $(m, n) = (5, 2)$

$$-0.0047258 = (-1)^1 \times 0.47258 \times 10^{(-1)^1 \times 2} \quad \text{is a machine number}$$

$$\begin{array}{|c|c|c|c|} \hline 1 & 47258 & 1 & 02 \\ \hline \end{array}$$

Example. Take $(m, n) = (5, 2)$

$\frac{10}{3} = (-1)^0 \times 0.33333333 \dots \times 10^{(-1)^0 \times 1}$ is NOT a machine number

| 0 | 33333333 ... | 0 | 01 |

Definition. Floating-point representation and rounding error.

- $rd(x)$ is the closest machine number to x , called *floating-point* representation of x .
- We call $\frac{|x - rd(x)|}{|x|}$ the rounding error.
- The smallest ϵ such that $\frac{|x - rd(x)|}{|x|} \leq \epsilon \forall x$ is called the *machine precision*.

Example. Consider $m = 5, n = 2$, and $\pi = 3.14159265 \dots$

$rd(\pi) = 0.31416 \times 10^1 \Rightarrow$ | 0 | 31416 | 0 | 01 |

$$\left| \frac{\pi - rd(\pi)}{\pi} \right| \approx 2.34 \times 10^{-6}$$

```
1 print(pi)
```

```
3.141592653589793
```

```
1 x = 0.111111111111111111 # 16 ones
2 y = 0.111111111111111111 # 17 ones
3 print('Is x equal to y?', x==y)
```

```
Is x equal to y? True
```

```
1 z = 0.111111111111111111 # 15 ones
```

```
2 print('Is x equal to z?', x==z)
Is x equal to z? False
```

```
1 lie = 0.1
2 print(lie)
```

0.1

Remark.

$$[0.1]_{10} = \underbrace{[0.110011001100 \dots]_2}_{\infty \text{ many times } 1100} \times 2^{-3}$$

$$\text{lie} = rd([0.1]_{10}) = \underbrace{[0.110011001100 \dots]_2}_{\text{only 52 digits}} \times 2^{-3}$$

Floating-point arithmetic

Example.

Are the following numbers machine numbers for $m = 5$ and $n = 2$?

- $x = 3.1416$
- $y = 0.00011$
- $s = x + y$

Answer.

- $x = 0.31416 \times 10^{01}$ ✓
- $y = 0.11000 \times 10^{-03}$ ✓
- $s = x + y = 0.314171 \times 10^{01}$ ✗

Definition.

Floating-point arithmetic (finite-digits arithmetic)

$$\begin{aligned}x \oplus y &= rd(rd(x) + rd(y)), & x \ominus y &= rd(rd(x) - rd(y)) \\x \otimes y &= rd(rd(x) \times rd(y)), & x \oslash y &= rd(rd(x) / rd(y))\end{aligned}$$

Example. For $m = 5$ and $n = 2$.

- $x = 3.1416$
- $y = 0.00011$
- $s = x + y = 3.14171$
- $s_{num} = x \oplus y = rd(x + y) = 3.1417$

$$\frac{|s - s_{num}|}{|s|} \approx 3.18 \times 10^{-6}$$

Example.

Addition of a large and a small number.

- $x = 1/3$
- $y = 6/7 \times 10^4$
- $s = x + y = 8571.7619 \dots$

Answer. For $m = 5$ and $n = 2$.

$$rd(x) = 0.33333 \times 10^0$$

$$rd(y) = 0.85714 \times 10^4$$

and

$$\begin{aligned}x \oplus y &= rd(0.33333 \times 10^0 + 0.85714 \times 10^4) \\&= rd(0.85717\textcolor{red}{3333} \times 10^4) \\&= 8571.7\end{aligned}$$

Example.

Subtraction of nearly equal numbers.

- $y = 6/7 \times 10^4$
- $z = 0.85717 \times 10^4$
- $t = z - y = 0.2714285 \dots$

Answer. For $m = 5$ and $n = 2$.

$$rd(y) = 0.85714 \times 10^4$$

$$rd(z) = 0.85717 \times 10^4$$

and

$$\begin{aligned} z \ominus y &= rd(0.85717 \times 10^4 - 0.85714 \times 10^4) \\ &= 0.3 \end{aligned}$$

The relative error with the exact value is quite large $\approx 10\%$

Example. The order of operations matters for floating-point arithmetic!

- $x = 1/3$
- $y = 6/7 \times 10^4$
- $s = x + y - y = 0.333333 \dots$

Answer.

$$\begin{aligned} s_1 &= (x \oplus y) \ominus y \\ &= 0.85717 \times 10^4 \ominus y \\ &= 0.3 \quad (\text{relative error} \approx 10\%) \end{aligned}$$

Answer.

$$\begin{aligned}s_2 &= x \oplus (y \ominus y) \\ &= 0.33333 \ominus 0 \\ &= 0.33333 \quad (\text{relative error} \approx 0.001\%)\end{aligned}$$

```
1 ## x + y - x = y ?
2
3 x = 10**15
4 y = 0.1
5 s1 = (x + y) - x
6 s2 = (x - x) + y
7 s = y
8 print('s1 =',s1)
9 print('s2 =',s2)
10 print('The "exact" answer s =',s)
```

s1 = 0.125

s2 = 0.1

The "exact" answer s =
0.1

Mathematical approximations: truncation error

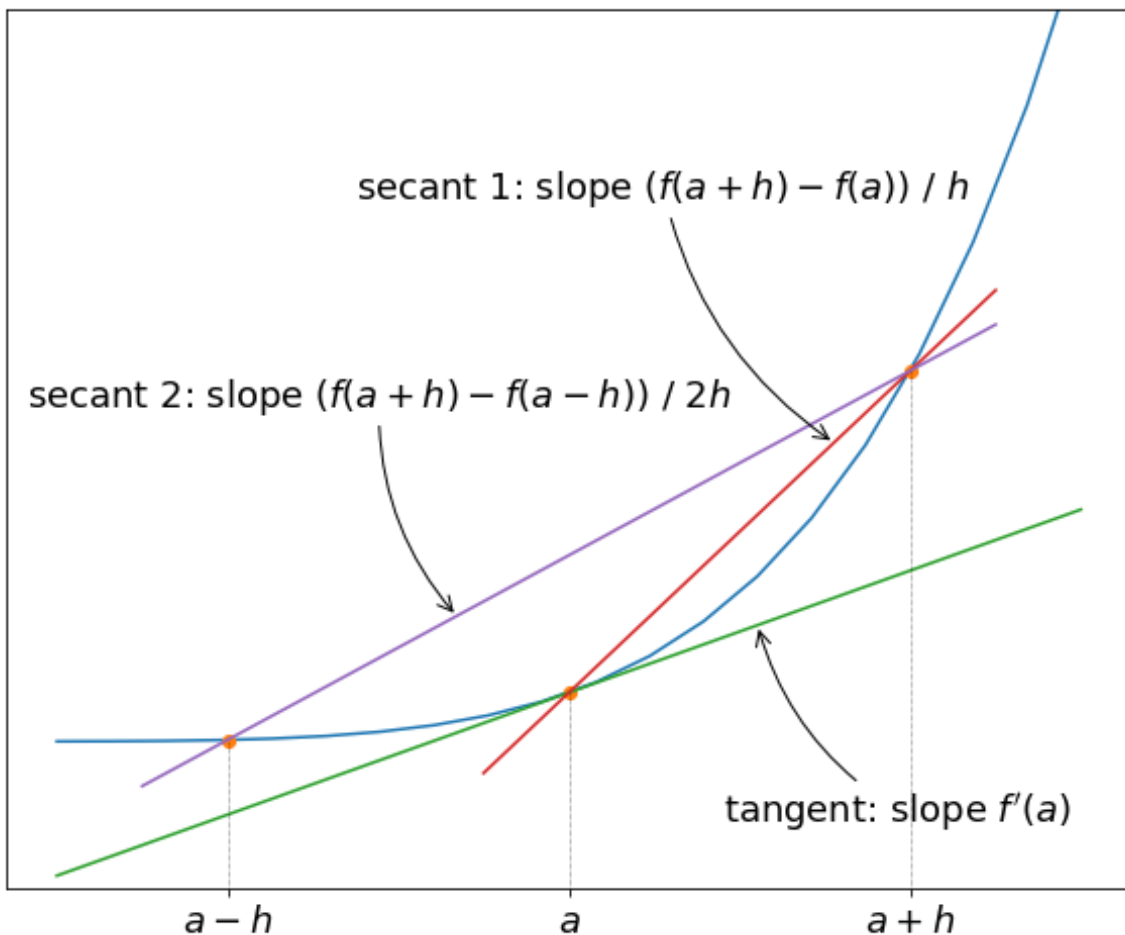
- We ignore rounding errors for the moment and assume that the computer can compute with real numbers exactly.
- In most problems, there is still an important issue: there is no explicit and computable formula for the solution.

- Therefore, we replace the initial problem by an easier one, for which we can actually compute the solution.
- This leads to mathematical approximations, or truncation errors, that also have to be controlled.

Example. Approximate computation of a derivative:

- f a given derivable function, but no formula known for f' .
- a a given real number.

How to approximate $x^* = f'(a)$



Algorithm. Two options:

$$x_h = \frac{f(a+h) - f(a)}{h} \quad \text{and} \quad \bar{x}_h = \frac{f(a+h) - f(a-h)}{2h}.$$

Definition.

Truncation error.

$$e_h = |f'(a) - x_h| \quad \text{and} \quad \bar{e}_h = |f'(a) - \bar{x}_h|.$$

Questions.

- Does each algorithm converge, meaning that $e_h, \bar{e}_h \xrightarrow{h \rightarrow 0} 0$?
- If yes, can we predict the *rate* or *speed* at which they converge?
- Can we make some quantitative prediction about the accuracy for a given h ?

```
1 def f(x):
2     return x**5
3
4 def ApproxDerivative1(f, a, h):
5     return (f(a+h) - f(a))/h
6
7 def ApproxDerivative2(f, a, h):
8     return (f(a+h) - f(a-h))/(2*h)
```

```
1 a = 1
2 h = 0.000000000000001
3 x1 = ApproxDerivative1(f, a, h)
4 x2 = ApproxDerivative2(f, a, h)
5 print('a =', a, ', h =', h, ', exact =', f'(a))
6 print('First algorithm:', x1)
7 print('Second algorithm:', x2)
```

a = 1 , h = 1e-12 , exact value: f'(1)=5
First algorithm: 5.000444502911705
Second algorithm: 5.00

Questions.

- It looks like $e_h, \bar{e}_h \xrightarrow{h \rightarrow 0} 0$. Can we prove it? Under which assumptions?
- It also looks like \bar{e}_h goes to 0 faster than e_h does. Can we also prove that?

Remark.

Usually, it is hard to exactly compute the error e_h . In practice, one often tries to find an error estimator β_h such that

- $e_h \leq \beta_h$,
- $\beta_h \xrightarrow{h \rightarrow 0} 0$.

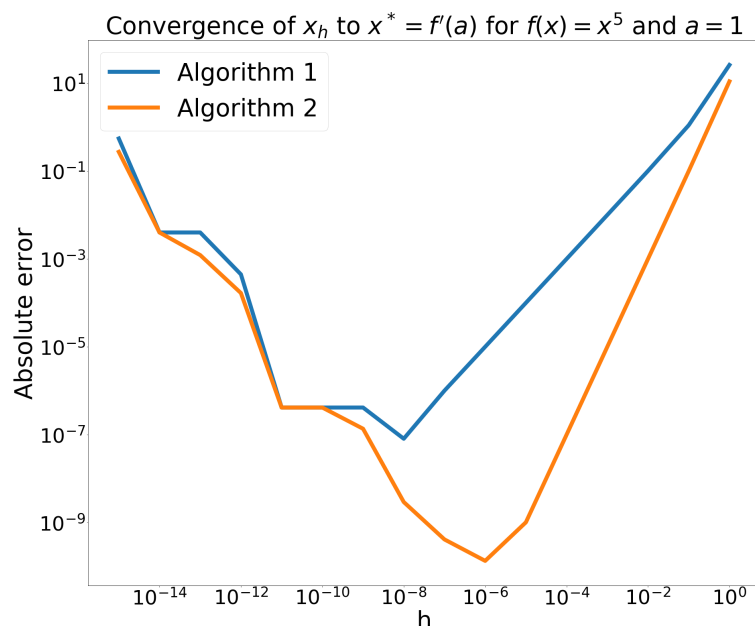
Total numerical error

```
1 a = 1
2
3 n = np.arange(16)
4 h = 10.**(-n) #sequence of h = [1
5
6 Der1 = ApproxDerivative1(f, a, h)
7 Err1 = abs(Der1 - 5.)
```

```

8
9 Der2 = ApproxDerivative2(f, a, h)
10 Err2 = abs(Der2 - 5.)
11
12 # plot of the errors versus h
13 fig = plt.figure(figsize=(30, 25))
14 plt.loglog(h, Err1, linewidth=10,
15 plt.loglog(h, Err2, linewidth=10,
16 plt.legend(loc='upper left', fonts
17 plt.xlabel('h', fontsize=60)
18 plt.ylabel('Absolute error', fonts
19 plt.title('Convergence of  $x_h$  to
20
21 plt.show()
22

```



```

1 # execute this part to modify the
2 from IPython.core.display import

```

```
3 def css_styling():  
4     styles = open("./style/custom_".  
5     return HTML(styles)  
6 css_styling()
```

```
1
```