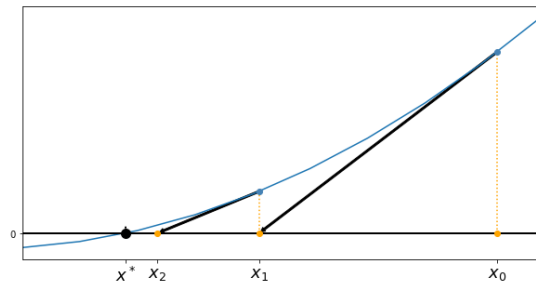




# Rootfinding of equations in one variable



In this chapter, we consider the problem of finding **roots of an equation in one variable**: find  $x$  such that  $f(x) = 0$ . We discuss numerical methods to approximate solutions of this kind of problems to an arbitrarily high accuracy. First, we formalize the notion of convergence and order of convergence for **iterative methods**. Then, we focus on three iterative algorithms approximating roots of functions: **bisection method**, **fixed point iterations** and **Newton Raphson method**. These methods are described, analysed and used to solve 3 problems coming from physics, finance and dynamics of population.

## 1 Table of contents

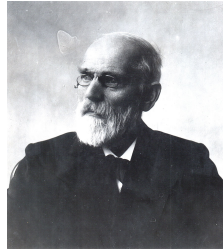
- [Introduction](#)
- [Iterative methods: errors and convergence](#)
- [The bisection method](#)
- [Fixed point iterations](#)
- [The Newton-Raphson method](#)
- [Back to the case studies](#)

```
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 # 2D plotting library
8 import numpy as np             # package for scientific computing
9
10 from math import *              # package for mathematics (pi, arctan, s
```

## 2 Introduction

Computing the zeros of a function  $f$  or equivalently the roots of equation  $f(x) = 0$  is a problem that can be encountered in various situations. In lot of these situations, the solution cannot be computed exactly and one has to design numerical algorithms to approximate the solutions. We give below a few examples of such situations.

## 2.1 Case study 1: State equation of a gaz



**Johannes Diderik van der Waals (1837-1923).** He is a Dutch theoretical physicist. He was primarily known for his thesis work (1873) in which he proposed a state equation for gases to take into account their non-ideality and the existence of intermolecular interactions. His new equation of state revolutionized the study of the behavior of gases. This work was followed by several other researches on molecules that has been fundamental for the development of molecular physics.

The state equation of a gaz relating the pressure  $p$ , the volume  $V$  and the temperature  $T$  proposed by van der Waals can be written

$$\left[ p + a \left( \frac{N}{V} \right)^2 \right] (V - Nb) = kNT$$

where  $N$  is the number of molecules of the gaz,  $k$  is the Boltzmann-constant and  $a$  and  $b$  are coefficients depending on the gaz. To determine the volume occupied by a gaz at pressure  $p$  and temperature  $T$ , we need to solve this equation whose root is  $V$ .

Suppose one wants to find the volume occupied by 1000 molecules of  $\text{CO}_2$  at temperature  $T = 300 \text{ K}$  and pressure  $p = 3.5 \cdot 10^7 \text{ Pa}$ . Then, the previous equation has to be solved for  $V$ , with the following values of parameters  $a$  and  $b$  corresponding to carbon dioxide:  $a = 0.401 \text{ Pa m}^6$  and  $b = 42.7 \cdot 10^{-6} \text{ m}^3$ . The Boltzmann constant is  $k = 1.3806503 \cdot 10^{-23} \text{ J K}^{-1}$ .

## 2.2 Case study 2: Investment found

Suppose someone wants to have a saving account valued at  $S = 30\,000$  euros upon retirement in 10 years. He can deposit  $d = 30$  euros each month on its account. The rate of interest is  $i$  and  $S_n$  the capital after  $n$  months. If the interest is computed monthly, we have:

$$S_n = \sum_{k=0}^{n-1} d(1+i)^k = d \frac{(1+i)^n - 1}{i}$$

If this person wants to know the minimal rate interest needed to achieve his goal, he has to solve the following equation for  $i$ :

$$S = d \frac{(1+i)^{n_{\text{end}}} - 1}{i} \quad \text{where} \quad n_{\text{end}} = 120$$

## 2.3 Case study 3: A first population model



**Thomas Robert Malthus (1766-1834).** He is a British economist. He is mainly known for his works about the links between a population dynamics and its productions. He published anonymously in 1798 an *Essay on the principle of populations*. It is based on the idea that the growth of a population is essentially geometric while the growth of the production is arithmetic. This leads to the so-called Malthusianism doctrine suggesting that the population size has to be controlled to avoid a catastrophe.

Population dynamics is a branch of mathematical biology that gave rise to a great amount of research and is still very active nowadays. The objective is to study the evolution of the size and composition of populations and how the environment drives them. The first model that can be derived is a natural exponential growth model. It depends on two parameters:  $\beta$  and  $\delta$ , the average numbers of births and deaths per individual and unit of time. If we suppose that these parameters are the same for all individuals and do not depend on the size of the population, we can denote the growth rate of the population by  $\lambda = \beta - \delta$  and write:

$$\frac{dN}{dt} = \lambda N$$

where  $N$  is the population size. This model leads to exponentially increasing ( $\lambda > 0$ ) or decreasing populations ( $\lambda < 0$ ). Of course, this model can be enriched to derive more realistic models such as the logistic population growth model where the growth rate  $\lambda$  depends on the size of the population as follows :  $\lambda(N) = \lambda_* - cN$ . This way, too large populations have a negative growth rate, leading to population regulation. When the population is not isolated, one has to take into account immigration or emigration. If we denote by  $r$  the average number of individuals joining the community per unit of time, a new model can be written as

$$\frac{dN}{dt} = \lambda N + r,$$

whose solution is (if  $\lambda \neq 0$ )

$$N(t) = N(0) \exp(\lambda t) + \frac{r}{\lambda}(\exp(\lambda t) - 1).$$

If one wants to estimate the natural growth rate  $\lambda$  in France, one can use the following (evaluated) data:

Population 01/01/2016	Population 01/01/2017	migratory balance in 2016
66 695 000	66 954 000	67 000

and solve the corresponding equation for  $\lambda$  (unit of time = year)

$$N(2017) = N(2016) \exp(\lambda) + \frac{r}{\lambda}(\exp(\lambda) - 1).$$

### 3 Iterative methods: errors and convergence

#### 3.1 Convergence / order of convergence

All the previous problems have the same characteristic: the exact solution cannot be computed through an explicit formula and they have to be approximated through numerical methods.

Let us write these problems under the following generic rootfinding problem:

$$\text{given } f : [a, b] \rightarrow \mathbb{R}, \quad \text{find } x^* \in [a, b] \quad \text{such that } f(x^*) = 0.$$

Methods for approximating the root  $x^*$  of  $f$  are often iterative: algorithms generate sequences  $(x_k)_{k \in \mathbb{N}}$  that are supposed to converge to  $x^*$ . Given such a sequence, the two questions one has to answer are:

- Does the sequence converge to  $x^*$  ?
- if it converges, how fast does it converge to  $x^*$  ?

Before going further, we formalize below the notions of convergence and convergence speed.

**Definition.**

**Convergence.** Suppose that a sequence  $(x_k)_k$  is generated to approximate  $x^*$ . The error at step  $k$  is defined as

$$e_k = |x_k - x^*|$$

where  $|\cdot|$  denotes the absolute value. The sequence  $(x_k)_k$  is said to *converge* to  $x^*$  if

$$e_k \longrightarrow 0 \quad \text{when } k \rightarrow \infty$$

Most of the time, several sequences can be generated and converge to  $x^*$ . One has to choose which one will be used by comparing their properties such as the computational time or the speed of convergence.

**Example.** Let us consider the three following sequences converging to  $x^* = 0$ :

$$x_k = \left(\frac{1}{2}\right)^k, \quad \bar{x}_k = \left(\frac{1}{7}\right)^k, \quad \text{and} \quad \hat{x}_k = \left(\frac{1}{2}\right)^{2^k}$$

The values obtained for the first terms of these sequences are

k	0	1	2	3	4	5
$x_k$	1	0.5	0.25	0.125	0.0625	0.03125
$\bar{x}_k$	1	0.14285	0.02041	0.00291	4.164 e -4	5.94 e -5
$\hat{x}_k$	0.5	0.25	0.0625	0.00390.	1.52 e -5	2.328 e -10

The three sequences converge to zero but  $\hat{x}_k$  seems to converge to zero faster than  $\bar{x}_k$ , itself converging faster than  $x_k$ .

A way to quantify the convergence speed of a sequence is to estimate its order of convergence:

#### Definition.

**Order of convergence for iterative algorithms.** Suppose that the sequence  $(x_k)_k$  converges to  $x^*$ . It is said to converge to  $x^*$  with order  $\alpha > 1$  if

$$\exists k_0 > 0, \quad \exists C > 0, \quad \forall k \geq k_0, \quad \frac{e_{k+1}}{(e_k)^\alpha} \leq C.$$

The convergence is said to be *linear* if  $\alpha = 1$  and *quadratic* if  $\alpha = 2$ .

Of course,

- The bigger is  $\alpha$ , the better is the convergence: the number of exact digits is multiplied by  $\alpha$  at each step.
- $\alpha$  being given, the smaller is  $C$ , the better is the convergence.

**Do it yourself.** Consider again the three following sequences converging to  $x^* = 0$ :

$$x_k = \left(\frac{1}{2}\right)^k, \quad \bar{x}_k = \left(\frac{1}{7}\right)^k, \quad \text{and} \quad \hat{x}_k = \left(\frac{1}{2}\right)^{2^k}$$

Explain the results given in the previous example by studying the order of convergence of the three sequences. Justify your answers.

## 3.2 Graphical study of convergence

### 3.2.1 Study of $e_k$ versus $k$

First, the convergence of a sequence can be observed plotting  $e_k$  versus  $k$ :

```
1 N = np.arange(0,6,1)
2 err1 = (1./2) ** N
3 err2 = (1./7) ** N
4 err3 = (1./2) ** (2**N)
5
6 fig = plt.figure(figsize=(10, 8))
7 plt.plot(N, err1, marker="o", label='error for $x_k$')
8 plt.plot(N, err2, marker="o", label='error for $\bar{x}_k$')
9 plt.plot(N, err3, marker="o", label='error for $\hat{x}_k$')
10 plt.legend(loc='upper right', fontsize=18)
11 plt.xlabel('k', fontsize=18)
12 plt.ylabel('$e_{k}$', fontsize=18)
13 plt.title('Convergence', fontsize=18)
14 plt.show()
```

As expected, the error decrease when  $k$  increase for the three sequences (run the previous cell to see the corresponding plot...).

### 3.2.2 Study of $\log(e_k)$ versus $k$

The  $\log$  function is also of great help to better understand the behaviour of the error.

For example, since  $x \rightarrow \log(x)$  is an increasing function with derivative going to infinity when  $x$  goes to zero, it allows to "zoom" on the smallest values of the error by plotting  $\log(e_k)$  versus  $k$  and check that the error is still decreasing and not stagnating for big values of  $k$  (which can not be affirmed using the previous plot).

**Do it yourself.** Modify the following cell to plot the logarithm of the error versus  $k$ . To do so use the `plt.yscale` function to modify the scale for the y-axis. This function will allow you to plot the logarithm of the error while keeping the values of the error as ticks in the y-axis.

```
1 N = np.arange(0,6,1)
2 x1 = (1./2) ** N
3 x2 = (1./7) ** N
4 x3 = (1./2) ** (2**N)
5
6 fig = plt.figure(figsize=(10, 8))
7 plt.plot(N, x1, marker="o", label='error for $x_k$')
8 plt.plot(N, x2, marker="o", label='error for $\bar{x}_k$')
9 plt.plot(N, x3, marker="o", label='error for $\hat{x}_k$')
10 plt.legend(loc='lower left', fontsize=18)
11 plt.xlabel('k', fontsize=18)
12 plt.ylabel('$e_{k}$', fontsize=18)
13 plt.title('Convergence', fontsize=18)
14 plt.show()
```

Plotting  $\log(e_k)$  versus  $k$  can in fact provide more information on the behaviour of the error than its convergence to zero:

**Do it yourself.** You shall see in the previous figure that two of the plots are linear. What can you deduce about the error from that observation? Check the consistency of this observation with the explicit values of the error for the two corresponding sequences.

### 3.2.3 Study of $\log(e_{k+1})$ versus $\log(e_k)$

Finally, since the  $\log$  is increasing, one have that, for a method of order  $\alpha$ ,

$$\log e_{k+1} \leq \log(Ce_k^\alpha) = \log C + \alpha \log e_k.$$

As a consequence, the convergence rate of a sequence can be a graphically observed by plotting  $\log e_{k+1}$  versus  $\log e_k$ .

**Do it yourself.** Run the following cell and explain the resulting plot. You can add some plots to confirm the slopes of the lines.

```
1 N = np.arange(0,6,1)
2 x1 = (1./2) ** N
3 x2 = (1./7) ** N
4 x3 = (1./2) ** (2**N)
5
6 fig = plt.figure(figsize=(10, 8))
7 plt.loglog(x1[:-1:], x1[1:], marker="o", label='slope ?? for $x_k$') #lo
8 plt.loglog(x2[:-1:], x2[1:], marker="o", label='slope ?? for $\bar{x}_k$')
9 plt.loglog(x3[:-1:], x3[1:], marker="o", label='slope ?? for $\hat{x}_k$')
10 plt.legend(loc='lower right', fontsize=18)
11 plt.axis('equal')
12 plt.xlabel('$e_k$', fontsize=18)
13 plt.ylabel('$e_{k+1}$', fontsize=18)
14 plt.title('Order of convergence', fontsize=18)
15
16 plt.show()
```

**To go further (not graded).** Consider the sequence  $x_{k+1} = 1 - \cos(x_k)$ . Using numerical experiments, conjecture the limit and the order of convergence of this sequence.

### 3.3 Error estimator

To finish, notice that, most of the time, since  $x^*$  is not known, we cannot compute the value of the true error at step  $k$ . Instead we try to find a (calculable) bound for the error, which gives us a "worst-case" error:

**Definition.**

**Error estimator.** Suppose that a sequence  $(x_k)_k$  is generated to approximate  $x^*$ . The sequence  $(\beta_k)_k$  is an error estimator if

- $\beta_k > 0$  is computable
- $\beta_k$  is a bound for the error:  $e_k < \beta_k$  for all  $k$

In that case, if the estimator  $\beta_k \rightarrow 0$  when  $k \rightarrow \infty$ , we obtain that

- the sequence  $x_k$  converges to  $x^*$
- the error goes to zero at least as fast as the sequence  $\beta_k$ .

One has to take care that an estimator only provides an upper bound on the error. As a consequence, the error can go to zero faster than the estimator.

## 4 The bisection method

The first method to approximate the solution to  $f(x) = 0$  is based on the Intermediate Value Theorem (see Appendix). Suppose  $f$  is a continuous function on the interval  $[a, b]$  where  $f(a)$  and  $f(b)$  have opposite signs:  $f(a)f(b) < 0$ . Then, there exists  $x^*$  in  $]a, b[$  such that  $f(x^*) = 0$ .

Starting from an interval  $I_0 = [a_0, b_0]$  such that  $f(a_0)f(b_0) < 0$ . Let  $x_0$  be the midpoint of  $I_0$ :

$$x_0 = \frac{a_0 + b_0}{2}.$$

Then, the bisection method iterates by choosing  $I_1 = [a_1, b_1]$  and  $x_1$  as follows:

- if  $f(x_0) = 0$  then  $x^* = x_0$  and the algorithm terminates
- if  $f(a_0)f(x_0) < 0$  then there exists a zero of  $f$  in  $[a_0, x_0]$ : set

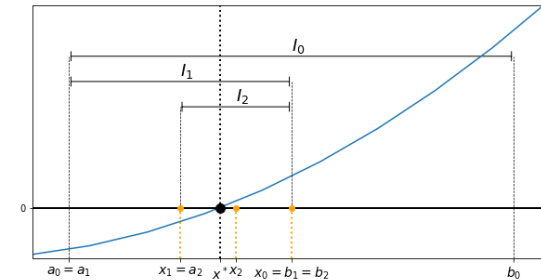
$$a_1 = a_0, \quad b_1 = x_0 \quad \text{and} \quad x_1 = \frac{a_1 + b_1}{2}$$

- if  $f(x_0)f(b_0) < 0$  then there exists a zero of  $f$  in  $[x_0, b_0]$ : set

$$a_1 = x_0, \quad b_1 = b_0 \quad \text{and} \quad x_1 = \frac{a_1 + b_1}{2}$$

The method iterates until a stopping criterion that will be discussed later.

An example of the first two iterations is illustrated on an example in the figure below.



The bisection method leads to the following algorithm:

### Algorithm.

**Bisection method.** Computes a sequence  $(x_k)_k$ , approximating  $x^*$  solution to  $f(x^*) = 0$ .

```
INPUT :  $f, a, b$ 
DO :  $x = (a + b)/2$ 
    While stopping criterion is not achieved do
        If  $f(a)f(x) < 0$ ,  $b = x$  else  $a = x$ 
         $x = (a + b)/2$ 
    end while
RETURN :  $x$ 
```

In the following, we implement the bisection method and test it to approximate  $x^*$ , the unique solution in  $\mathbb{R}$  to  $f(x) = x^3 - 2 = 0$ . In this first version, the stopping criterion is: stop if the requested number of iteration is achieved or if the zero was found.

**Do it yourself.** Complete the following function encoding  $f$ .

```
1 ## Function f:  $x \rightarrow x^3 - 2$ 
2
3 def ftest(x):
4     return ---
```

**Do it yourself.** Complete the following function. It shall compute the sequence generated using the bisection algorithm for a given function  $f$  and initialized by an interval  $[a_0, b_0]$ .

The algorithm terminates when the zero is found or when a given maximal number  $K$  of iterations has been achieved. The output is a vector  $x$  with size  $K + 1$ . It contains the values of the sequence:  $x[k] = x_k$ .

```
1 ## Bisection algorithm for function f
2 ## input : f = name of the function
3 ##         a0, b0 = initial interval I_0 with f(a)f(b)<0
4 ##         K = number of iterations
5 ## output : x = sequence approximating the zero of f
6 ##         x[k]=x_k for k=0..K
7
8 def Bisection(f,a0,b0,K):
9     x = --- # create vector x of zeros with size K+1
10    k = 0 # initialize k
11    a = a0 # initialize a
12    b = b0 # initialize b
13    x[0] = --- # initialize x_0
14    # computation of x_k for k>0
15    # stops if f(x[k])=0 or if the number of iterations is achieved
16    while --- and --- : #test the two stopping criterion
17        # do not stop => enter the loop and iterate the bisection algorithm
18        if --- :
19            --- #do something
20        else:
21            --- #do something
22        k = k+1
23        x[k] = --- #compute and store x_k
24    return x
```

**Do it yourself.** Test the bisection method to compute  $x^* = 2^{1/3}$  solution to  $f(x) = 0$ . Initialize with  $[a_0, b_0] = [1, 2]$  and compute the first 20 iterations. Plot the error  $e_k$  versus  $k$ . Use a log scale for the error (y-axis). Do not forget to add a title to the figure and a label to the axes (see the graphical study in the previous section as example).

```
1 xstar = 2**(1.0/3)
2
3 # parameters
4 a0 = ---
5 b0 = ---
6 K = ---
7
8 # compute the first 20 iterations of the bisection method for I0=[1,2]
9 x=---
10
11 #print x^* and x
12 print('xstar =',xstar)
13 print('x =',x)
14
15 # compute the error
16 # err is a vector, err[k]=abs(x[k]-x^*)
17 err = ---
18
19 # create the vector tabk : tabk[k]=k for k=0..K
20 tabk = ---
21
22 # plot the error versus k
23 fig = plt.figure(figsize=(10, 5))
24 plt.plot(---, ---, marker="o")
25 # set log scale for the error (y-axis)
26 ---
27 # set title of the figure and labels of the axis
28 ---
29
30 plt.show()
```

**Do it yourself.** Comment the previous plot.

## 4.1 Error estimator and stopping criterion

In the previous example, the stopping criterion is simply based on the number of iterations the user wants to achieve. However, when one wants to approximate  $x^*$ , one has in mind the maximal error allowed and therefore, fixing the number of iterations has no sense as a stopping criterion. A criterion based on the error at the current step would be much more meaningful.

Suppose that a parameter  $\epsilon$  is given, fixing the precision needed. We give below three classical stopping criteria:

$$\begin{aligned} & 1. \quad |x_k - x_{k-1}| < \epsilon \\ & 2. \quad |f(x_k)| < \epsilon \\ & 3. \quad \frac{|x_k - x_{k-1}|}{|x_k|} < \epsilon \end{aligned}$$

Unfortunately, each of these criteria can induce difficulties. For example, criterion 1 can be fulfilled even for non-converging sequences (think e.g. at  $x_k = \sum_{j=1}^k \frac{1}{j}$ ). Criterion 2 is also non-relevant for some functions  $f$  for which  $f(x)$  can be close to zero while  $x$  is still far from  $x^*$ : the test will be satisfactory if  $f'(x^*) \approx 1$ , not reliable if  $f'(x^*) \ll 1$  and too restrictive if  $f'(x^*) \gg 1$ .

Without any further information on  $f$  or on the convergence of the sequence, one should make criterion 3 its first choice.

In order to use a more precise stopping criterion, related to the true error, one should know more about the way the sequence converges to  $x^*$ . To do so, error estimators are very useful. Concerning the bisection method we have the following result:

### Proposition.

**Convergence of the bisection method.** Let  $f$  be a continuous function on  $[a, b]$  with  $f(a)f(b) < 0$ . Suppose  $(x_k)_k$  is the sequence generated by the bisection method to approximate  $x^*$ , solution to  $f(x) = 0$  on  $[a, b]$ .

Then, the sequence  $(x_k)_k$  converges to  $x^*$  and the following estimation holds:

$$\forall k \geq 0, \quad |x_k - x^*| \leq \frac{b-a}{2^k}.$$

**Proof.** Since the interval is divided by 2 at each step of the method, we have

$$\forall k \geq 0 \quad |b_k - a_k| \leq \frac{b-a}{2^k}$$

Remarking that both  $x^*$  and  $x_k$  are in  $I_k = [a_k, b_k]$ , we obtain

$$\forall k \geq 0 \quad |x_k - x^*| \leq \frac{b-a}{2^k}$$

This proves the convergence of  $x_k$  to  $x^*$  and provides the requested estimation.

**Remark.** The bisection method is said to be *globally convergent*. Indeed, the initialization of  $a$  and  $b$  doesn't need to be close to  $x^*$ . Whatever the choice for these parameters is, the generated sequence will converge to  $x^*$ , provided that  $f(a)f(b) < 0$ .

This proposition provides a new stopping criterion: if one wants the error to be less than  $\epsilon$ , one should stop at iteration  $k$  such that

$$\frac{b-a}{2^k} \leq \epsilon.$$

We rewrite the code for the bisection method using this criterion. Note that we still ask for a maximal number of iterations in order to avoid infinite loops in case the convergence of the method is too slow to lead to the requested precision in a reasonable time.

**Do it yourself.** Rewrite the bisection algorithm so that it terminates when the stopping criterion  $\frac{b-a}{2^k} \leq \epsilon$  is verified or when a maximal number  $K_{max}$  of iterations have been achieved. If  $k_{end}$  is the number of iterations needed to fulfil this criteria, we have  $k_{end} \leq K$ .

The function returns a tuple of two elements: the vector  $x$  containing the computed iterations together with  $k_{end}$ , the number of iterations achieved.

$x$  is initialized as vector with size  $K_{max} + 1$ . At the end of the algorithm, it contains the computed values of the sequence:  $x[k] = x_k$  for  $0 \leq k \leq k_{end} + 1$  and the other elements of  $x$  are equal to 0. Take care to remove the zeros and return a vector with size  $k_{end} + 1$ .

```

1  ## Bisection algorithm for function f
2  ## input : f = name of the function
3  ##      a0, b0 = initial intervall I_0 with f(a0)f(b0)<0
4  ##      eps = tolerance
5  ##      Kmax = maximal number of iterations allowed
6  ## output : x = sequence approximating the zero of f
7  ##      k = total number of iterations that has been achieved (lower
8
9  def Bisection2(f, a0, b0, eps, Kmax):
10     ---
11     return (x, k)

```

**Do it yourself.** Test this new function to compute  $2^{1/3}$  with precision at least  $\epsilon = 10^{-3}$ . Use  $K_{\max} = 20$ ,  $I_0 = [1, 2]$ . Plot on the same figure the error versus  $k$  and the corresponding estimator. Do not forget the title, the labels of axes and the legend. Take care that the output  $x$  of the bisection function is of size  $K_{\max}$  while the number of iterations  $k_{\text{end}}$  can be strictly smaller than  $K_{\max}$ .

```

1  # parameters
2  a0 = ---
3  b0 = ---
4  eps = ---
5  Kmax = ---
6
7  xstar = 2**(1.0/3)
8
9  # run the bisection method
10 x, kend = ---
11 print('precision: eps =', eps)
12 print('number of iterations =', kend)
13
14 # compute the error
15 # err is a vector, err[k]=abs(x[k]-x^*) for k=0..kend
16 err = ---
17
18 # create the vector tabk : tabk[k]=k for k=0..kend
19 tabk = ---
20
21 # compute the error estimator, errEstim[k]=(b-a)/2^k for k=0..kend
22 # use tabk / no loop on k
23 errEstim = ---
24
25 # plot the error versus k
26 fig = plt.figure(figsize=(10, 5))
27 plt.plot(---, ---, marker="o", label="Error")
28 # plot the error estimator versus k
29 plt.plot(---, ---, marker="o", label="Error estimator")
30 # set log scale for the error (y-axis)
31 ---
32 # set title of the figure, labels of the axis and the legend
33 ---
34
35
36 plt.show()

```

**Do it yourself.** Comment the previous plot.



**Luitzen Egbertus Jan Brouwer (1881 – 1966) and Stefan Banach (1892-1945).** Brouwer is a Dutch mathematician and philosopher. He proved a lot of results in topology. One of his main theorem is his fixed point theorem (1909). One of its simpler form says that a continuous function from an interval to itself has a fixed point. The proof of the theorem does not provide a method to compute the corresponding fixed point. Among lot of other fixed point results, Brouwer's theorem became very famous because of its use in various fields of mathematics or in economics. In 1922, a polish mathematician, Stefan Banach, stated a contraction mapping theorem, proving in some case the existence of a unique fixed point and providing a constructive iterative method to approximate this fixed point. Banach is one of the founders of modern analysis and is often considered as one of the most important mathematicians of the 20-th century.

## 5 Fixed point iterations



A fixed point for a function  $g$  is a number  $x$  such that  $g(x) = x$ . In this section we consider the problem of finding solutions of fixed point problems. This kind of problem is equivalent to rootfinding problems in the following sense:

- If  $x^*$  is a solution to  $f(x) = 0$ , we can find a function  $g$  such that  $x^*$  is a fixed point of  $g$ . For example, one can choose  $g(x) = f(x) + x$ .
- If  $x^*$  is a solution to  $g(x) = x$ , then,  $x^*$  is also a solution to  $f(x) = 0$  where  $f(x) = g(x) - x$ .

If the two kind of problems are equivalent, the fixed point problem is easier to analyze. In this section, we will focus on such problems in order to understand how to use them the best way for solving rootfinding problems. In the following, functions  $f$  will be used for rootfinding problems and  $g$  for corresponding fixed point problems.

First, note that, given a function  $f$ , the choice of  $g$  is not unique. For example, any function  $g$  of the form  $g(x) = G(f(x)) + x$  where  $G(0) = 0$  is suitable for solving the problem. Let us consider again the problem of computing an approximation of  $x^* = 2^{1/3}$  as the root of  $f(x) = x^3 - 2$ . The five following functions  $g$  can be chosen:

- $g_1(x) = x^3 - 2 + x$
- $g_2(x) = \sqrt{\frac{x^5 + x^3 - 2}{2}}$
- $g_3(x) = -\frac{1}{3}(x^3 - 2) + x$
- $g_4(x) = -\frac{1}{20}(x^3 - 2) + x$
- $g_5(x) = \frac{2}{3}x + \frac{2}{3x^2}$

From a numerical point a view, solutions to fixed point problems can be approximated by choosing an initial guess  $x_0$  for  $x^*$  and generate a sequence by iterating function  $g$ :

$$x_{k+1} = g(x_k), \quad \text{for } k \geq 0.$$

Indeed, suppose that  $g$  is continuous and that the sequence  $(x_k)_k$  converges to  $x_\infty$ , then, passing to the limit in the previous equation gives

$$x_\infty = g(x_\infty)$$

and  $x_\infty$  is a fixed point of  $g$ . This leads to the following algorithm:

#### Algorithm.

**Fixed point iterations method.** Computes a sequence  $(x_k)_k$ , approximating  $x^*$  solution to  $g(x^*) = x^*$ .

```

INPUT :  g, x0
DO :    x = x0
        While stopping criterion is not achieved do
            x = g(x)
        end while
RETURN :  x

```

Now, for a given function  $g$ , one has to answer the following questions:

- does  $g$  have a fixed point ?
- does the sequence generated using fixed point iterations converge ?
- if the sequence converges, how fast does it converge ?

## 5.1 Graphical investigation

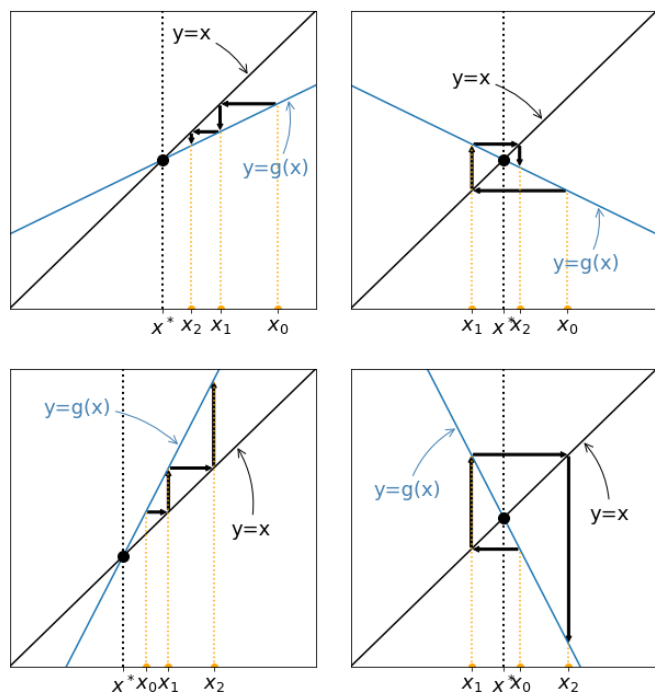
In order to better understand the behaviour of fixed point iterations, one can try to visualize them on a graph.

First, the fixed point of a function  $g$  can be found graphically searching for the intersection between the graph of  $g$  and the graph of function  $\phi(x) = x$ .

Then, suppose  $x_0$  is given and place it on the abscissa axis. To place  $x_1 = g(x_0)$  on the same axis, proceed as follows:

- from  $(x_0, 0)$ , go up to find the point  $(x_0, g(x_0)) = (x_0, x_1)$ , when crossing the graph of  $g$
- from  $(x_0, x_1)$  move horizontally to find the point  $(x_1, x_1)$ , when crossing the graph of  $\phi$
- finally, go down towards the abscissa axis to place the point  $(x_1, 0)$

Then iterate the procedure to visualize the generated sequence. Four examples are given below:



Cases with increasing functions  $g$  are given on the left and leads to monotonous sequences. On the contrary, oscillating sequences are generated for non increasing functions  $g$  (right). The two examples given at the top converge. Remark that they correspond to cases where  $|g'(x)| < 1$ .

## 5.2 Convergence analysis

### Theorem.

**Existence of a fixed point.** Let  $g : [a, b] \rightarrow \mathbb{R}$ . Suppose

- $g \in C[a, b]$
- $g : [a, b] \rightarrow [a, b]$  (i.e.  $[a, b]$  is stable for  $g$ )

Then,  $g$  has a fixed point in  $[a, b]$ :

$$\exists x^* \in [a, b], \quad g(x^*) = x^*$$

**Do it yourself.** Complete the proof of the previous theorem.

**Proof.**

### Theorem.

**Existence of a unique fixed point.** Let  $g : [a, b] \rightarrow \mathbb{R}$ . Suppose

- $g \in C[a, b]$
- $g : [a, b] \rightarrow [a, b]$  (i.e.  $[a, b]$  is stable for  $g$ )
- $g'$  exists on  $[a, b]$  and  
 $\exists K < 1$  such that  $\forall x \in [a, b], \quad |g'(x)| \leq K$  (i.e.  $g$  is a contraction)

Then,  $g$  has a unique fixed point in  $[a, b]$ :

$$\exists! x^* \in [a, b], \quad g(x^*) = x^*$$

**Proof.** The existence of a fixed point  $x^*$  is given by the previous theorem. The fact that  $g$  is a contraction mapping ensures the uniqueness of the fixed point. Indeed, suppose that  $x^1$  and  $x^2$  are two fixed points of  $g$  and write the Taylor Lagrange expansion of  $g$  around  $x^1$  at order 1:

$$\exists \xi \in I_{x^1, x^2}, \quad \text{such that} \quad g(x^2) = g(x^1) + (x^2 - x^1) g'(\xi)$$

where  $I_{x^1, x^2} = [x^1, x^2]$  if  $x^1 < x^2$  and  $I_{x^1, x^2} = [x^2, x^1]$  otherwise.

Using  $g(x^1) = x^1$  and  $g(x^2) = x^2$ , we obtain

$$x^2 - x^1 = (x^2 - x^1) g'(\xi)$$

and using the contraction:

$$|x^2 - x^1| \leq K |x^2 - x^1|$$

which gives  $x^2 = x^1$  since  $K < 1$ .

**Theorem.**

**Convergence of fixed point iterations.** Let  $g : [a, b] \rightarrow \mathbb{R}$ . Consider the sequence  $x_{k+1} = g(x_k)$  for  $k \geq 0$ ,  $x_0$  being given. Suppose

- $g \in C([a, b])$
- $g : [a, b] \rightarrow [a, b]$  (i.e.  $[a, b]$  is stable for  $g$ )
- $g' \in C^1([a, b])$  and
 
$$\exists K < 1 \quad \text{such that} \quad \forall x \in [a, b], \quad |g'(x)| \leq K$$

Then,  $g$  has a unique fixed point  $x^*$  in  $[a, b]$  and the sequence  $(x_k)_k$  converges to  $x^*$  for any choice of  $x_0 \in [a, b]$ . Moreover we have

$$\lim_{k \rightarrow \infty} \frac{x_{k+1} - x^*}{x_k - x^*} = g'(x^*)$$

so that the sequence converges at least with order 1.

**Proof.** The existence and uniqueness of the fixed point is given by the previous theorem. The convergence analysis is given again using a Taylor expansion:

$$\forall k \geq 0, \quad \exists \xi_k \in I_{x^*, x_k}, \quad \text{such that} \quad g(x_k) = g(x^*) + (x_k - x^*) g'(\xi_k).$$

This, together with  $g(x_k) = x_{k+1}$  and  $g(x^*) = x^*$  gives

$$\forall k \geq 0, \quad \exists \xi_k \in I_{x^*, x_k}, \quad \text{such that} \quad x_{k+1} - x^* = (x_k - x^*) g'(\xi_k)$$

From this we obtain that

$$|x_{k+1} - x^*| \leq K |x_k - x^*| \leq K^{k+1} |x_0 - x^*| \rightarrow 0 \quad \text{when } k \rightarrow \infty$$

and the sequence converges to  $x^*$ . Moreover, since  $x_k$  converges to  $x^*$ , we have that  $\xi_k$  converges to  $x^*$  and from the continuity of  $g'$  we obtain  $g'(\xi_k) \rightarrow g'(x^*)$  when  $k$  goes to infinity. Then, we have

$$\frac{x_{k+1} - x^*}{x_k - x^*} = g'(\xi_k) \rightarrow g'(x^*) \quad \text{when } k \rightarrow \infty$$

which ends the proof.

**Remark.** Note that these theorems provide sufficient but not necessary condition for convergence.

- If  $|g'(x^*)| > 1$ , if  $x_k$  is sufficiently close to  $x^*$  we have that  $g'(\xi_k) > 1$  and then  $|x_{k+1} - x^*| > |x_k - x^*|$ . The sequence cannot converge.

**Do it yourself.** Complete the following function. It shall compute the sequence generated using the fixed point algorithm for a given function  $g$ . The algorithm terminates when a given number  $K$  of iterations have been achieved.

```
1  ## Fixed point algorithm for function g
2  ## input : g = name of the function
3  ##         x0 = initialization
4  ##         K = number of iterations
5  ## output : x = sequence generated using the fixed point iteration for g
6
7  def FixedPoint(g, x0, K):
8      # create vector x
9      x = np.zeros(K+1)
10     k = 0
11     x[0] = x0
12     # computation of x_k
13     while ---:
14         ---
15         k=k+1
16     return x
```

**Do it yourself.** Run the two following cells to test the fixed point algorithm for the functions:

- $\phi_1(x) = x - x^3$
- $\phi_2(x) = x + x^3$

What can you conclude for the case  $|g'(x^*)| = 1$  ?

```
1  # phi1(x) = x-x^3.
2
3  def phi1(x):
4      return x - x**3
5
6  x0 = 0.1
7  K = 20
8  x=FixedPoint(phi1,x0,K)
9  print('x =',x)
```

```
1  # phi2(x) = x+x^3.
2
3  def phi2(x):
4      return x + x**3
5
6  x0 = 0.1
7  K = 20
8  x=FixedPoint(phi2,x0,K)
9  print('x =',x)
```

**Remark.** The fixed point theorem ensures the convergence of the sequence for any choice of  $x_0 \in [a, b]$  and then presents a global convergence result.

However, in practice, even if  $|g'(x^*)| < 1$ , finding a stable interval on which  $g$  is a contracting mapping is not so easy.

In fact, one can prove that, if  $g$  is continuous and differentiable and if  $|g'(x^*)| < 1$ , such an interval exists: more precisely, there exists a neighbourhood  $I$  of  $x^*$  such that, for any  $x_0 \in I$ , the fixed point iterations converge to  $x^*$ . This local convergence result is stated in the following theorem:

#### Theorem.

**Local convergence for fixed point iterations.** Let  $g : [a, b] \rightarrow \mathbb{R}$ . Consider the sequence  $x_{k+1} = g(x_k)$  for  $k \geq 0$ ,  $x_0$  being given. Suppose

- $x^*$  is a fixed point of  $g$
- $g \in C([a, b])$
- $g$  is differentiable on  $[a, b]$  and  $|g'(x^*)| < 1$

Then, there exists a neighbourhood  $I$  of  $x^*$  such that, for any  $x_0 \in I$ , the fixed point iterations converge to  $x^*$ .

From the previous estimations, we remark that the smaller is the constant  $|g'(x^*)|$ , the faster is the convergence. In the next theorem, we prove (among others) that for  $|g'(x^*)| = 0$ , the convergence is quadratic.

#### Theorem.

**"Better than linear" speed of convergence of fixed point iterations.** Let  $g : [a, b] \rightarrow \mathbb{R}$  and suppose that the hypothesis of the previous theorem are fulfilled. If

- $g \in C^{p+1}(I)$  where  $I$  is a neighbourhood of  $x^*$  and  $p$  is an integer  $p \geq 0$
- $g^{(i)}(x^*) = 0$  for  $0 \leq i \leq p$
- $g^{(p+1)}(x^*) \neq 0$

Then, the fixed point iteration method with function  $g$  has order  $p + 1$  and

$$\lim_{k \rightarrow \infty} \frac{x_{k+1} - x^*}{(x_k - x^*)^{p+1}} = \frac{g^{(p+1)}(x^*)}{(p+1)!}.$$

This proves that the sequence converges at least with order  $p + 1$ .

**Proof.** Again, we expand  $g$  around  $x^*$  at order  $p + 1$ :

$$\forall k \geq 0, \quad \exists \xi_k \in I_{x^*, x^k}, \quad \text{such that} \quad g(x_k) = g(x^*) + \frac{(x_k - x^*)^{p+1}}{(p+1)!} g^{(p+1)}(\xi_k)$$

and we obtain

$$\frac{x_{k+1} - x^*}{(x_k - x^*)^{p+1}} = \frac{g^{(p+1)}(\xi_k)}{(p+1)!} \rightarrow \frac{g^{(p+1)}(x^*)}{(p+1)!} \text{ when } k \rightarrow \infty$$

## 5.3 Numerical tests

**Do it yourself.** We consider again the 5 iteration functions proposed at the beginning of the section to compute  $x^* = 2^{1/3}$ . Run the following cells to observe the behaviour of the algorithm for these 5 functions and comment in light of the previous theorems.

```
1 xstar = 2**(1.0/3)
```

- $g_1(x) = x^3 - 2 + x$

```
1 def g1(x):
2     return x**3 - 2 + x
3
4 x0 = xstar + 0.001
5 #x0 = xstar - 0.001
6 K = 10
7 x = FixedPoint(g1, x0, K)
8 print('x =', x)
```

- $g_2(x) = \sqrt{\frac{x^5 + x^3 - 2}{2}}$

```
1 def g2(x):
2     return np.sqrt((x**5 + x**3 - 2) / 2)
3
4 x0 = xstar - 0.001
5 #x0 = xstar + 0.001
6 K = 10
7 x = FixedPoint(g2, x0, K)
8 print('x =', x)
```

- $g_3(x) = -\frac{1}{3}(x^3 - 2) + x$

```

1 def g3(x):
2     return - (x**3-2)/3 + x
3
4 x0 = xstar + 1
5 #x0 = xstar + 2
6 K = 10
7 x = FixedPoint(g3,x0,K)
8 print('xstar =',xstar)
9 print('x =',x)
10 err3 = abs(x-xstar)
11 print('error =',err3)

```

$$\bullet g_4(x) = -\frac{1}{20}(x^3 - 2) + x$$

```

1 def g4(x):
2     return - (x**3-2)/20 + x
3
4 x0 = xstar + 1
5 #x0 = sqrt(2) + 4
6 K = 10
7 x = FixedPoint(g4,x0,K)
8 print('xstar =',xstar)
9 print('x =',x)
10 err4 = abs(x-xstar)
11 print('error =',err4)

```

$$\bullet g_5(x) = \frac{2}{3}x + \frac{2}{3x^2}$$

```

1 def g5(x):
2     return 2*x/3 + 2/(3*x**2)
3
4 x0 = xstar + 1
5 K = 5
6 x = FixedPoint(g5,x0,K)
7 print('xstar =',xstar)
8 print('x =',x)
9 err5 = abs(x-xstar)
10 print('error =',err5)

```

**Do it yourself.** Compare graphically the convergence for iterations of  $g_3$ ,  $g_4$ , and  $g_5$ :

- On the same figure, plot the three errors versus  $k$  with log-scale for the error.
- On the same figure, plot the  $e_{k+1}$  versus  $e_k$  in log-log scale for the three methods.

Do not forget titles, labels and legends.

```

1 # initialization
2 x0 = xstar + 0.2
3
4 # g3 and g4: compute 10 iterations
5 K=10
6 tabk1 = np.arange(0,K+1,1)
7 x3 = ---
8 err3 = ---
9 x4 = ---
10 err4 = ---
11
12 # g5: compute 3 iterations (if K is too big, the error reaches 0 and log-
13 K=3
14 tabk2 = np.arange(0,K+1,1)
15 x5 = ---
16 err5 = ---
17
18 fig = plt.figure(figsize=(20, 10))
19
20 plt.subplot(121) # plot of e_k versus k for the three methods
21 plt.plot(---)
22 plt.plot(---)
23 plt.plot(---)
24
25 plt.subplot(122) # plot of e_{k+1} versus e_k for the three methods
26 plt.loglog(---) #log-log scale
27 plt.loglog(---) #log-log scale
28 plt.loglog(---) #log-log scale
29 plt.axis('equal')
30
31 plt.show()

```

**Do it yourself.** Comment the previous figures.

## 5.4 Stopping criterion

In general, fixed point iterations are terminated using criterion 1: for  $\epsilon$  given, the computation terminates when

$$|x_{k+1} - x_k| < \epsilon$$

This is justified by the fact that, using again a Taylor expansion, we have:

$$\exists \xi_k \in I_{x^*, x_k}, \quad \text{such that} \quad g(x_k) = g(x^*) + (x_k - x^*) g'(\xi_k)$$

From this, together with  $g(x_k) = x_{k+1}$  and  $g(x^*) = x^*$  we get

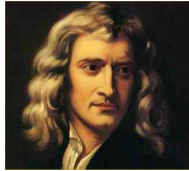
$$x^* - x_k = (x^* - x_{k+1}) + (x_{k+1} - x_k) = -(x_k - x^*) g'(\xi_k) + (x_{k+1} - x_k)$$

and finally we obtain:

$$x^* - x_k = \frac{1}{1 - g'(\xi_k)} (x_{k+1} - x_k)$$

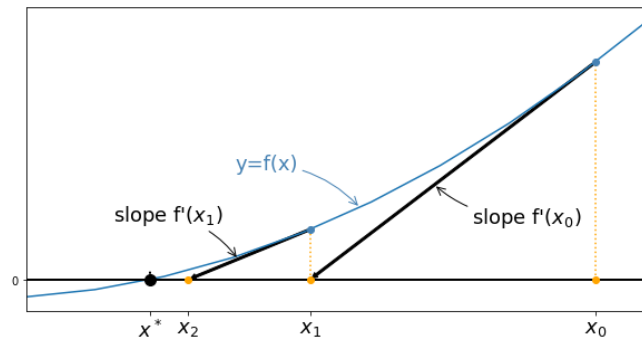
Consequently, if  $g'(x^*) = 0$  (which is the case for methods of order 2),  $x_{k+1} - x_k$  is a good estimator for the error. In the case  $g'(x^*)$  is close to 1, it is not satisfactory...

## 6 The Newton-Raphson method



**Isaac Newton (1643 – 1727).** English mathematician, astronomer, theologian, author and physicist, Isaac Newton is known as one of the most important scientists. He made breaking contributions to classical mechanics, optic and also contributed to infinitesimal calculus. In particular, he described in an unpublished work in 1671 a method to find zeros of polynomials now known as the Newton-Raphson method. Indeed, it was first published (with a reference to Newton) by another English mathematician, Joseph Raphson in 1690. Newton finally published his analysis in 1736. Both of them focused on zeros of polynomial functions but the basis of the general method was already present in their works.

The Newton-Raphson (or simply Newton's) method is one of the most powerful and well-known method to solve rootfinding problems  $f(x) = 0$ . The simplest way to describe it is to see it as a graphical procedure:  $x_{k+1}$  is computed as the intersection with the  $x$ -axis of the tangent line to the graph of  $f$  at point  $(x_k, f(x_k))$ .



So that the Newton's method starts with an initial approximation  $x_0$  and generates the sequence of approximations  $(x_k)_k$  defined by

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)},$$

which leads to the following algorithm:

### Algorithm.

**Newton-Raphson method.** Computes a sequence  $(x_k)_k$ , approximating  $x^*$  solution to  $f(x^*) = 0$ .

```

INPUT :  $f, x_0$ 
DO :  $x = x_0$ 
    While stopping criterion is not achieved do
         $x = x - \frac{f(x)}{f'(x)}$ 
    end while
RETURN :  $x$ 

```

Interpreting Newton's method as a fixed point iteration method, one can prove the following **local** convergence theorem:

### Theorem.

**Local convergence of Newton's method.** Let  $f : [a, b] \rightarrow \mathbb{R}$ . Consider the sequence  $(x_k)_k$  generated by Newton's method for  $k \geq 0$ ,  $x_0$  being given. Suppose

- $x^*$  is a root of  $f$  in  $[a, b]$
- $f \in C^2([a, b])$
- $f'(x^*) \neq 0$  ( $x^*$  is a simple root of  $f$ )

Then, there exists a neighbourhood  $I$  of  $x^*$  such that, for any  $x_0 \in I$ , Newton's iterations converge to  $x^*$  and the convergence is of order 2.

**Proof.** Let us consider function  $g(x) = x - \frac{f(x)}{f'(x)}$ , such that  $x_{k+1} = g(x_k)$ . Using continuity of  $f'$ ,  $g$  is defined in a neighbourhood  $I$  of  $x^*$ . Moreover  $g \in \mathcal{C}(I)$  and we have

$$g'(x^*) = \frac{f(x^*)f''(x^*)}{(f'(x^*))^2} = 0$$

so that the fixed point local convergence theorem provides a neighbourhood  $\bar{I} \subset I$  of  $x^*$  for which the sequence converges towards  $x^*$  if  $x_0 \in \bar{I}$ .

If we suppose that  $f \in \mathcal{C}^3(I)$ , one can prove the quadratic convergence using the corresponding result of the fixed point iterations of  $g \in \mathcal{C}^2(I)$ . In fact, the result is still true for  $f \in \mathcal{C}^2(I)$ . Indeed, a Taylor expansion of  $f$  gives

$$0 = f(x^*) = f(x_k) + f'(x_k)(x^* - x_k) + \frac{f''(\xi_k)}{2}(x^* - x_k)^2 \quad \text{with } \xi_k \in I_{x^*, x_k}$$

and then using that  $\xi_k \rightarrow x^*$  we have

$$\frac{x_{k+1} - x^*}{(x_k - x^*)^2} = \frac{f''(\xi_k)}{2f'(x_k)} \rightarrow \frac{f''(x^*)}{2f'(x^*)} \text{ when } k \rightarrow \infty$$

which proves the quadratic convergence.

**Do it yourself.** In view of this result, comment the order of convergence you observed for the fixed point method applied to  $g_5$ .

**Remark.** One of the main drawback of Newton's method is that the convergence result is a local convergence result. As a consequence, the sequence has to be carefully initialized with an approximation  $x_0$  close to  $x^*$ , which is not so easy to do in practice. A method to do that is to run a bisection method to compute a rough approximation of  $x^*$  and then to initialize Newton methods with this approximation in order to make it much more precise.

**Remark.** Another drawback of Newton's method is that it necessitates the evaluation of the derivative of  $f$  at each iteration. Most of the time,  $f'$  is much more difficult to evaluate than  $f$  and it can even be unknown... To skip this difficulty, the derivative can be approximated by

$$f'(x_k) \approx \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}.$$

The corresponding algorithm is called the **secant method**.

**Remark.** Another main difficulty with Newton's method is the case where  $f'(x^*)$  is close to (or equal to) zero. Suppose that it is the case but that the sequence is still defined for any  $x \geq 0$  (i.e.  $f'(x_k) \neq 0$  for all  $k \geq 0$ ). Then

- if  $f'(x^*) < 1$  but  $f'(x^*) \neq 0$ . The convergence is still quadratic but is very deteriorated due to the big constant  $\frac{f''(x^*)}{2f'(x^*)}$

- if  $f'(x^*) = 0$ ,  $x^*$  is a multiple root and we do not have anymore  $g'(x^*) = 0$ .

One can prove that  $g'(x^*) = 1 - \frac{1}{m}$  where  $m$  is the multiplicity of the root.

From  $|g'(x^*)| < 1$  we obtain the local convergence of the algorithm with order 1. The quadratic convergence can be recovered using fixed point iterations with  $g^{new}(x) = x - m \frac{f(x)}{f'(x)}$ .

We are now going to use Newton's method to solve case study 2 and 3. To do so, we first implement Newton method and test it to approximate  $x^* = 2^{1/3}$ , the unique solution in  $\mathbb{R}$  to  $f(x) = x^3 - 2 = 0$ .

**Do it yourself.** Implement Newton method and test it to approximate  $x^* = 2^{1/3}$ , the unique solution in  $\mathbb{R}$  to  $f(x) = x^3 - 2 = 0$ . Check that you recover the results obtained using the fixed point iteration with function  $g_5$ . In this version, the stopping criterion is: stop if the maximal number of iteration is achieved, if the zero was found or if  $|x_{k+1} - x_k| < \epsilon$  with  $\epsilon$  given.

```
1  ## Newton's algorithm for function f
2  ## input : f = name of the function
3  ##         df = name of the derivative of function f
4  ##         x0 = initial guess for x**
5  ##         eps = precision for stopping criterion
6  ##         Kmax = maximal number of iterations
7  ## output : x = sequence approximating the zero of f
8
9  def Newton(f, df, x0, eps, Kmax):
10     ---
11     return (x, k)
```

```
1  ## Test of the newton algorithm for f(x) = x^3 - 2
2  ## comparison with the results given by the fixed point iterations for f
3
4  x0 = xstar + 1
5  ---
6
```

**Do it yourself.** Check on the previous example that  $|x^{k+1} - x^k|$  is a good estimator for the error  $|x^* - x^k|$  (case of a fixed point of order 2). To do so, plot the two quantities versus  $k$  on the same figure and comment. Explain why, when the algorithm stops, the precision is much better than expected.

```

1 err = ---
2 criterion = ---
3
4 fig = plt.figure(figsize=(10, 5))
5 ---
6
7 plt.show()

```

## 7 Back to the case studies

We come back here to the case studies described in the introduction and try to solve them using the methods presented above.

### 7.1 Case study 1: State equation of a gaz, a solution using bisection

We use the bisection method to solve case study 1 and compute the volume of 1000 molecules of  $\text{CO}_2$  at temperature  $T = 300 \text{ K}$  and pressure  $p = 3.5 \cdot 10^7 \text{ Pa}$ . We want to compute the corresponding volume with tolerance  $10^{-12}$ .

To do so, we have to solve the following equation for  $V$ :

$$f(V) = \left[ p + a \left( \frac{N}{V} \right)^2 \right] (V - Nb) - kNT = 0$$

with  $N = 1000$ ,  $k = 1.3806503 \cdot 10^{-23} \text{ J K}^{-1}$ ,  $a = 0.401 \text{ Pa m}^6$  and  $b = 42.7 \cdot 10^{-6} \text{ m}^3$ .

**Do it yourself.** Solve the problem using the bisection method. Feel free to provide any result/remark/plot you consider as interesting to solve this problem. Suggestions:

- plot of the function  $V \rightarrow f(V)$
- initialization for the algorithm
- approximation of the solution computed
- behaviour of the algorithm (precision, number of iteration...)

Write your remarks/answer in a text cell.

```

1 ## Function f
2
3 def fgaz(V):
4     k = 1.3806503e-23
5     a = 0.401
6     b = 42.7e-6
7     N = 1000.0
8     T = 300.0
9     p = 3.5e7
10    return ---

```

### 7.2 Case study 2: Investment found, solutions using bisection or Newton's method

Here, we use the bisection method to solve case study 2. We recall that we have to find  $i$  solution to

$$f(i) = d \frac{(1+i)^{n_{\text{end}}} - 1}{i} - S = 0 \quad \text{where} \quad S = 30\,000, \quad d = 30, \quad \text{and} \quad n_{\text{end}}$$

**Do it yourself.** Solve the problem using the bisection method and/or Newton's method. Feel free to provide any result/remark/plot you consider as interesting to solve this problem. Suggestions:

- plot of the function  $i \rightarrow f(i)$
- initialization for the algorithms
- approximation of the solutions computed
- behaviour of the algorithms (precision, number of iteration...)
- comparison of the algorithms

Write your remarks/answer in a text cell.

```

1 ## Function f
2
3 def finterest(i):
4     d = 30.0
5     S = 30000.0
6     n = 120.0
7     return ---

```

### 7.3 Case study 3: A first population model, a solution using Newton's method



We want to find an approximation for the natural growth rate  $\lambda$  in France. To do so, we have to solve the following non-linear equation for  $\lambda$  (we know that  $\lambda \neq 0$  since the population increases more than the migratory balance):

$$f(\lambda) = N(2017) - N(2016) \exp(\lambda) - \frac{r}{\lambda}(\exp(\lambda) - 1)$$

where  $N(2016)=66\,695\,000$ ,  $N(2017)=66\,954\,000$  and  $r=67\,000$ .

**Do it yourself.** Solve the problem using the bisection method and/or Newton's method. Feel free to provide any result/remark/plot you consider as interesting to solve this problem. Suggestions:

- initialization for the algorithm
- approximation of the solution computed
- behaviour of the algorithm (precision, number of iteration...)
- using the value of  $\lambda$  you computed, and assuming that the migratory balance will be the same in 2017, you can compute an estimation of the population in France at the beginning of year 2018

Write your remarks/answer in a text cell.

```
1 def fpop(l):
2     N0 = 66695000.0
3     N1 = 66954000.0
4     r = 67000.0
5     return ---
```

## 8 Appendix

### 8.1 Intermediate value thm

#### Theorem.

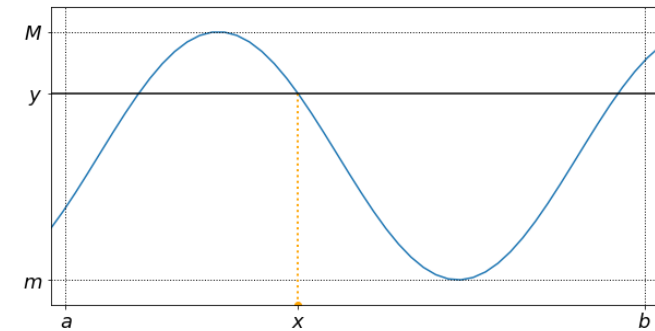
#### Intermediate value Theorem

Suppose  $f : [a, b] \mapsto \mathbb{R}$  is continuous on  $[a, b]$ . Define  $m = \min\{f(a), f(b)\}$  and  $M = \max\{f(a), f(b)\}$ . Then,

$$\forall y \in ]m, M[, \quad \exists x \in ]a, b[, \quad \text{such that} \quad f(x) = y.$$

As a consequence, if a continuous function has values of opposite signs in an interval, it has a root in this interval.

The following figure provides an example of choice for  $x$  guaranteed by this theorem. In this case, the choice is not unique.



```
1 # execute this part to modify the css style
2 from IPython.core.display import HTML
3 def css_styling():
4     styles = open("./style/custom3.css").read()
5     return HTML(styles)
6 css_styling()
```

1