

# Solving non linear systems

Aksel Hiorth, the National IOR Centre & Institute for Energy  
Resources,

University of Stavanger

Aug 4, 2020

## Contents

<b>1</b>	<b>Nonlinear equations</b>	<b>2</b>
<b>2</b>	<b>Example: van der Waals equation of state</b>	<b>3</b>
1:	van der Waal EOS and CO <sub>2</sub> . . . . .	5
<b>3</b>	<b>Fixed-point iteration</b>	<b>5</b>
2:	Implement the fixed point iteration . . . . .	5
3:	Finding the molar volume from the van der Waal EOS by fixed point iteration . . . . .	6
3.1	When do the fixed point method fail? . . . . .	7
3.2	What to do when the fixed point method fail . . . . .	8
4:	Solve $x = e^{1-x^2}$ using fixed point iteration . . . . .	8
<b>4</b>	<b>Rate of convergence</b>	<b>9</b>
<b>5</b>	<b>The bisection method</b>	<b>9</b>
<b>6</b>	<b>Newtons method</b>	<b>12</b>
5:	Compare Newtons, Bisection and the Fixed Point method . . . . .	14
<b>7</b>	<b>Newton Rapson method</b>	<b>16</b>
<b>8</b>	<b>Gradient Descent</b>	<b>16</b>
6:	Gradient descent solution of linear regression . . . . .	16
<b>9</b>	<b>Other Useful Methods</b>	<b>19</b>

In this chapter we will cover some theory related to the solution of nonlinear equations, and introduce the most used methods. A nonlinear problem is represented as a single equation or a system of equations, where the response is not changing proportionally to the input. Almost all physical systems are nonlinear, and one frequent use of the methods presented in this chapter is to determine model parameters by matching a nonlinear model to data.

## 1 Nonlinear equations

A nonlinear equation is simply an equation that is not linear. That means that when the variables changes the response is not changing proportional to the values of the variables. Solving a nonlinear equation always proceeds by *iterations*, we start with one or several initial guesses and then search for the solution. In many cases we do not know beforehand if the equation actually has a solution, or multiple solutions. An example of a nonlinear problem is:

$$e^{-x} = x^2. \quad (1)$$

Traditionally one collect all the terms on one side, to solve an equation of the form

$$f(x) = x^2 - e^{-x} = 0. \quad (2)$$

In figure 1, the solution is shown graphically. Note that in one case the solution is when the graph of  $e^{-x}$ , and  $x^2$  intersect, whereas in the other case the root is located when  $x^2 - e^{-x}$  intersect the  $x$ -axis.

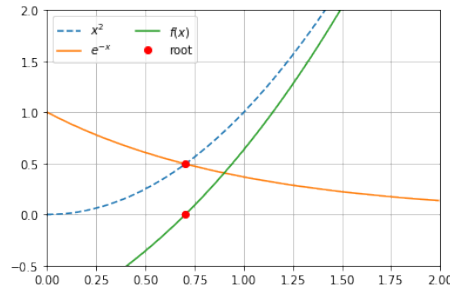


Figure 1: Notice that the root is located at the same place ( $x = 0.703467417$ )

In the case of more than one unknown, or a set of equations that must be satisfied simultaneously, equation (2) is replaced with a vector equation

$$\mathbf{f}(\mathbf{x}) = \mathbf{0}. \quad (3)$$

Although this equation looks quite similar to equation (2), this equation is *much* harder to solve. The only methods we will cover is the Newton Rapson method, which is a very good method if a good starting point is given. If you have a multidimensional problem, the advice is to try Newton-Raphson, if this method fails you need to try more advanced method, see e.g. [3].

## 2 Example: van der Waals equation of state

Before we begin with the numerical algorithms, let us consider an example: the van der Waals equation of state. The purpose is to illustrate some of the typical challenges. You are probably familiar with the ideal gas law:

$$P\nu = R_g T, \quad (4)$$

where  $\nu = V/n$  is the molar volume of the gas,  $P$  is the pressure,  $V$  is the volume,  $T$  is the temperature,  $n$  is the number of moles of the gas, and  $R_g$  is the ideal gas constant. This equation is an example of an *equation of state* (EOS), it relates  $P$ ,  $T$ , and  $\nu$ . Thus if we know the pressure and temperature of the gas, we can calculate  $\nu$ . Equation (4) assumes that there are no interactions between the molecules in the gas. Clearly, this is too simplistic, and because of this one normally uses an EOS that better reflect the physical properties of the substance. A very famous EOS is the van der Waal EOS, which is a slight modification of equation (4):

$$\left(P + \frac{a}{\nu^2}\right)(\nu - b) = R_g T. \quad (5)$$

$a$  and  $b$  are material constants that needs to be determined experimentally. This equation is *not* used in industrial design, but most equations used in practice are based on equation (5). Multiplying equation (5) with  $\nu^2$ , we get a non linear equation that is cubic in the molar volume. It turns out that cubic EOS are a class of equations that are quite successful in modeling the behavior of real systems [2]. However equation (5) is a good starting point for more complex and realistic equations.

It is common practice to rescale EOS with respect to the critical point. At the critical point we have [ref]:

$$\left.\frac{\partial P}{\partial \nu}\right|_{T_c, P_c} = 0 \quad (6)$$

$$\left.\frac{\partial^2 P}{\partial \nu^2}\right|_{T_c, P_c} = 0 \quad (7)$$

From equation (6), (7), and (5), it follows:

$$\nu_c = 3b \quad , P_c = \frac{a}{27b^2} \quad , R_g T_c = \frac{8a}{27b^2}. \quad (8)$$

Inserting these equations into equation (5), and defining the *reduced* quantities  $\hat{P} = P/P_c$ ,  $\hat{T} = T/T_c$ ,  $\hat{\nu} = \nu/\nu_c$ , we get

$$\left(\hat{P} + \frac{3}{\hat{\nu}^2}\right)(3\hat{\nu} - 1) = 8\hat{T}. \quad (9)$$

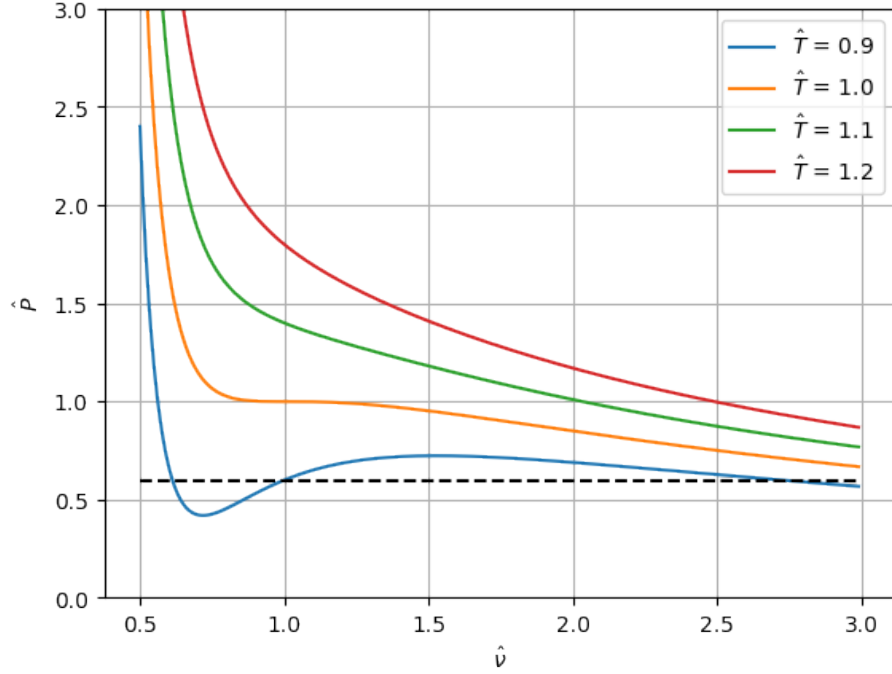


Figure 2: van der Waal isotherms.

In figure 2, we have plotted the isotherms. Note that if  $\hat{T} < 1$  ( $T < T_c$ ), there might be more than one solution for the molar volume. This is clearly unphysical and additional constraints are needed. For the curve  $\hat{T} = 0.9$ , the dashed lined shows that for  $\hat{P} = 0.7$ , there are three solutions. This is a typical behavior of the cubic EOS, and physically it corresponds to the saturated case, where the vapor and liquid phase co-exist. The left root is the liquid state and the right root is the vapor state. The root in the middle represents a meta stable state.

#### It never hurts to look at your function.

The example in figure 2 illustrates some important points. Solving a nonlinear problem might be very easy in part of the parameter space (e.g. when  $T > T_c$  there are only one solution), but extremely hard in other

part of the parameter space (e.g. when  $T < T_c$ , where there are multiple solutions). However, much of the trick to find a solution is to choose a good starting point. When there are multiple solutions we need to start close to the physical solution.

### Exercise 1: van der Waal EOS and CO<sub>2</sub>

Use equation (5), and the parameters for CO<sub>2</sub>:  $a=3.640 \text{ L}^2\text{bar/mol}$ , and  $b=0.04267 \text{ L/mol}$ , to test the van der Waal EOS in equation (5). Use that at 2 MPa and 100 °C, CO<sub>2</sub> has a specific volume of  $0.033586 \text{ m}^3/\text{kg}$ .

**Solution.** The calculation is straight forward, but it is easy to get an error due to units. We will use SI units:  $a=0.3640 \text{ m}^6\text{Pa/mol}$ ,  $b=4.267 \cdot 10^{-5} \text{ m}^3/\text{mol}$ ,  $R=8.314 \text{ J/mol K}$ . The molar volume is obtained by multiplying by the molar weight of CO<sub>2</sub>:  $M_w = 44 \text{ g/mol}$ , hence  $\nu = 1.478 \cdot 10^{-3} \text{ m}^3/\text{mol}$ . Using  $P = RT/(\nu - b) - a/\nu^2 = 1.993 \text{ MPa}$ , or an error of 0.3%.

## 3 Fixed-point iteration

A simple (but not always possible) way of solving a nonlinear equation is to reformulate the problem  $f(x) = 0$  to a problem of the form

$$x = g(x). \quad (10)$$

The algorithm for solving this equation is to guess at a starting point,  $x_0$ , evaluate  $x_1 = g(x_0)$ ,  $x_2 = g(x_1)$ , and so on. In some circumstances we might end up at a stable point, where  $x$  does not change. This point is termed a *fixed point*.

Note that the form of  $g(x)$  is not uniquely determined. For our function defined in equation (1), we can solve for  $x$  directly

$$x = e^{-x/2}, \quad (11)$$

or we could write:

$$x = x - x^2 + e^{-x}. \quad (12)$$

These functions are illustrated in figure 3, by visual inspection they look very similar, but as we will show in the next exercise the convergence is quite different.

### Exercise 2: Implement the fixed point iteration

Write a Python function that utilizes the fixed point algorithm in the previous section, find the root of  $f(x) = x^2 - e^{-x}$ . In one case use  $g(x) = e^{-x/2}$ , and in the other case use  $g(x) = x - x^2 + e^{-x}$ . How many iterations does it take in each case?

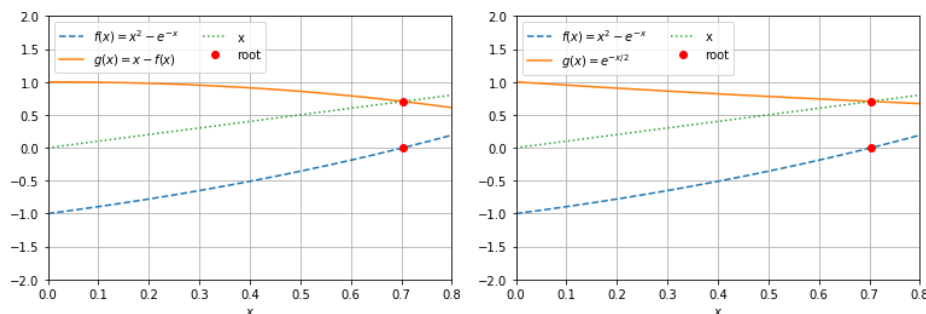


Figure 3: Two examples of iterative functions, that will give the same solution.

**Solution.** Below is a straight forward (vanilla) implementation:

```
def iterative(x,g,prec=1e-8, MAXIT=1000):
    '''Approximate solution of x=g(x) by fixed point iterations.
    x : starting point for iterations
    eps : desired precision
    Returns x when x does not change more than prec
    and number of iterations MAXIT are not exceeded
    '''
    eps = 1
    n=0
    while eps>prec and n < MAXIT:
        x_next = g(x)
        eps = np.abs(x-x_next)
        x = x_next
        n += 1
        if(np.isinf(x)):
            print('Quitting .. maybe bad starting point?')
            return x
    if (n<MAXIT):
        print('Found solution: ', x, ' After ', n, 'iterations')
    else:
        print('Max number of iterations exceeded')
    return x
```

If we start at  $x = 0$ , it will take 174 iterations using  $x - x^2 + e^{-x}$  ( $g(x)$ ) and only 19 for  $e^{-x/2}$  ( $h(x)$ ), the root is  $x=0.70346742$ .

### Exercise 3: Finding the molar volume from the van der Waal EOS by fixed point iteration

Extend the code above to take as argument the van der Waal EOS. For simplicity we will use the rescaled EOS in equation (9). Show that for the reduced temperature,  $\hat{T}=1.2$ , and pressure,  $\hat{P}=1.5$ , the reduced molar volume  $\hat{n}u$  is 1.3522091.

**Solution.** First we rewrite equation (9) in a more useful form

$$\hat{\nu} = \frac{1}{3} \left( 1 + \frac{8\hat{T}}{\hat{P} + 3/\hat{\nu}^2} \right) \quad (13)$$

The right hand side will play the same role as  $g(x)$  above, where  $x$  now is the reduced molar volume, and can be implemented in Python as:

```
def dvdwEOS(nu,t,p):
    return (1+8*t/(p+3/nu**2))/3
```

Note that this function requires the values of  $\hat{P}$  and  $\hat{T}$ , in addition to  $\hat{\nu}$  to return a value. Thus in order to use the fixed point iteration method implemented above, we need to pass arguments to our function. This can easily be achieved by taking advantage of Python's `*args` functionality. By simply rewriting our implementation slightly:

```
def iterative(x,g,*args,prec=1e-8):
    MAX_ITER=1000
    eps = 1
    n=0
    while eps>prec and n < MAX_ITER:
        x_next = g(x,*args)
        eps = np.abs(x-x_next)
        x = x_next
        n += 1
    print('Number of iterations: ', n)
    return x
```

We can find the root by calling the function as:

```
iterative(1,dvdwEOS,1.2,1.5)
```

The program returns the correct solution after 71 iterations.

### 3.1 When do the fixed point method fail?

If we replace  $e^{-x}$  with  $e^{1-x^2}$  in equation (12), our method will not give a solution. You can easily verify that the  $x = 1$  is a solution, so why does our method fail? To investigate this in a bit more detail, we turn to Taylor's formula (once again). Assume that the root is located at  $x^*$ , and our guess is  $x_k$ , then the next  $x$ -value will be

$$x_{k+1} = g(x_0) = g(x^*) + g'(x^*)(x_k - x^*) + \dots \quad (14)$$

The true solution is  $x^*$ , hence  $x^* = f(x^*)$ , and we can write

$$x_{k+1} - x^* = g'(x^*)(x_k - x^*), \quad (15)$$

where we have neglected higher order terms. The point is: at each iteration we want the distance  $x_1 - x^*$  to decrease, i.e. to be smaller than  $x_0 - x^*$ . This can

only be achieved if

$$|g'(x^*)| < 1. \quad (16)$$

In our example above we saw that if  $g(x) = x - x^2 + e^{-x}$ , we used 172 iterations and only 19 iterations if we replaced  $g(x)$  with  $h(x) = e^{-x/2}$  to converge to the *same* root  $x=0.70346742$ . We can now understand this, because  $g'(x) = 1 - 2x - e^{-x}$  and  $g(x^*) \simeq -0.90$ , whereas  $h'(x) = -e^{-x/2}/2$ , and  $h'(x^*) \simeq 0.35$ . We expect the number of iterations,  $n$ , needed to reach a certain precision,  $\varepsilon$ , to scale as

$$|g'(x^*)|^n = \varepsilon. \quad (17)$$

We expect to use  $\log |h'(x^*)| / \log |g'(x^*)| \simeq 10$  more iterations using  $g(x)$  compared to  $h(x)$ , which is close to the observed value of  $172/19 \simeq 9$ .

### 3.2 What to do when the fixed point method fail

As discussed in [1], there might be an elegant solution whenever  $|g'(x^*)| > 1$ . If it is possible to invert the  $g(x)$ , we can show that the derivative of the inverse function  $g'^{-1}(x^*) = 1/g'(x^*)$ . Why is this useful? Because if  $x^* = g(x^*)$  is the solution we are searching for, then this is equivalent to  $x^* = g^{-1}(x^*)$  *if and only if* we can invert  $g(x)$ . Note that in many cases it is not possible to invert  $g(x)$ . Let us first show that  $g'^{-1}(x^*) = 1/g'(x^*)$ . For simplicity write

$$y = g(x) \Leftarrow x = g^{-1}(y), \quad (18)$$

taking the derivative with respect to  $x$  gives

$$\frac{d}{dx} g^{-1}(y) = \frac{dx}{dy} = 1, \quad (19)$$

$$\frac{dg^{-1}(y)}{dy} \frac{dy}{dx} = \frac{dx}{dx} = 1, \quad (20)$$

$$\frac{dg^{-1}(y)}{dy} = \frac{1}{\frac{dy}{dx}} = \frac{1}{g'(x)} = \frac{1}{g'(g^{-1}(y))}. \quad (21)$$

Going from equation (19) to (20), we have used the chain rule. Equation (21) is general, let us now specify to our fixed point iteration. Then we can use  $x^* = g(x^*) = y^*$ , and  $x^* = g^{-1}(y^*) = g^{-1}(x^*)$  hence we can write the last equation as

$$\frac{d}{dx} g^{-1}(x^*) = \frac{1}{g'(x^*)}. \quad (22)$$

#### Exercise 4: Solve $x = e^{1-x^2}$ using fixed point iteration

The solution to  $x = e^{1-x^2}$  is clearly  $x = 1$ .

- First try the fixed point method using  $g(x) = e^{1-x^2}$  to find the root  $x = 1$ . Try to start very close to the true solution  $x = 1$ . What is the value of  $g'(x^*)$ ?



- Next, invert  $g(x)$ , what is the derivative of  $g^{-1}(x^*)$ ? Try the fixed point method using  $g^{-1}(x^*)$

**Solution.** First, we calculate the derivative of  $g(x)$ ,  $g'(x) = -2xe^{1-x^2}$ , hence  $g'(x^*) = -2$  and  $|g'(x^*)| > 1$ . This is an unstable fixed point, and if we start a little bit off from this point we will spiral away from it.

Inverting  $y = g(x)$  gives us  $g^{-1}(y) = \sqrt{1 - \ln y}$ . Note that  $y^* = x^* = 1$  is a solution to this equation as it should be. The derivative is

$$g^{-1}'(y) = -\frac{1}{2\sqrt{1 - \ln y}}, \quad (23)$$

and  $g^{-1}'(y^*) = -1/2$ . It takes about 30 iterations to reach the correct solution  $y^* = 1$ , when the starting point is  $y = 0$ .

## 4 Rate of convergence

The rate of convergence is the speed at which a *convergent* sequence approach the limit. Assume that our sequence  $x_k$  converges to the number  $x^*$ , the sequence is said to *converge linearly* to  $x^*$  if there exists a number  $\mu \in (0, 1)$ , such that

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|} = \mu \quad (24)$$

Inserting equation (15) in equation (24), we get:

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - x_k|}{x_k - x^*} = \frac{|g'(x^*)(x_k - x^*)|}{|x_k - x^*|} = |g'(x^*)|. \quad (25)$$

Hence the fixed point iteration is expected to converge *linearly* to the correct solution. The definition in equation (24), can be extended to include the definition of quadratic, cubic, etc. convergence:

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|^q} = \mu. \quad (26)$$

If  $q = 2$  the convergence is said to be quadratic and so on.

## 5 The bisection method

The idea behind bisection is that the root is bracketed, i.e. that there exists two points  $a$  and  $b$ , such that  $f(a) \cdot f(b) < 0$ . In practice it might be a challenge to find these two points. However, if you know that the function has a only root between two values, and that speed is not a big issue this method guarantees that the root will be found within a finite number of steps. The basic idea behind the method is to divide the interval into two (i.e. bisecting the interval). The method only works if the function is continuous on the interval.

The algorithm is as follows:

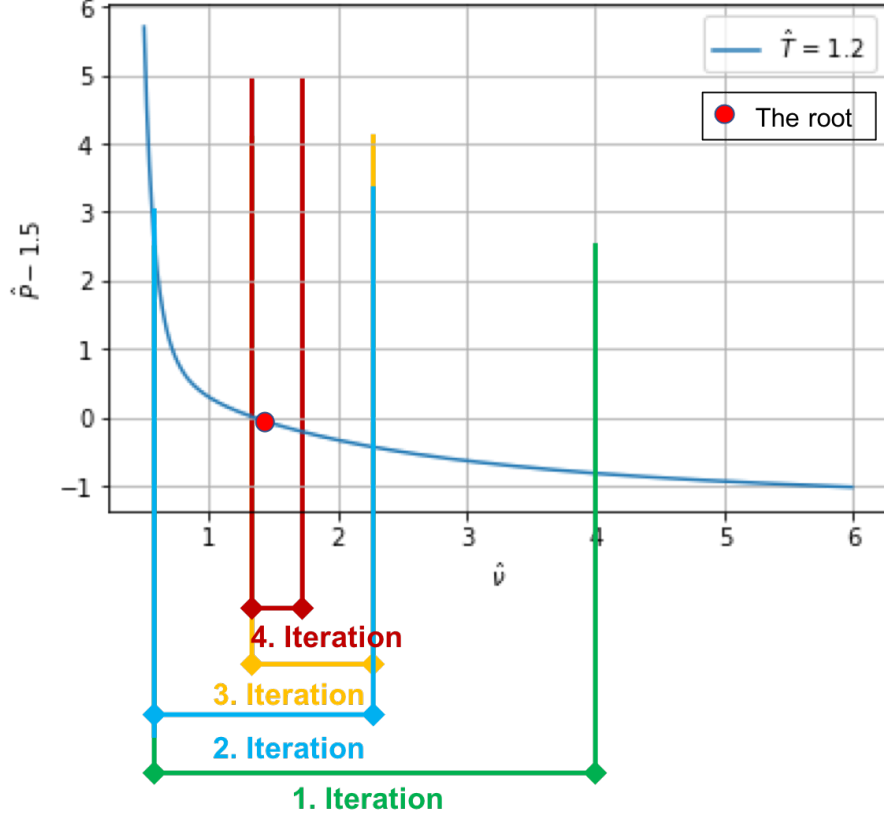


Figure 4: Illustration of the bisection method for the van der Waal EOS.

- Test if  $f(a) \cdot f(b) < 0$ , if not return an error message
- Calculate the midpoint  $c = (a + b)/2$ . If  $f(a) \cdot f(c) < 0$  the root is in the interval  $[a, c]$ , else the root is in the interval  $[c, b]$
- Half the interval, and test in which interval the root lies, and continue until a convergence criterion.

In figure 4, a graphical illustration of the method is shown. If  $c_n$  is the midpoint after  $n$  steps, the difference between the solution  $x^*$  and  $c_n$  is

$$|c_n - x^*| \leq \frac{|b - a|}{2^n} \quad (27)$$

Using our previous definition in equation (26), we find that

$$\lim_{k \rightarrow \infty} \frac{|c_{k+1} - x^*|}{|c_k - x^*|} \leq \frac{|b - a|/2^{n+1}}{|b - a|/2^n} = \frac{1}{2}, \quad (28)$$

hence the bisection method converges linearly. Below is an implementation of the bisection method.

```
def bisection(f,a,b,prec=1e-8,MAXIT=100):
    '''Approximate solution of  $f(x)=0$  on interval  $[a,b]$  by bisection.

    f :  $f(x)=0$ .
    a,b : brackets the root  $f(a)*f(b)$  has to be negative
    eps : desired precision

    Returns the midpoint when it is closer than eps to the root,
    unless MAX_ITERATIONS are not exceeded
    '''
    MAX_ITERATIONS=MAXIT
    if f(a)*f(b) >= 0:
        print('You need to bracket the root,  $f(a)*f(b) >= 0$ ')
        return None
    an = a
    bn = b
    c_old = 0.5*(an + bn)
    for n in range(MAX_ITERATIONS):
        cn = c_old
        f_cn = f(cn)
        if f(an)*f_cn < 0:
            bn = cn
        elif f(bn)*f_cn < 0:
            an = cn
        elif f_cn == 0:
            print('Found exact solution ', cn,
                  ' after ', n, 'iterations' )
            return cn
        else:
            print('Bisection method fails.')
            return None
        cn = 0.5*(an+bn)
        if(abs(cn-c_old)<prec):
            print('Found solution ', cn,
                  ' after ', n, 'iterations' )
            return cn
        c_old=cn
    print('Max number of iterations: ', MAXIT, ' reached.')
    print('Try to increase MAXIT or decrease prec')
    print('Returning best guess, value of function is: ', f_cn)
    return cn
```

### Warnings.

Note that the implementation of the bisection algorithm is only a few lines of code, and most of the code is to give warnings to the user. In this case it is important to do additional checking, and give the user warnings. If  $f(c)=0$ , then we must stop and return the exact solution. If we only test if  $f(a) \cdot f(c)$  is greater or lower than zero the algorithm would fail.

## 6 Newtons method

Newtons method is one of the most used methods. If it converges, it converges quadratically to the correct solution. The drawback is that contrary to the bisection method it may fail if a bad starting point is given. Newtons method for finding the root of a function  $f(x) = 0$  is illustrated in figure 5. The main idea is to use more information about the function in the search of the root. In this case we want to find the point where the tangent of the function in  $x_k$  intersect the  $x$ -axis, and take that as our next point,  $x_{k+1}$ .

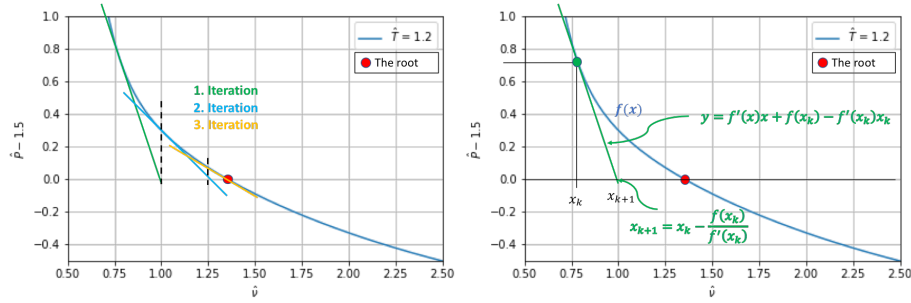


Figure 5: Illustration of Newtons method for the van der Waals EOS.

We can easily derive the algorithm by finding the formula for the tangent line. Using  $y = ax + b$  for the tangent line, we immediately know that  $a = f'(x_k)$ .  $b$  can be found as we know that the line intersects  $(x_k, f(x_k))$ :  $f(x_k) = f'(x_k)x_k + b$ , hence the equation for the tangent line is  $y = f'(x_k)x + f(x_k) - f'(x_k)x_k$ . The next point is located where  $y$  crosses the  $x$ -axis, hence  $0 = f'(x_k)x_{k+1} + f(x_k) - f'(x_k)x_k$ . Rearranging this equation, we can write Newtons method in the standard form

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}. \quad (29)$$

Note that the derivative of  $f(x)$  enters in equation (29), which means that if our function has a extremal value in our search domain, Newtons method most likely will fail. In particular  $x_1$ , and  $x_4$  in the figure to the right in figure 6 are bad starting point for Newtons method. Newtons method

Before we show how to implement Newtons method, let us first investigate the convergence and rate of convergence. Comparing figure 4 and 5, you immediately get the sense that Newtons method converges faster. Newtons method is similar

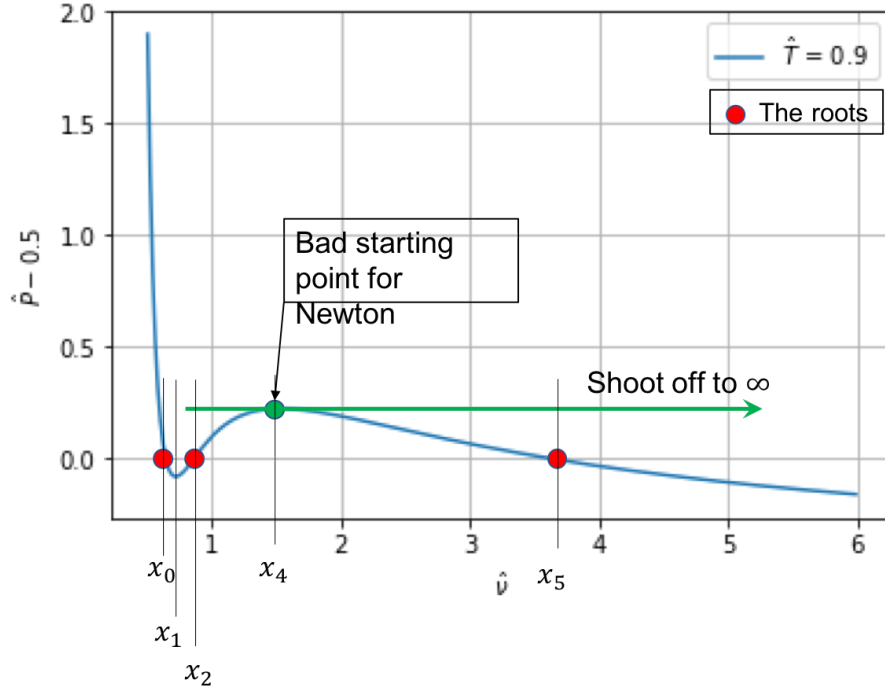


Figure 6: Illustration of some of the possible challenges with Newton's method. Note that if the derivative is zero somewhere in the search interval, Newton's method will fail.

to the fixed point method, but where we do not use  $g(x) = x - f(x)$ , but  $g(x) = x - \frac{f(x)}{f'(x)}$ . We will now analyze Newton's method, using the same approach as in section 3.1. First we expand  $g(x)$  around the root  $x^*$

$$x_{k+1} = g(x_k) = g(x^*) + g'(x^*)(x_k - x^*) + \frac{1}{2}g''(x^*)(x_k - x^*)^2, \quad (30)$$

where we have skipped all higher order terms. You can easily verify that

$$g'(x) = \frac{f''(x)f(x)}{f'(x)^2} \quad (31)$$

$$g''(x) = \frac{(f'''(x)f'(x) - 2f''(x)^2f'(x))f(x) + f''(x)f'(x)^2}{f'(x)^4}. \quad (32)$$

$x^*$  is a solution, hence  $f(x^*) = 0$ , we then find from equation (31) and (32) that  $g'(x^*) = 0$ , and  $g''(x^*) = f'''(x^*)/f'(x^*)^2$ . Thus from equation (30) we get

$$x_{k+1} = x^* + \frac{1}{2} \frac{f'''(x^*)}{f'(x^*)^2} (x_k - x^*)^2, \quad (33)$$

or equivalently:

$$\frac{x_{k+1} - x^*}{(x - x^*)^2} = \frac{1}{2} \frac{f''(x^*)}{f'(x^*)^2}. \quad (34)$$

The denominator has a power of two, and hence Newtons method is *quadratic* convergent (assuming that the sequence  $x_{k+1}$  is a convergent sequence). Note that it also follows from the analyses above that Newtons method will fail if the derivative at the root,  $f'(x^*)$ , is zero. An implementation is shown below.

```
def newton(f,x, prec=1e-8,MAXIT=500):
    '''Approximate solution of f(x)=0 by Newtons method.
    The derivative of the function is calculated numerically
    f    : f(x)=0.
    x    : starting point
    eps  : desired precision

    Returns x when it is closer than eps to the root,
    unless MAX_ITERATIONS are not exceeded
    '''
    MAX_ITERATIONS=MAXIT
    x_old = x
    h      = 1e-4
    for n in range(MAX_ITERATIONS):
        x_new = x_old - 2*h*f(x_old)/(f(x_old+h)-f(x_old-h))
        if(abs(x_new-x_old)<prec):
            print('Found solution:', x_new,
                  ', after:', n, 'iterations.' )
            return x_new
        x_old=x_new
    print('Max number of iterations: ', MAXIT, ' reached.')
    print('Try to increase MAXIT or decrease prec')
    print('Returning best guess, value of function is: ', f(x_new))
    return x_new
```

## Exercise 5: Compare Newtons, Bisection and the Fixed Point method

Find the root of  $f(x) = x^2 - e^{-x}$  using bisection, fixed point, and Newtons method, start at  $x = 0$ . How many iterations do you need to use reach a precision of  $10^{-8}$ ? What happens if you widen the search domain or start further away from the root?

**Solution.** The root is located at  $x^* = 0.70346742$ .

- Fixed point method: we saw earlier that using  $g(x) = x - f(x)$  used 174 iterations, and  $g(x) = \sqrt{x^2 - f(x)}$  used 19 iterations. If we start at  $x = -100$ ,  $g(x) = x - f(x)$  fails, and  $g(x) = \sqrt{x^2 - f(x)}$  uses only 21 iterations, and at  $x = 100$  we use 20 iterations.

- Bisection method: it use 25 iterations for  $a = 0$ , and  $b = 1$  (implementation shown earlier in the chapter). Choosing  $a = -b = -100$  we use 33 iterations.
- Newtons method: it use only 5 function evaluations (implementation above) starting at  $x = 0$ . Starting at  $x = -100$ , it uses 106 iterations. Newtons method is slow in this case because the function is very steep around the starting point, see figure 7. Starting at  $x = 100$ , we only use 10 iterations.

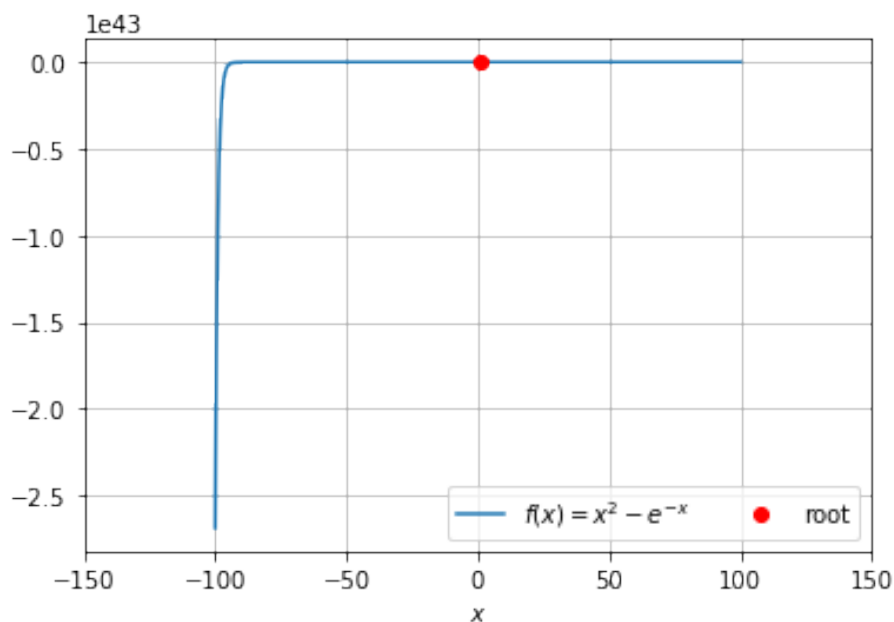


Figure 7: Newtons method performs poorly far away due to the shape of the function close to  $x = -100$ , bisection performs much better while the fixed point method fails.

#### A good starting point is crucial.

Note that it is not given which method is best, but if we are "close" to the root Newtons method is usually superior. If we are far away, other methods might work better. In many cases one uses a more stable method far away from the root, and then "polish up" the root by a couple of Newton iterations [3]. See also Brents method which combines bisection and linear interpolation (secant method) [3].

## 7 Newton Rapson method

The derivation of Newtons method, equation (29), done in the previous section was based on figure 5. We will now derive it using a slightly different approach, but which lends itself easier to extend Newtons method to higher dimensions. The starting point is to expand the function around  $x_k$ , using Taylors formula

$$f(x) = f(x_k) + f'(x_k)(x - x_k) + \dots \quad (35)$$

Equation (29) can be derived from equation (35) by simply demanding that we keep the linear terms, and that the next point  $x_{k+1}$  is located where the linear approximation intersects the  $x$ -axis, i.e. simply set  $f(x) = 0$ , and  $x = x_{k+1}$  in equation (35).

In higher order dimensions, we solve equation (3), and equation (35) is

$$\mathbf{f}(\mathbf{x}) = \mathbf{f}(\mathbf{x}_k) + \mathbf{J}(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) + \dots \quad (36)$$

$\mathbf{J}(\mathbf{x}_k)$  is the Jacobian. As before, we simply set  $\mathbf{f}(\mathbf{x}) = \mathbf{0}$ ,  $\mathbf{x} = \mathbf{x}_{k+1}$ , and keep the linear terms, hence

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{J}^{-1}(\mathbf{x}_k)\mathbf{f}(\mathbf{x}_k). \quad (37)$$

To make the mathematics a bit more clear, let us specify to  $2D$ . Assume that  $\mathbf{f}(\mathbf{x}) = [f_x(x, y), f_y(x, y)]$ , then the Jacobian is

$$\mathbf{J}(\mathbf{x}_k) = \begin{pmatrix} \frac{\partial f_x}{\partial x} & \frac{\partial f_x}{\partial y} \\ \frac{\partial f_y}{\partial x} & \frac{\partial f_y}{\partial y} \end{pmatrix}. \quad (38)$$

## 8 Gradient Descent

This method used is to minimize functions (does not work for root finding). In many nonlinear problems, we would like to minimize (or maximize) a function. An ideal 2D example is shown in figure 8. The algorithm moves in the direction of steepest descent. Note that the step size might change towards the search.

Assume that we have a function  $\mathbf{f}(\mathbf{x})$ , that we would like to minimize. The gradient descent algorithm is simply to update parameters according to the derivative (gradient) of  $\mathbf{f}$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma \nabla \mathbf{f}. \quad (39)$$

$\gamma$  is the learning rate, and a good choice of  $\gamma$  is important.  $\gamma$  might also change from one iteration to the other, and does not have to be constant.

### Exercise 6: Gradient descent solution of linear regression

A very typical example is if we have a model and we would like to fit some parameters of the model to a data set (e.g. linear regression). Assume that we have observations  $(x_i, y_i)$  and model predictions  $f(x_i, \beta)$ , the model parameters



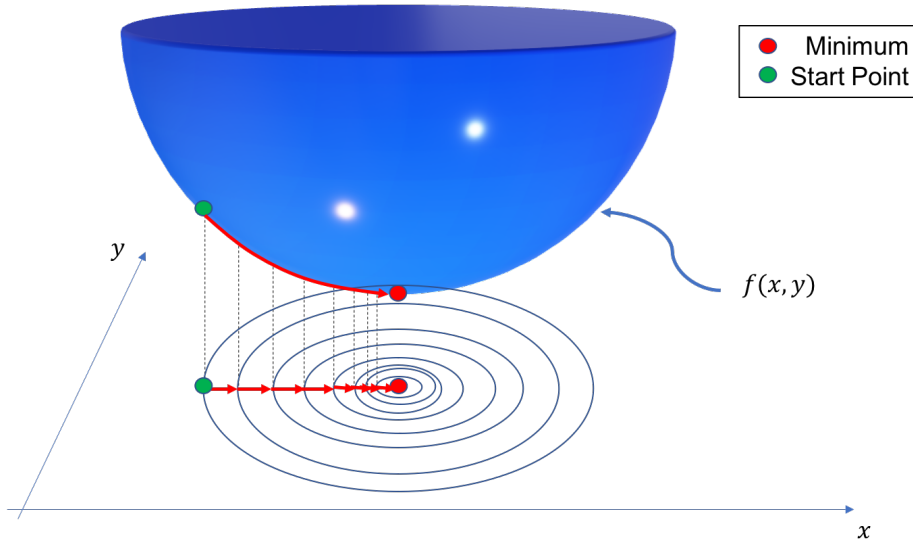


Figure 8: A very simple example of the gradient descent method.

are contained in the vector  $\beta$ . The *least square*,  $S$ , is the square of the sum of all the *residuals*, i.e. the difference between the observations and model predictions

$$S = \sum_i (y_i - f(x_i, \beta))^2. \quad (40)$$

Specializing to linear regression, we choose the model to be linear

$$f(x_i, \beta) = b_0 + b_1 x_i. \quad (41)$$

Equation (40) now takes the form

$$S = \sum_i (y_i - b_0 + b_1 x_i)^2. \quad (42)$$

The gradients are:

$$\begin{aligned} \frac{\partial S}{\partial b_0} &= -2 \sum_i (y_i - b_0 + b_1 x_i), \\ \frac{\partial S}{\partial b_1} &= -2 \sum_i (y_i - b_0 + b_1 x_i) x_i. \end{aligned} \quad (43)$$

- Implement the gradient descent method using a constant learning rate of  $10^{-3}$ , to minimize the least square function
- Test the linear regression on the data set  $x_i = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]$ , and  $y = [1, 3, 2, 5, 7, 8, 8, 9, 10, 12]$ , choose a starting value  $(b_0, b_1) = (0, 0)$ . What happens if you increase the learning rate?

**Solution.** Below is an implementation of the gradient descent method with a constant learning rate

```
def gradient_descent(f,x,df, g=.001, prec=1e-8,MAXIT=10000):
    '''Minimize f(x) by gradient descent.
    f      : min(f(x))
    x      : starting point
    df     : derivative of f(x)
    g      : learning rate
    prec: desired precision

    Returns x when it is closer than eps to the root,
    unless MAX_ITERATIONS are not exceeded
    '''
    MAX_ITERATIONS=MAXIT
    x_old = x
    for n in range(MAX_ITERATIONS):
        x_new = x_old - g*df(x_old)
        if(abs(np.max(x_new-x_old))<prec):
            print('Found solution:', x_new,
                  ', after:', n, 'iterations.' )
            return x_new
        x_old=x_new
    print('Max number of iterations: ', MAXIT, ' reached.')
    print('Try to increase MAXIT or decrease prec')
    print('Returning best guess, value of function is: ', f(x_new))
    return x_new
```

The linear regression is implemented as below

```
x_obs_ = np.array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
y_obs_ = np.array([1, 3, 2, 5, 7, 8, 8, 9, 10, 12])
def S(b,x=x_obs_,y=y_obs_):
    return np.sum((y-b[0]-b[1]*x)**2)

def dS(b,x=x_obs_,y=y_obs_):
    return np.array([-2*np.sum(y-b[0]-b[1]*x),
                     -2*np.sum((y-b[0]-b[1]*x)*x)])

b=np.array([0,0])
```

The first four iterations are shown in figure 9. If we choose a learning rate that is too high, we will move past the minimum, and the solution will oscillate. This can be avoided by lowering the learning rate as we iterate, by e.g. replacing  $g$  with  $g/(n+1)$  in the implementation above.

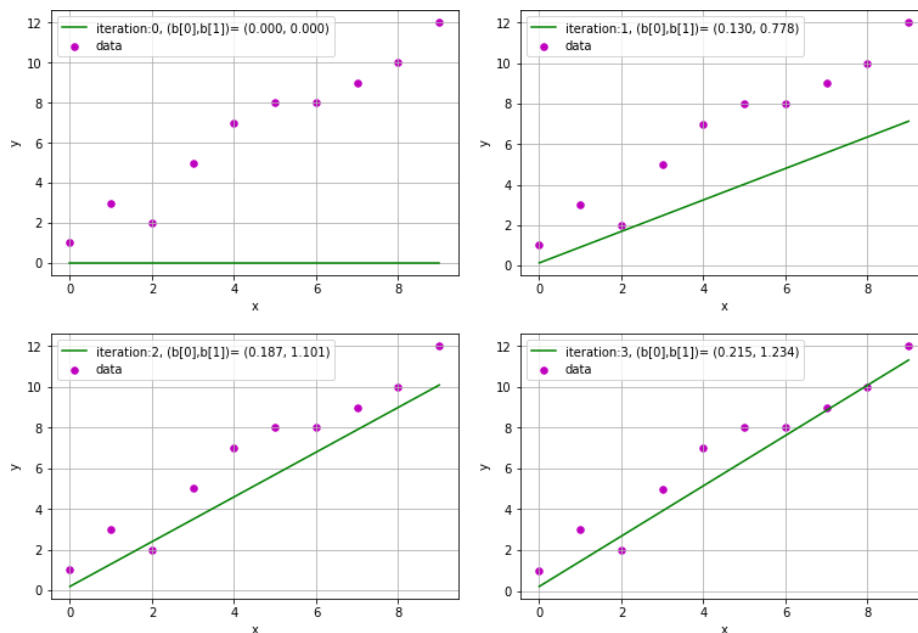


Figure 9: First four iterations of the gradient descent solution of linear regression.

## 9 Other Useful Methods

In this chapter we have covered the *basic*, but you should now be well equipped to dive into other methods. We highly recommend [3] as a starting point, although the code examples are written in C++, the theory is presented in a very accurate, but informal way.

- Brents method: uses root bracketing, bisection, and inverse quadratic interpolation. The 1D method of choice if the function and not its derivative is known
- Secant method: uses linear interpolation to find the root (only function values, not derivatives)

## References

- [1] Mark Newman. *Computational Physics*. CreateSpace Independent Publ., 2013.
- [2] Ding-Yu Peng and Donald B. Robinson. A new two-constant equation of state. *Industrial & Engineering Chemistry Fundamentals*, 15(1):59–64, 1976.

- [3] William H. Press, William T. Vetterling, Saul A. Teukolsky, and Brian P. Flannery. *Numerical Recipes in C++: the Art of Scientific Computing*. Cambridge University Press, New York, NY, USA, 2nd edition, 2002.