

# MA 8019: Numerical Analysis I

## Solution of Nonlinear Equations



Suh-Yuh Yang (楊肅煜)

Department of Mathematics, National Central University  
Jhongli District, Taoyuan City 320317, Taiwan

First version: May 4, 2018    Last updated: September 4, 2024

## Introduction

---

- **A nonlinear equation:**

Let  $f : \emptyset \neq A \subseteq \mathbb{R} \rightarrow \mathbb{R}$  be a nonlinear real-valued function in a single variable  $x$ . *We are interested in finding the roots (solutions) of the equation  $f(x) = 0$ , i.e., zeros of the function  $f(x)$ .*

- **A system of nonlinear equations:**

Let  $F : \emptyset \neq A \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^n$  be a nonlinear vector-valued function in a vector variable  $X = (x_1, x_2, \dots, x_n)^\top$ . *We are interested in finding the roots (solutions) of the equation  $F(X) = \mathbf{0}$ , i.e., zeros of the vector-valued function  $F(X)$ .*

## Example: zeros of polynomial

---

- Let us look at three functions (polynomials):
  - (1)  $f(x) = x^4 - 12x^3 + 47x^2 - 60x$
  - (2)  $f(x) = x^4 - 12x^3 + 47x^2 - 60x + 24$
  - (3)  $f(x) = x^4 - 12x^3 + 47x^2 - 60x + 24.1$
- Find the zeros of these polynomials is not an easy task.
  - (1) The first function has *real zeros 0, 3, 4, and 5*.
  - (2) The real zeros of the second function are *1 and 0.888...*
  - (3) The third function *has no real zeros at all*.
  - (4) MATLAB: see `polyzeros.m`
- *The  $n$  roots of a polynomial of degree  $n$  depend continuously on the coefficients.* (see Complex Analysis)
  - (1) This result implies that the eigenvalues of a matrix depend continuously on the matrix. (see Tyrtyshnikov's book).
  - (2) However, the problem of approximating the roots given the coefficients is *ill-conditioned*, see Wilkinson's polynomial.  
[https://en.wikipedia.org/wiki/Wilkinson%27s\\_polynomial](https://en.wikipedia.org/wiki/Wilkinson%27s_polynomial)

## Objectives

---

Consider the nonlinear equation  $f(x) = 0$  or  $F(X) = \mathbf{0}$ .

- The basic questions:
  - (1) Does the solution exist?
  - (2) Is the solution unique?
  - (3) *How to find it?*
- We will mainly focus on the third question and we always assume that the problem under considered has a solution  $x^*$ .
- *We will study iterative methods for finding the solution:* first find an initial guess  $x_0$ , then a better guess  $x_1, \dots$ , in the end we hope that  $\lim_{n \rightarrow \infty} x_n = x^*$ .
- **Iterative methods:** bisection method; Newton's method; secant method; fixed-point method; continuation method; special methods for zeros of polynomials.

## Bisection method (method of interval halving)

---

- **An observation:** *If  $f(x)$  is a continuous function on an interval  $[a, b]$ , and  $f(a)$  and  $f(b)$  have different signs such that  $f(a)f(b) < 0$ , then  $f(x)$  must have a zero in  $(a, b)$ , i.e., a root of the equation  $f(x) = 0$ .*

*(ensured by the Intermediate-Value Theorem for continuous functions)*

- **The basic idea:** assume that  $f(a)f(b) < 0$ .
  - (1) compute  $c = \frac{1}{2}(a + b) = a + \frac{1}{2}(b - a)$ .
  - (2) if  $f(a)f(c) = 0$ , then  $f(c) = 0$  and  $c$  is a zero of  $f(x)$ .
  - (3) if  $f(a)f(c) < 0$ , then the zero is in  $[a, c]$ ; otherwise the zero is in  $[c, b]$ . In either case, a new interval containing the root is produced, and the size of the new interval is half of the original one.
  - (4) repeat the process until the interval is very small then any point in the interval can be used as approximations of the zero.

## What do we need?

---

- We need an initial interval  $[a, b]$ . This is often the hardest thing to find.
- We need some stopping criteria: given  $\varepsilon > 0$  and  $\delta > 0$  are tolerances,  $k$  is the number of iterations.
  - (1) if  $|f(c)| < \varepsilon$ , we stop.
  - (2) if  $|b - a| < \delta$ , we stop.
  - (3) if  $k > M$ , we stop to avoid infinite loop.

## A pseudocode for the bisection algorithm

---

```
input  $a, b, M, \delta, \varepsilon$   
 $u \leftarrow f(a), v \leftarrow f(b), e \leftarrow b - a$   
output  $a, b, u, v$   
if  $\text{sign}(u) = \text{sign}(v)$  then stop  
for  $k = 1$  to  $M$  do  
     $e \leftarrow e/2, c \leftarrow a + e, w \leftarrow f(c)$   
    output  $k, c, w$   
    if  $|e| < \delta$  or (and)  $|w| < \varepsilon$  then stop  
    if  $\text{sign}(w) \neq \text{sign}(u)$  then  
         $b \leftarrow c, v \leftarrow w$   
    else  
         $a \leftarrow c, u \leftarrow w$   
    end if  
end do
```

---

### Note:

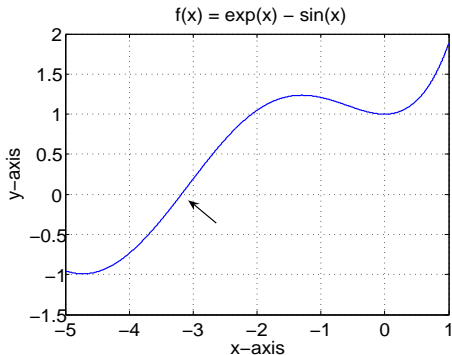
- $\text{sign}(w) \neq \text{sign}(u)$  is better than  $wu < 0$ . (why?)
- compute midpoint as  $c = a + \frac{b-a}{2}$  rather than  $c = \frac{a+b}{2}$ . (why?)

## An example

---

Use the bisection method to find the root of  $e^x = \sin(x)$ .

A rough plot of  $f(x) = e^x - \sin(x)$  shows there are no positive zeros, and the first zero to the left of 0 is somewhere in the interval  $[-4, -3]$ .



see [functiongraph1.m](#)



## Numerical results

---

The output obtained by bisection algorithm running a MATLAB M-file, `bisection.m`

*Starting with  $a = -4$  and  $b = -3$ :*

$k$	$c$	$f(c)$
1	-3.500000000000000	-0.32058584426730
2	-3.250000000000000	-0.06942092669839
3	-3.125000000000000	0.06052882585276
4	-3.187500000000000	-0.00461629388698
⋮	⋮	⋮
13	-3.18298339843750	0.00008284596304
14	-3.18304443359375	0.00001933261037
15	-3.18307495117188	-0.00001242395017
16	-3.18305969238281	0.00000345432045
⋮	⋮	⋮

See the details of the M-file: `bisection.m`

## Theorem (on bisection method)

Suppose that  $[a_0, b_0] := [a, b], [a_1, b_1], \dots, [a_n, b_n], \dots$  are the intervals in the bisection method. Then

- (1)  $\lim_{n \rightarrow \infty} a_n$  and  $\lim_{n \rightarrow \infty} b_n$  exist and the limits are equal.
- (2) Let  $r = \lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} b_n$ . Then  $f(r) = 0$ .
- (3) Let  $c_n = a_n + \frac{1}{2}(b_n - a_n)$ . Then  $\lim_{n \rightarrow \infty} c_n = r$  and  $|r - c_n| \leq 2^{-(n+1)}(b_0 - a_0)$ .

*Proof:*

(1) Notice that  $a_0 \leq a_1 \leq a_2 \leq \dots \leq b_0$  and  $b_0 \geq b_1 \geq b_2 \geq \dots \geq a_0$ .

$\therefore \{a_n\}$  is monotonically nondecreasing (*i.e., increasing, but may not be strictly increasing*) and bounded above by  $b_0 \quad \therefore \lim_{n \rightarrow \infty} a_n$  exists

$\therefore \{b_n\}$  is monotonically nonincreasing (*i.e., decreasing, but may not be strictly decreasing*) and bounded below by  $a_0 \quad \therefore \lim_{n \rightarrow \infty} b_n$  exists

$\therefore b_{n+1} - a_{n+1} = \frac{1}{2}(b_n - a_n) \quad \forall n \geq 0 \quad \therefore b_n - a_n = 2^{-n}(b_0 - a_0)$

$\therefore \lim_{n \rightarrow \infty} b_n - \lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} (b_n - a_n) = (b_0 - a_0) \lim_{n \rightarrow \infty} 2^{-n} = 0$

$\therefore \lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} b_n$ , say  $\lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} b_n = r$ .

## Proof of the theorem

---

(2)

$\because f(x)$  is continuous

$$\therefore \lim_{n \rightarrow \infty} f(a_n) = f(\lim_{n \rightarrow \infty} a_n) = f(r) \text{ and } \lim_{n \rightarrow \infty} f(b_n) = f(\lim_{n \rightarrow \infty} b_n) = f(r)$$

$$\therefore f(a_n)f(b_n) < 0$$

$$\therefore 0 \geq \lim_{n \rightarrow \infty} f(a_n)f(b_n) = f(r)f(r)$$

$$\therefore f(r) = 0$$

(3)

$$\therefore r \in [a_n, b_n] \text{ and } c_n = \frac{1}{2}(a_n + b_n) = a_n + \frac{1}{2}(b_n - a_n)$$

$$\therefore |r - c_n| \leq \frac{1}{2}(b_n - a_n) = 2^{-(n+1)}(b_0 - a_0) \quad \square$$

**Note:** Is it true that  $|c_0 - r| \geq |c_1 - r| \geq |c_2 - r| \geq \dots$ ?

Answer: No!  $\Rightarrow$  *not linear convergence!*

**linear:** if  $\exists 0 < C < 1$  and  $\exists n_0 \in \mathbb{N}$  s.t.  $|x_{n+1} - x^*| \leq C|x_n - x^*|, \forall n \geq n_0$ .

## An example

---

If we start with the initial interval  $[50, 63]$ , how many steps do we need in order to have a relative accuracy less than or equal to  $10^{-12}$ ?

This is what we want

$$\frac{|r - c_n|}{|r|} \leq 10^{-12}.$$

Since we know  $r \geq 50$ , thus it is sufficient to have

$$\frac{|r - c_n|}{50} \leq 10^{-12}.$$

Using the above estimate, all we need is

$$2^{-(n+1)} \frac{63 - 50}{50} \leq 10^{-12}.$$

*That means  $n \geq 37$ .*

## Some major problems with the bisection method

---

- Finding the initial interval is not easy.
- Often slow.
- Doesn't work for higher dimensional problems:  $F(X) = \mathbf{0}$ .

## Newton's method

---

- **Motivation:** we know how to solve  $f(x) = 0$  if  $f$  is linear. For nonlinear  $f$ , we can always approximate it with a linear function.
- Let  $x^*$  be a root of  $f(x) = 0$  and  $x$  an approximation of  $x^*$ . Let  $x^* = x + h$ . Using Taylor's expansion, we have

$$0 = f(x^*) = f(x + h) = f(x) + hf'(x) + O(h^2).$$

If  $h$  is small, then we can drop the  $O(h^2)$  term,  $0 \approx f(x) + hf'(x)$ , which means

$$h \approx -\frac{f(x)}{f'(x)}, \quad \text{provided } f'(x) \neq 0.$$

Thus, if  $x$  is an approximation of  $x^* = x + h$ , then

$$x^* = x + h \approx x - \frac{f(x)}{f'(x)}, \quad \text{provided } f'(x) \neq 0.$$

- Newton's method can be defined as follows: for  $n = 0, 1, \dots$

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}, \quad \text{provided } f'(x_n) \neq 0.$$

## An example

---

Find the root of  $f(x) = e^x - 1.5 - \tan^{-1}(x)$ .

Note that  $f(0) = -0.5$ ,  $\lim_{x \rightarrow \infty} f(x) = \infty$ , and  $\lim_{x \rightarrow -\infty} f(x) > 0.07$ .

Therefore,  $\exists c^+ \in (0, \infty)$  and  $c^- \in (-\infty, 0)$  are zeros of  $f$ .

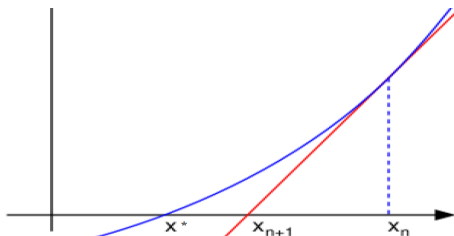
Suppose we start with  $x_0 = -7.0$ , then the results of Newton iterations are

$$\begin{aligned}x_0 &= -7.0, & f(x_0) &= -0.7 \times 10^{-1} \\x_1 &= -10.7, & f(x_1) &= -0.2 \times 10^{-1} \\x_3 &= -14.0, & f(x_3) &= -0.2 \times 10^{-3} \\x_4 &= -14.1, & f(x_4) &= -0.8 \times 10^{-6}\end{aligned}$$

The output shows rapid convergence of the iterations.

## Geometrical interpretation

---



- This is an illustration of one iteration of Newton's method. The function  $f$  is shown in blue and the tangent line is in red. We see that  $x_{n+1}$  is a better approximation than  $x_n$  for the root  $x^*$  of the function  $f$ .
- What is the geometrical meaning of  $f'(x_n) = 0$ ?



## Some stopping criteria

---

- Using the residual information  $f(x_n)$ :
  - (1) if  $|f(x_n)| < \varepsilon$  then stop (absolute residual criterion).
  - (2) if  $|f(x_n)| < \varepsilon|f(x_0)|$  then stop (relative residual criterion).
- Using the step size information  $|x_{n+1} - x_n|$ :
  - (1) if  $|x_{n+1} - x_n| < \delta$  then stop (approximate absolute error criterion).
  - (2) if  $\frac{|x_{n+1} - x_n|}{|x_{n+1}|} < \delta$  then stop (approximate relative error criterion).
- Maximum number of iterations  $M$ .

## Newton's algorithm including stopping criteria

---

```
input  $x_0, M, \varepsilon, \delta$   
 $v \leftarrow f(x_0)$   
if  $|v| < \varepsilon$  then stop  
for  $k = 1$  to  $M$  do  
     $x_1 = x_0 - v/f'(x_0)$   
     $v \leftarrow f(x_1)$   
    if  $|x_1 - x_0| < \delta$  or  $|v| < \varepsilon$  then stop  
     $x_0 \leftarrow x_1$   
end do
```

---

See the details of the M-file `newton.m` for  $f(x) = e^x - \sin(x)$

---

**Note:** if  $f'(x_0)$  is too small, then  $1/f'(x_0)$  may overflow.

## Convergence analysis

Assume that  $f''$  is continuous and  $x^*$  is a simple zero of  $f$ , i.e.,  $f(x^*) = 0$  and  $f'(x^*) \neq 0$ . Define the error as  $e_n = x_n - x^*$ . Then

$$\begin{aligned}e_{n+1} &= x_{n+1} - x^* = x_n - \frac{f(x_n)}{f'(x_n)} - x^* \\ &= e_n - \frac{f(x_n)}{f'(x_n)} = \frac{e_n f'(x_n) - f(x_n)}{f'(x_n)}.\end{aligned}$$

Using Taylor's expansion,

$$0 = f(x^*) = f(x_n - e_n) = f(x_n) - e_n f'(x_n) + \frac{1}{2} e_n^2 f''(\xi_n),$$

for some  $\xi_n$  between  $x_n$  and  $x^*$ . Therefore, we have

$$(\star) \quad e_{n+1} = \frac{1}{2} \frac{f''(\xi_n)}{f'(x_n)} e_n^2 \left( \approx \frac{1}{2} \frac{f''(x^*)}{f'(x^*)} e_n^2 := C e_n^2, \text{ provided } x_n \approx x^* \right).$$

Define a quantity  $c_\delta$  for  $\delta > 0$  by

$$c_\delta := \frac{1}{2} \left( \max_{|x-x^*| \leq \delta} |f''(x)| \right) / \left( \min_{|x-x^*| \leq \delta} |f'(x)| \right) \geq 0.$$

We can select  $\delta > 0$  such that  $\rho := \delta c_\delta < 1$ . (why?)

## Theorem on Newton's method

---

Assume that  $|e_0| = |x_0 - x^*| < \delta$ . Then  $|\xi_0 - x^*| < \delta$  and we have  $\frac{1}{2}|f''(\xi_0)/f'(x_0)| \leq c_\delta$ . Therefore,

$$|x_1 - x^*| = |e_1| \leq e_0^2 c_\delta = |e_0||e_0|c_\delta < |e_0|\delta c_\delta = |e_0|\rho < |e_0| < \delta.$$

Repeating this argument, we have

$$|e_1| < \rho|e_0|, |e_2| < \rho|e_1| < \rho^2|e_0|, \dots, |e_n| < \rho^n|e_0|.$$

Since  $0 \leq \rho < 1$ , we have  $\lim_{n \rightarrow \infty} \rho^n = 0$  which implies that  $\lim_{n \rightarrow \infty} e_n = 0$ .

Finally, since  $|e_n| = |x_n - x^*| < \delta$  and  $|\xi_n - x^*| < \delta$ , we have from (\*) that

$$|e_{n+1}| = \frac{1}{2} \frac{|f''(\xi_n)|}{|f'(x_n)|} |e_n|^2 \leq \frac{1}{2} c_\delta |e_n|^2 \leq \frac{1}{2} (c_\delta + 1) |e_n|^2 := C |e_n|^2,$$

which implies the quadratic convergence.  $\square$

## Theorem on Newton's method

---

**Theorem on Newton's method:** *Let  $f''$  be continuous and let  $x^*$  be a simple zero of  $f$ . Then there exist  $\delta > 0$  and  $C > 0$  such that if the initial guess  $x_0 \in N(x^*, \delta)$  (i.e.,  $|x_0 - x^*| < \delta$ ) then Newton's method converges and satisfies*

$$|x_{n+1} - x^*| \leq C|x_n - x^*|^2 \quad (\forall n \geq 0).$$

**Good:** *the convergence is quadratic.*

**Bad:** *the initial guess  $x_0$  has to be close to the solution  $x^*$ .*

## Example

---

Find the root of  $f(x) = \alpha - 1/x$ , for any given  $\alpha > 0$  (we know the exact solution is  $x^* = 1/\alpha$ ). Using Newton's method, we have

$$x_{n+1} = x_n - \frac{\alpha - \frac{1}{x_n}}{1/x_n^2},$$

which is same as

$$x_{n+1} = 2x_n - \alpha x_n^2, \quad n = 0, 1, 2, \dots$$

### Questions:

- Does the sequence  $x_0, x_1, x_2, \dots$  converge? ( $\iff 0 < x_0 < \frac{2}{\alpha}$ )
- How fast? (quadratic)
- Does the convergence depend on the initial guess  $x_0$ ? (Yes)

## Example (cont'd)

Let us define the error  $e_n = x^* - x_n = \frac{1}{\alpha} - x_n$ . Then

$$e_{n+1} = \frac{1}{\alpha} - x_{n+1} = \frac{1}{\alpha} - 2x_n + \alpha x_n^2 = \alpha \left( \frac{1}{\alpha} - x_n \right)^2 = \alpha e_n^2.$$

Thus, if it converges, then the rate is quadratic. We now have

$$\begin{aligned} e_{n+1} &= \alpha e_n^2 = \alpha (\alpha e_{n-1}^2)^2 = \alpha^3 (e_{n-1}^2)^2 = \frac{1}{\alpha} (\alpha^2 e_{n-1}^2)^2 = \frac{1}{\alpha} (\alpha e_{n-1})^2 \\ &= \frac{1}{\alpha} (\alpha \alpha e_{n-2}^2)^2 = \frac{1}{\alpha} (\alpha^2 e_{n-2}^2)^2 = \frac{1}{\alpha} (\alpha e_{n-2})^2 = \dots = \frac{1}{\alpha} (\alpha e_0)^{2^{n+1}}, \end{aligned}$$

which implies that

$$\begin{aligned} x_n \text{ converges to } x^* &\iff \lim_{n \rightarrow \infty} e_n = 0 \iff |\alpha e_0| < 1 \iff |e_0| < \frac{1}{\alpha} \\ &\iff \left| \frac{1}{\alpha} - x_0 \right| < \frac{1}{\alpha} \iff -\frac{1}{\alpha} < \frac{1}{\alpha} - x_0 < \frac{1}{\alpha} \\ &\iff 0 < x_0 < \frac{2}{\alpha}. \end{aligned}$$

## Some remarks on Newton's method

---

### Advantages:

- The convergence is **quadratic**.
- Newton's method works for higher dimensional problems.

### Disadvantages:

- Newton's method converges only **locally**; i.e., the initial guess  $x_0$  has to be close enough to the solution  $x^*$ .
- It needs the first derivative of  $f(x)$ .



## Newton's method for systems of nonlinear equations

- We wish to solve

$$\begin{cases} f_1(x_1, x_2) = 0, \\ f_2(x_1, x_2) = 0, \end{cases}$$

where  $f_1$  and  $f_2$  are nonlinear functions of  $x_1$  and  $x_2$ .

- Assume that  $(x_1 + h_1, x_2 + h_2)$  is a solution of the nonlinear system of equations. Applying Taylor's expansion in two variables around  $(x_1, x_2)$ , we obtain

$$\begin{cases} 0 = f_1(x_1 + h_1, x_2 + h_2) \approx f_1(x_1, x_2) + h_1 \frac{\partial f_1(x_1, x_2)}{\partial x_1} + h_2 \frac{\partial f_1(x_1, x_2)}{\partial x_2}, \\ 0 = f_2(x_1 + h_1, x_2 + h_2) \approx f_2(x_1, x_2) + h_1 \frac{\partial f_2(x_1, x_2)}{\partial x_1} + h_2 \frac{\partial f_2(x_1, x_2)}{\partial x_2}. \end{cases}$$

- Putting it into the matrix form, we have

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} \approx \begin{bmatrix} f_1(x_1, x_2) \\ f_2(x_1, x_2) \end{bmatrix} + \begin{bmatrix} \frac{\partial f_1(x_1, x_2)}{\partial x_1} & \frac{\partial f_1(x_1, x_2)}{\partial x_2} \\ \frac{\partial f_2(x_1, x_2)}{\partial x_1} & \frac{\partial f_2(x_1, x_2)}{\partial x_2} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \end{bmatrix}.$$

## Newton's method for systems of nonlinear equations (cont'd)

---

- To simplify the notation we introduce the **Jacobian matrix**:

$$J(x_1, x_2) = \begin{bmatrix} \frac{\partial f_1(x_1, x_2)}{\partial x_1} & \frac{\partial f_1(x_1, x_2)}{\partial x_2} \\ \frac{\partial f_2(x_1, x_2)}{\partial x_1} & \frac{\partial f_2(x_1, x_2)}{\partial x_2} \end{bmatrix}.$$

- Then we have

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} \approx \begin{bmatrix} f_1(x_1, x_2) \\ f_2(x_1, x_2) \end{bmatrix} + J(x_1, x_2) \begin{bmatrix} h_1 \\ h_2 \end{bmatrix}.$$

- If  $J(x_1, x_2)$  is nonsingular then we can solve for  $[h_1, h_2]^T$ :

$$J(x_1, x_2) \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} \approx - \begin{bmatrix} f_1(x_1, x_2) \\ f_2(x_1, x_2) \end{bmatrix}.$$

## Newton's method for systems of nonlinear equations (cont'd)

- Newton's method for the system of nonlinear equations is defined as follows: for  $k = 0, 1, \dots$ ,

$$\begin{bmatrix} x_1^{(k+1)} \\ x_2^{(k+1)} \end{bmatrix} = \begin{bmatrix} x_1^{(k)} \\ x_2^{(k)} \end{bmatrix} + \begin{bmatrix} h_1^{(k)} \\ h_2^{(k)} \end{bmatrix}$$

with

$$J(x_1^{(k)}, x_2^{(k)}) \begin{bmatrix} h_1^{(k)} \\ h_2^{(k)} \end{bmatrix} = - \begin{bmatrix} f_1(x_1^{(k)}, x_2^{(k)}) \\ f_2(x_1^{(k)}, x_2^{(k)}) \end{bmatrix}.$$

- Exercise:**

Solve the following nonlinear system by using Newton's method with the initial guess  $\mathbf{x}^{(0)} = (x_1^{(0)}, x_2^{(0)})^\top = (0, 1)^\top$ . Perform two iterations.

$$\begin{cases} 4x_1^2 - x_2^2 = 0, \\ 4x_1x_2^2 - x_1 = 1. \end{cases}$$

## Newton's method for higher dimensional problems

- In general, we can use Newton's method for  $F(X) = \mathbf{0}$ , where  $X = (x_1, x_2, \dots, x_n)^\top$  and  $F = (f_1, f_2, \dots, f_n)^\top$ .
- For higher dimensional problem, the first derivative is defined as a matrix (the Jacobian matrix)

$$DF(X) := \begin{bmatrix} \frac{\partial f_1(X)}{\partial x_1} & \frac{\partial f_1(X)}{\partial x_2} & \cdots & \frac{\partial f_1(X)}{\partial x_n} \\ \frac{\partial f_2(X)}{\partial x_1} & \frac{\partial f_2(X)}{\partial x_2} & \cdots & \frac{\partial f_2(X)}{\partial x_n} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial f_n(X)}{\partial x_1} & \frac{\partial f_n(X)}{\partial x_2} & \cdots & \frac{\partial f_n(X)}{\partial x_n} \end{bmatrix}.$$

- Newton's method: given  $X^{(0)} = [x_1^{(0)}, \dots, x_n^{(0)}]^\top$ , define

$$X^{(k+1)} = X^{(k)} + H^{(k)},$$

where

$$DF(X^{(k)})H^{(k)} = -F(X^{(k)}),$$

which requires the solving of a large linear system of equations at every iteration.

## Operations involved in Newton's method

---

- vector operations: not expensive.
- function evaluations: can be expensive.
- compute the Jacobian: can be expensive.
- solving matrix equations (linear system): very expensive – **topic of the next chapter!**

## Methods without using derivatives

---

- “Finite difference Newton’s method” and “secant method.”
- **Basic idea:**

$$x \leftarrow x - \frac{f(x)}{f'(x)}.$$

If  $f'(x)$  is too hard or too expensive to compute, we can use an approximation.

- **Questions:** how to obtain an approximation? Do we lose the fast convergence?

## Finite difference Newton's method

---

- Let  $h$  be a small nonzero parameter, then

$$a := \frac{f(x_n + h) - f(x_n)}{h}$$

can be a good approximation of  $f'(x_n)$ .

- FD-Newton's method:**

(1) compute  $a = \frac{f(x_n + h) - f(x_n)}{h}$ .

(2) compute  $x_{n+1} = x_n - \frac{f(x_n)}{a}$ .

- Remarks:**

- (1) the method needs **an extra parameter  $h$** . What shall we use?
- (2) the method needs **two function evaluations** per iteration.
- (3) what is the convergence rate?

## Secant method

---

- Since  $h$  can be any small number in the FD-Newton's method, why don't we simply use  $h = x_n - x_{n-1}$ , which may be positive or negative, but usually not zero.

- **Secant method:**

(1) compute  $a = \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}$ .

(2) compute  $x_{n+1} = x_n - \frac{f(x_n)}{a}$ .

- **Remarks:**

- (1) now we need **only one function evaluation** per iteration.
- (2)  $x_{n+1}$  depends on two previous iterations. For example, to compute  $x_2$ , we need both  $x_1$  and  $x_0$ .
- (3) how do we obtain  $x_1$ ? We need to use FD-Newton: pick a small parameter  $h$ , compute  $a_0 = (f(x_0 + h) - f(x_0))/h$ , then  $x_1 = x_0 - f(x_0)/a_0$ .



## Which of the three methods is better?

An example:  $f(x) = x^2 - 1$ , and we take  $x_0 = 2.0$ .

Stopping parameters:  $\delta = 10^{-10}$ ,  $\varepsilon = 10^{-10}$ .

$h = 10^{-7}$  in FD-Newton method.

Iter.	Newton	FD-Newton	Secant
$x_0$	2.0	2.0	2.0
$x_1$	1.2500000000000000	1.25000001709125	1.25000001709125
$x_2$	1.0250000000000000	1.02500001222170	1.07692308177740
$x_3$	1.00030487804878	1.00030487955710	1.00826446381851
$x_4$	1.00000004646115	1.00000004647732	1.00030487810437
$x_5$	1.0000000000000000	1.0000000000000000	1.00000125445212
$x_6$			1.00000000019120
$x_7$			1.0000000000000000

See the details of the M-files: [comparisonnewton.m](#),  
[comparisonFDnewton.m](#), [comparisonsecant.m](#)

## Convergence rates

---

- If  $|h_n| \leq C|x_n - x^*|$ , then the convergence of FD-Newton is **quadratic**.
- *The convergence of secant method is superlinear (i.e., better than linear).* More precisely, we have (see Textbook, pp. 96-97)

$$|e_{n+1}| \leq C|e_n|^{(1+\sqrt{5})/2}, \quad (1 + \sqrt{5})/2 \approx 1.62 < 2.$$

- **Remark:** when selecting algorithms for a particular problem, one should consider not only the rate (order) of convergence, but also the cost of computing  $f(x_n)$  and  $f'(x_n)$ .

## An informal convergence analysis of the secant method

---

Let  $e_n := x_n - x^*$ . Under suitable assumptions, it can be shown that  $e_{n+1} \approx C e_n e_{n-1}$  (Textbook, p. 96) and  $\lim_{n \rightarrow \infty} e_n = 0$  (cf. analysis for Newton's method).

To discover the order of convergence, we assume that for large  $n$ ,  $|e_{n+1}| \approx \lambda |e_n|^\alpha$ . Thus,  $|e_n| \approx \lambda |e_{n-1}|^\alpha \Rightarrow |e_{n-1}| \approx \lambda^{-1/\alpha} |e_n|^{1/\alpha}$ .

$$\therefore \lambda |e_n|^\alpha \approx |e_{n+1}| \approx |C| |e_n| \lambda^{-1/\alpha} |e_n|^{1/\alpha}$$

$$\therefore |e_n|^\alpha \approx |C| \lambda^{-1/\alpha - 1} |e_n|^{1+1/\alpha}$$

$$\therefore |e_n|^{\alpha - 1 - 1/\alpha} \approx |C| \lambda^{-1/\alpha - 1}$$

$\therefore$  the right side of this relation is a nonzero constant while  $e_n \rightarrow 0$

$$\therefore \alpha - 1 - 1/\alpha = 0$$

$$\therefore \alpha^2 - \alpha - 1 = 0$$

$$\therefore \alpha = \frac{1 + \sqrt{5}}{2} \approx 1.62 > 0 \quad \square$$

## Steffensen's method – method without using derivative

**Steffensen's method:**

$$x_{n+1} = x_n - \frac{f(x_n)}{g(x_n)}, \text{ where } g(x_n) := \frac{f(x_n + f(x_n)) - f(x_n)}{f(x_n)}.$$

Under suitable hypotheses, the method is **quadratically** convergent (p. 90, # 4).

**An informal convergence analysis:** Assume that  $f \in C^2$ . By Taylor expansion, we have

$$f(x + f(x)) = f(x) + f(x)f'(x) + \frac{f(x)^2}{2}f''(\xi),$$

for some  $\xi$  between  $x$  and  $x + f(x)$ . Therefore,

$$g(x) := \frac{1}{f(x)} \{f(x + f(x)) - f(x)\} = f'(x) + \frac{f(x)}{2}f''(\xi) \approx f'(x), \text{ if } f(x) \approx 0.$$

$$\text{Let } e_n := x_n - x^*. \text{ Then, } e_{n+1} = e_n - \frac{f(x_n)}{g(x_n)} = \frac{1}{g(x_n)} \{e_n g(x_n) - f(x_n)\}.$$

## Steffensen's method (cont'd)

---

$$\therefore 0 = f(x^*) = f(x_n - e_n) = f(x_n) - e_n f'(x_n) + \frac{e_n^2}{2} f''(\xi_n),$$

for some  $\xi_n$  between  $x_n$  and  $x_n - e_n$

$$\therefore f(x_n) - e_n g(x_n) \approx -\frac{e_n^2}{2} f''(\xi_n)$$

$$\therefore e_{n+1} \approx \frac{e_n^2 f''(\xi_n)}{2 g(x_n)} \left( \approx \frac{f''(x^*)}{2 f'(x^*)} e_n^2, \text{ provided } x_n \approx x^* \right)$$

(cf. analysis of Newton's method).  $\square$

### Remarks:

- Bisection algorithms is **global**, and all the other Newton-type algorithms are **local**.
- Local algorithms are often **fast**, and global algorithms are often **slow**.

## Fixed points

---

- A function  $F : x \mapsto F(x)$  is often called a mapping from  $x$  to  $F(x)$  ( $F$  takes an input value  $x$  and generates an output value  $F(x)$ ).

*If there is a point  $p$ , at which the output is the same as the input, then that point is called a fixed point of  $F$ , i.e.,  $p = F(p)$ .*

- Finding the fixed points of  $F$  has many applications. For example, if

$$F(x) := x - \frac{f(x)}{f'(x)},$$

then the fixed point of  $F$  is simply the root of  $f(x) = 0$ .

*“root-finding problem”  $\implies$  “fixed point problem”*

## Fixed point iterations

---

- Fixed point iterations:

$$x_{n+1} = F(x_n), \quad n = 0, 1, \dots$$

Assume that  $F$  is continuous and  $\lim_{n \rightarrow \infty} x_n = p$ . Then

$$F(p) = F(\lim_{n \rightarrow \infty} x_n) = \lim_{n \rightarrow \infty} F(x_n) = \lim_{n \rightarrow \infty} x_{n+1} = p.$$

Therefore,  $p$  is a fixed point of the function  $F$ .

- The following three fixed point iterations can be considered for solving  $x^3 - x - 5 = 0$ :

$x_{n+1} = F(x_n)$ ,  $n = 0, 1, \dots$ , where

(1)  $F(x) = x^3 - 5$ .

(2)  $F(x) = (x + 5)^{1/3}$ .

(3)  $F(x) = \frac{5}{x^2 - 1}$ .

Do the iterations converge?

## A fixed point theorem

- If  $F \in C[a, b]$  and  $F(x) \in [a, b], \forall x \in [a, b]$ , then  $F$  has a fixed point in  $[a, b]$ .
- If, in addition,  $F'$  exists on  $(a, b)$  and  $\exists 0 < k < 1$  such that  $|F'(x)| \leq k, \forall x \in (a, b)$ , then the fixed point is unique in  $[a, b]$ .
- Then, for any  $x_0 \in [a, b]$  and  $x_{n+1} := F(x_n), n \geq 0$ , the sequence converges to the unique fixed point  $p \in [a, b]$  and
  - (1)  $|x_n - p| \leq k^n \max\{x_0 - a, b - x_0\}, \forall n \geq 1;$
  - (2)  $|x_n - p| \leq \frac{k^n}{1-k} |x_1 - x_0|, \forall n \geq 1.$

*Proof.*

- If  $F(a) = a$  or  $F(b) = b$  then  $F$  has a fixed point in  $[a, b]$ . Suppose not, then  $a < F(a) \leq b$  and  $a \leq F(b) < b$ . Define  $H(x) := F(x) - x$ . Then  $H$  is continuous on  $[a, b]$  and  $H(a) > 0, H(b) < 0$ . By the Intermediate Value Theorem,  $\exists p \in (a, b)$  such that  $H(p) = 0$ , i.e.,  $F(p) = p$ .  $\square$
- Suppose that  $\exists p < q \in [a, b]$  are fixed points of  $F$ . Then  $F(p) = p$  and  $F(q) = q$ . By the Mean Value Theorem,  $\exists \xi \in (p, q)$  such that  $\frac{F(q) - F(p)}{q - p} = F'(\xi) \implies \frac{|F(q) - F(p)|}{|q - p|} = |F'(\xi)| \leq k < 1 \implies 1 = \frac{|q - p|}{|q - p|} \leq k < 1$ . This is a contradiction. Therefore, the fixed point is unique.  $\square$



## Proof of the fixed point theorem (cont'd)

- For  $n \geq 1$ , by the Mean Value Theorem,  $\exists \zeta \in (a, b)$  such that
$$0 \leq |x_n - p| = |F(x_{n-1}) - F(p)| = |F'(\zeta)| |x_{n-1} - p| \leq k |x_{n-1} - p|.$$
$$\implies 0 \leq |x_n - p| \leq k |x_{n-1} - p| \leq k^2 |x_{n-2} - p| \leq \cdots \leq k^n |x_0 - p|.$$
$$\implies \lim_{n \rightarrow \infty} |x_n - p| = 0 \Leftrightarrow \lim_{n \rightarrow \infty} x_n - p = 0 \Leftrightarrow \lim_{n \rightarrow \infty} x_n = p.$$

(1)  $\because |x_n - p| \leq k^n |x_0 - p|$  and  $p \in [a, b]$   
 $\therefore |x_n - p| \leq k^n \max\{x_0 - a, b - x_0\}, \forall n \geq 1$

(2) For  $n \geq 1$ ,  
 $|x_{n+1} - x_n| = |F(x_n) - F(x_{n-1})| \leq k |x_n - x_{n-1}| \leq \cdots \leq k^n |x_1 - x_0|.$   
 $\therefore$  For  $m > n \geq 1$ , we have

$$\begin{aligned} |x_m - x_n| &= |x_m - x_{m-1} + x_{m-1} - x_{m-2} + \cdots + x_{n+1} - x_n| \\ &\leq |x_m - x_{m-1}| + |x_{m-1} - x_{m-2}| + \cdots + |x_{n+1} - x_n| \\ &\leq k^{m-1} |x_1 - x_0| + k^{m-2} |x_1 - x_0| + \cdots + k^n |x_1 - x_0| \\ &= k^n (1 + k + \cdots + k^{m-n-1}) |x_1 - x_0|. \end{aligned}$$

$$\therefore \lim_{n \rightarrow \infty} x_n = p$$

$$\therefore |p - x_n| = \lim_{m \rightarrow \infty} |x_m - x_n| \leq k^n |x_1 - x_0| \sum_{i=0}^{\infty} k^i = k^n |x_1 - x_0| \frac{1}{1-k}$$

( $\because$  geometric series with  $0 < k < 1$ )

$$\therefore |p - x_n| \leq \frac{k^n}{1-k} |x_1 - x_0| \quad \square$$

## Contractive mappings

---

- **Definition:** A mapping (function)  $F$  is said to be contractive if  $\exists 0 < \lambda < 1$  such that  $|F(x) - F(y)| \leq \lambda|x - y|$ , for all  $x, y$  in the domain of  $F$ .
- **Note:** In the above theorem,  $F$  is contractive on  $[a, b]$ .
- **Example:**  $F(x) = 4 + \frac{1}{3} \sin(2x)$  is contractive on  $\mathbb{R}$ .

$$\begin{aligned}|F(x) - F(y)| &= \frac{1}{3} |\sin(2x) - \sin(2y)| \\ &= \frac{2}{3} |\cos(2\xi)| |x - y| \\ &\leq \frac{2}{3} |x - y|.\end{aligned}$$

## Contraction mapping principle

---

Let  $F$  be a contractive mapping from a complete metric space  $X \subseteq \mathbb{R}$  into itself. Then  $F$  has a unique fixed point  $p$  and the sequence  $\{x_n\}$  generated by  $x_{n+1} := F(x_n)$ ,  $n \geq 0$ , converges to  $p$  for any  $x_0 \in X$ .

*Proof:*

- show that  $\{x_n\}$  converges;
- let  $\lim_{n \rightarrow \infty} x_n = p$ . Then  $F(p) = p$ ;
- show that  $p$  is unique.  $\square$

**Note:** Let  $X$  be a closed subset of  $\mathbb{R}$ . Then  $X$  is a complete metric space.

**Example:** closed subsets of  $\mathbb{R}$ :  $[a, b]$ ,  $\mathbb{R}$ , etc.

## Error analysis

---

- Assume that  $F'$  exists and continuous. Consider the fixed point iterations,

$$x_{n+1} = F(x_n), \quad n \geq 0.$$

Assume that  $\{x_n\}$  converges to  $p$  ( $p$  is a fixed point). Let  $e_n := x_n - p$ . Then, by MVT, we have

$$e_{n+1} = x_{n+1} - p = F(x_n) - F(p) = F'(\xi_n)(x_n - p) = F'(\xi_n)e_n,$$

for some  $\xi_n$  between  $x_n$  and  $p$ . The condition  $|F'(x)| < 1$  for all  $x$  ensures that the errors decrease in magnitude. If  $e_n$  is small then  $\xi_n$  is near  $p$ , and  $F'(\xi_n) \approx F'(p)$ .

- One would expect rapid convergence if  $F'(p)$  is small. **Ideally,**  $F'(p) = 0$ .

## Error analysis (cont'd)

- Assume that  $F^{(k)}(p) = 0$  for  $1 \leq k < r$  but  $F^{(r)}(p) \neq 0$ . Then

$$\begin{aligned}e_{n+1} &= x_{n+1} - p = F(x_n) - F(p) = F(p + e_n) - F(p) \\&= \left\{ F(p) + e_n F'(p) + \frac{e_n^2}{2} F''(p) + \cdots + \frac{1}{r!} e_n^r F^{(r)}(\xi_n) \right\} - F(p) \\&= e_n F'(p) + \frac{e_n^2}{2} F''(p) + \cdots + \frac{e_n^{r-1}}{(r-1)!} F^{(r-1)}(p) + \frac{e_n^r}{r!} F^{(r)}(\xi_n) \\&= \frac{e_n^r}{r!} F^{(r)}(\xi_n).\end{aligned}$$

- If we know that the method converges and  $F^{(r)}$  is continuous then

$$\lim_{n \rightarrow \infty} \frac{|e_{n+1}|}{|e_n|^r} = \frac{1}{r!} |F^{(r)}(p)|$$

and the method converges with order  $r$ .

## Newton' method

---

**Newton' method:**  $F(x) = x - \frac{f(x)}{f'(x)}$ ,  $f(p) = 0$  and  $f'(p) \neq 0$ ,  $F(p) = p$ .

$$\therefore F'(x) = 1 - \frac{f'(x)f'(x) - f(x)f''(x)}{(f'(x))^2} = \frac{f(x)f''(x)}{(f'(x))^2}$$

$$\therefore F'(p) = 0$$

$\therefore$

$$F''(x) = \frac{(f'(x))^2 \{f(x)f'''(x) + f''(x)f'(x)\} - (f(x)f''(x))(2f'(x)f''(x))}{(f'(x))^4}$$

$$\therefore \text{we usually have } F''(p) = \frac{f''(p)}{f'(p)} \neq 0$$

$\therefore$  under suitable assumptions,  
the order (rate) of convergence of Newton's method is 2

## Roots of polynomials

---

- A general polynomial:  
 $p(z) = a_n z^n + a_{n-1} z^{n-1} + \cdots + a_2 z^2 + a_1 z^1 + a_0$ , where coefficients  $a_i \in \mathbb{C}$ ,  $i = 0, 1, \dots, n$ . If  $a_n \neq 0$  then we say *degree*( $p$ ) =  $n$ .
- **Fundamental Theorem of Algebra:** *Every nonconstant polynomial has at least one root in  $\mathbb{C}$ .*  
( $\iff$  A polynomial of degree  $n$  has exactly  $n$  roots in  $\mathbb{C}$ ).
- If  $p$  is a polynomial whose coefficients are all real,  $a_i \in \mathbb{R} \forall i$ , then its roots may be complex and if  $w = w_1 + iw_2$  is a complex root then its conjugate  $\bar{w} := w_1 - iw_2$  is also a root.

*In what follows, we consider polynomials with real coefficients.*

## Horner's algorithm

**Newton's method:**  $z_{k+1} := z_k - \frac{p(z_k)}{p'(z_k)}, k = 0, 1, 2, \dots$

We need function evaluations  $p(z_k)$  and  $p'(z_k)$  in Newton's method.

- Given a polynomial  $p(z) = a_n z^n + a_{n-1} z^{n-1} + \dots + a_1 z^1 + a_0$  and  $z_0 \in \mathbb{R}$ . **Horner's algorithm** will produce the number  $p(z_0)$  and the polynomial  $q(z)$  such that  $p(z) = (z - z_0)q(z) + p(z_0)$ .
- Assume that  $q(z) = b_{n-1} z^{n-1} + b_{n-2} z^{n-2} + \dots + b_1 z^1 + b_0$ . Then we have  $b_{n-1} = a_n, b_{n-2} = a_{n-1} + z_0 b_{n-1}, \dots, b_0 = a_1 + z_0 b_1, p(z_0) = p(z) - (z - z_0)q(z) = a_0 + z_0 b_0$ .
- Synthetic division:** (綜合除法)

$$\begin{array}{r|cccccc} & a_n & a_{n-1} & a_{n-2} & \cdots & a_0 & \\ z_0 & & z_0 b_{n-1} & z_0 b_{n-2} & \cdots & z_0 b_0 & \\ \hline & b_{n-1} & b_{n-2} & b_{n-3} & \cdots & b_{-1} & \leftarrow p(z_0) \end{array}$$

We have  $p(z_0) = b_{-1}$ .



## Example

---

Let  $p(z) = z^4 - 4z^3 + 7z^2 - 5z - 2$ . Evaluate  $p(3)$ .

$$\begin{array}{r|rrrrr} 3 & 1 & -4 & 7 & -5 & -2 \\ & & 3 & -3 & 12 & 21 \\ \hline & 1 & -1 & 4 & 7 & 19 & \leftarrow p(3) \end{array}$$

$\therefore p(3) = 19, q(z) = z^3 - z^2 + 4z + 7$ , and

$$p(z) = (z - 3)(z^3 - z^2 + 4z + 7) + 19.$$

## Complete Horner's algorithm

---

Given  $p(z) = a_n z^n + a_{n-1} z^{n-1} + \cdots + a_2 z^2 + a_1 z^1 + a_0$  and  $z_0 \in \mathbb{R}$ .

We wish to find  $c_i, i = 0, 1, \cdots, n$  such that

$$p(z) = c_n (z - z_0)^n + c_{n-1} (z - z_0)^{n-1} + \cdots + c_1 (z - z_0)^1 + c_0.$$

If so, by Taylor Theorem, we know that  $c_k = \frac{p^{(k)}(z_0)}{k!}$ .

$\therefore p(z_0) = c_0$  and  $p'(z_0) = c_1 = q(z_0)$

$\therefore$  We can apply Horner's algorithm again to  $q(z)$  with point  $z_0$

Repeat this process, we can obtain  $c_i, i = 0, 1, \cdots, n$ .

## Example

Let  $p(z) = z^4 - 4z^3 + 7z^2 - 5z - 2$  and  $z_0 = 3$ .

	1	-4	7	-5	-2	
3		3	-3	12	21	
	1	-1	4	7	19	$\leftarrow p(3)$
3		3	6	30		
	1	2	10	37		$\leftarrow p'(3)$
3		3	15			
	1	5	25			
3		3				
	1	8				

$\therefore p(3) = 19, p'(3) = 37$  and

$$p(z) = 1(z-3)^4 + 8(z-3)^3 + 25(z-3)^2 + 37(z-3)^1 + 19$$

## Newton's method with Horner's algorithm

---

```
program horner( $n, (a_i : 0 \leq i \leq n), z_0, \alpha, \beta$ )  
 $\alpha \leftarrow a_n$   
 $\beta \leftarrow 0$   
for  $k = n - 1 : -1 : 0$  do  
     $\beta \leftarrow \alpha + z_0\beta$   
     $\alpha \leftarrow a_k + z_0\alpha$   
end do  
output  $\alpha(= p(z_0)), \beta(= p'(z_0))$ 
```

---

```
program newton ( $n, (a_i : 0 \leq i \leq n), z_0, M, \delta$ )  
for  $k = 1 : 1 : M$  do  
    call horner( $n, (a_i : 0 \leq i \leq n), z_0, \alpha, \beta$ )  
     $z_1 \leftarrow z_0 - \alpha / \beta$   
    output  $\alpha, \beta, z_1$   
    if  $|z_1 - z_0| < \delta$  then stop  
     $z_0 \leftarrow z_1$   
end do
```

## Basic idea of continuation method (延拓法)

---

The basic idea of the continuation method is to embed the given problem in a one-parameter family of problems, using a parameter  $t$  that runs over  $[0, 1]$ , such that for  $t = 1$  we have the original problem, while for  $t = 0$  we have another problem with known solution.

Below is an example:

- Consider a root-finding problem:  $f(x) = 0$ . We extend the problem to a one-parameter family of problems:

$$h(t, x) = tf(x) + (1 - t)g(x),$$

where  $t \in [0, 1]$  and  $g(x)$  is given and have a known zero, say  $x_0$ .

- Select points  $0 = t_0 < t_1 < \cdots < t_{m-1} < t_m = 1$ . We then solve each equation  $h(t_i, x) = 0, i = 0, 1, \cdots, m$ . We say each solution  $x_i, i = 0, 1, \cdots, m$ .
- Assume that some iterative method such as Newton's method is used to solve  $h(t_i, x) = 0$ , we use the solution  $x_{i-1}$  of  $h(t_{i-1}, x) = 0$  as the starting point.

## Homotopy (同倫)

---

**Definition:** Let  $X$  and  $Y$  be two topological spaces and  $f, g : X \rightarrow Y$  be two continuous functions. A homotopy between  $f$  and  $g$  is defined to be a continuous function  $h : [0, 1] \times X \rightarrow Y$  such that, for all points  $x \in X$ ,  $h(0, x) = g(x)$  and  $h(1, x) = f(x)$ . If such a map exists, we say that  $f$  is homotopic to  $g$ .

A simple example that is often used in continuation method is

$$h(t, x) = tf(x) + (1 - t) \underbrace{(f(x) - f(x_0))}_{:=g(x)},$$

where  $x_0$  can be any point in  $X$ .

## Homotopy continuation method

---

- If  $h(t, x) = 0$  has a unique solution for each  $t \in [0, 1]$ , then the solution is a function of  $t$ , and we write  $x(t) \in X$ . The set  $\{x(t) : 0 \leq t \leq 1\}$  can be interpreted as a curve in  $X$ . *The continuation method attempts to determine this curve by computing points on it,  $x(t_0), x(t_1), \dots, x(t_m)$ .*
- **Homotopy continuation method:** Assume that  $x(t)$  and  $h(t, x)$  are differentiable functions. Then

$$0 = h(t, x(t)) \implies 0 = h_t(t, x(t)) + h_x(t, x(t))x'(t)$$
$$\implies x'(t) = -\left(h_x(t, x(t))\right)^{-1} h_t(t, x(t)).$$

This is an **ODE with a known initial value  $x(0)$** , it can be solved using numerical methods (cf. Chapter 8).

- If necessary, we can apply **Newton's iteration** starting at the point produced by the homotopy method to approximate the solution of  $h(1, x) = 0$  one more time.

## Example

---

Let  $X = Y = \mathbb{R}^2$  and define

$$f(x, y) = \begin{bmatrix} x^2 - 3y^2 + 3 \\ xy + 6 \end{bmatrix}, \quad (x, y) \in \mathbb{R}^2.$$

A homotopy is defined by

$$\begin{aligned} h(t, (x, y)) &= tf(x, y) + (1 - t)(f(x, y) - f(1, 1)) \\ &= f(x, y) + tf(1, 1) - f(1, 1), \quad t \in [0, 1], (x, y) \in \mathbb{R}^2, \end{aligned}$$

$$h_x(t, (x, y)) = Df(x, y) = \begin{bmatrix} \frac{\partial f_1}{\partial x}(x, y) & \frac{\partial f_1}{\partial y}(x, y) \\ \frac{\partial f_2}{\partial x}(x, y) & \frac{\partial f_2}{\partial y}(x, y) \end{bmatrix} = \begin{bmatrix} 2x & -6y \\ y & x \end{bmatrix},$$

$$h_t(t, (x, y)) = f(1, 1) = \begin{bmatrix} 1 \\ 7 \end{bmatrix}.$$



## Example (cont'd)

$$h_x^{-1}(t, (x, y)) = [Df(x, y)]^{-1} = \frac{1}{2x^2 + 6y^2} \begin{bmatrix} x & 6y \\ -y & 2x \end{bmatrix}.$$

The ODE is

$$\begin{bmatrix} x'(t) \\ y'(t) \end{bmatrix} = -\frac{1}{2x^2 + 6y^2} \begin{bmatrix} x & 6y \\ -y & 2x \end{bmatrix} \begin{bmatrix} 1 \\ 7 \end{bmatrix} = -\frac{1}{2x^2 + 6y^2} \begin{bmatrix} x + 42y \\ 14x - y \end{bmatrix}.$$

with initial condition  $(x(0), y(0))^T = (1, 1)^T$ . By the numerical method for initial-value problem, we have an approximation solution  $(-2.961, 1.978)^T$  of  $(x(1), y(1))^T$ . We can use this approximation as the initial guess in the Newton method:

$k$	$(x^{(k)}, y^{(k)})$	$\ f(x^{(k)}, y^{(k)})\ _2$
0	$(-2.96100000000000, 1.97800000000000)$	0.14626611680427
1	$(-3.00025328131376, 2.00012057060499)$	0.00087135657948
2	$(-3.00000001019155, 2.00000000338437)$	0.00000003679978
3	$(-3.00000000000000, 2.00000000000000)$	0.00000000000000

See the details of the M-file: [homotopynewton.m](#)

## Theorem on continuously differentiable solution

---

[Ortega and Rheinboldt, 1970]

*If  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is continuously differentiable and if  $\|[Df(x)]^{-1}\| \leq M$  on  $\mathbb{R}^n$ , then for any  $x_0 \in \mathbb{R}^n$  there is a unique curve  $\{x(t) : 0 \leq t \leq 1\}$  in  $\mathbb{R}^n$  such that  $f(x(t)) + (t - 1)f(x_0) = 0$ ,  $0 \leq t \leq 1$ . The function  $t \rightarrow x(t)$  is a continuously differentiable solution of the initial-value problem  $x'(t) = -[Df(x)]^{-1}f(x_0)$ , where  $x(0) = x_0$ .*

**Note:** 
$$tf(x(t)) + (1 - t) \underbrace{(f(x(t)) - f(x_0))}_{:=g(x(t))} = f(x(t)) + (t - 1)f(x_0).$$