

Alessio Hu (12345678)

# Numerical Analysis Summary

Numerical Analysis  
Prof. Perotto Simona  
- Politecnico di Milano

**Something:**

Something else



A.Y. 2022-2023

# Contents

<b>1</b>	<b>Intro</b>	<b>3</b>
1.1	Numerical Analysis and Errors . . . . .	3
1.2	Floating point representation . . . . .	4
<b>2</b>	<b>Nonlinear equations</b>	<b>5</b>
2.1	Bisection method . . . . .	6
2.1.1	Pros and Cons . . . . .	8
2.2	Newton method . . . . .	8
2.2.1	Taylor expansion . . . . .	9
2.2.2	Comparison with bisection . . . . .	9
2.2.3	Convergence . . . . .	9
2.2.4	Modified Newton scheme . . . . .	10
2.2.5	System of nonlinear equations, vector . . . . .	11
2.2.6	Bisection - Newton method . . . . .	12
2.3	Convergence order . . . . .	12
2.4	Stopping criteria/point . . . . .	13
2.4.1	Reliability of the residual . . . . .	14
2.5	Fixed Point Method . . . . .	15
2.5.1	Problems correlation . . . . .	16
2.5.2	The method with Newton and Bisection . . . . .	16
2.5.3	Convergence . . . . .	17
2.5.4	Global and Local Convergence Results . . . . .	18
2.5.5	Examples for the Local Result . . . . .	19
2.5.6	Fixed point method with any order of convergence . . . . .	19
2.5.7	Stopping criteria . . . . .	20
2.6	Consistency . . . . .	20
<b>3</b>	<b>Systems of Linear Equations: Direct Methods</b>	<b>21</b>
3.1	Cramer Method . . . . .	21
3.2	Numerical Approximations: Direct Methods . . . . .	21
3.2.1	LU factorization of a matrix . . . . .	21
3.2.2	Forward substitution method . . . . .	22
3.2.3	Backward substitution method . . . . .	23
3.2.4	Direct inspection . . . . .	23
3.2.5	Gaussian elimination method . . . . .	24
3.2.6	Necessary and sufficient criteria for GEM . . . . .	26
3.2.7	LU strategies for particular matrices . . . . .	28
3.2.8	Rows swapping: pivoting . . . . .	30
3.3	Accuracy . . . . .	31
3.3.1	Lemmas . . . . .	32
3.3.2	Perturbations relation . . . . .	32
3.3.3	Condition number . . . . .	33
3.3.4	Perturbation on the matrix A . . . . .	35
<b>4</b>	<b>Systems of Linear Equations: Iterative Methods</b>	<b>36</b>
4.1	Iterative methods . . . . .	36
4.1.1	Convergence . . . . .	36
4.1.2	Consistency . . . . .	36
4.1.3	Convergence analysis . . . . .	36
4.2	Richardson schemes . . . . .	37

4.3	Stationary Richardson Schemes . . . . .	38
4.3.1	Jacobi and Gauss-Seidel methods . . . . .	38
4.3.2	Convergence for Jacobi and Gauss-Seidel . . . . .	40
4.3.3	Optimal acceleration parameter for stationary schemes . . . . .	41
4.4	Dynamic Richardson Schemes . . . . .	43
4.4.1	The algorithm . . . . .	43
4.4.2	Stationary and dynamic: which one is the best . . . . .	44
4.4.3	Proof for the optimal acceleration parameter . . . . .	45
4.5	Conjugate gradient method . . . . .	45
4.6	Stopping Criteria . . . . .	46
4.6.1	Residual . . . . .	46
4.6.2	Increment . . . . .	47
<b>5</b>	<b>Approximation of Functions and Data</b>	<b>48</b>
5.1	Interpolation . . . . .	48
5.1.1	Polynomial interpolation . . . . .	48
5.1.2	Interpolation error . . . . .	50
5.1.3	Piecewise linear interpolation . . . . .	52
5.1.4	Cubic spline interpolation . . . . .	54
5.2	Least Squared Approximation . . . . .	54
5.2.1	Degree n . . . . .	54
5.2.2	Degree 1 . . . . .	55
5.2.3	Degree m generic . . . . .	55
5.3	Numerical Integration . . . . .	56
5.3.1	Newton-Cotes . . . . .	56

# 1 Intro

## 1.1 Numerical Analysis and Errors

M = Math

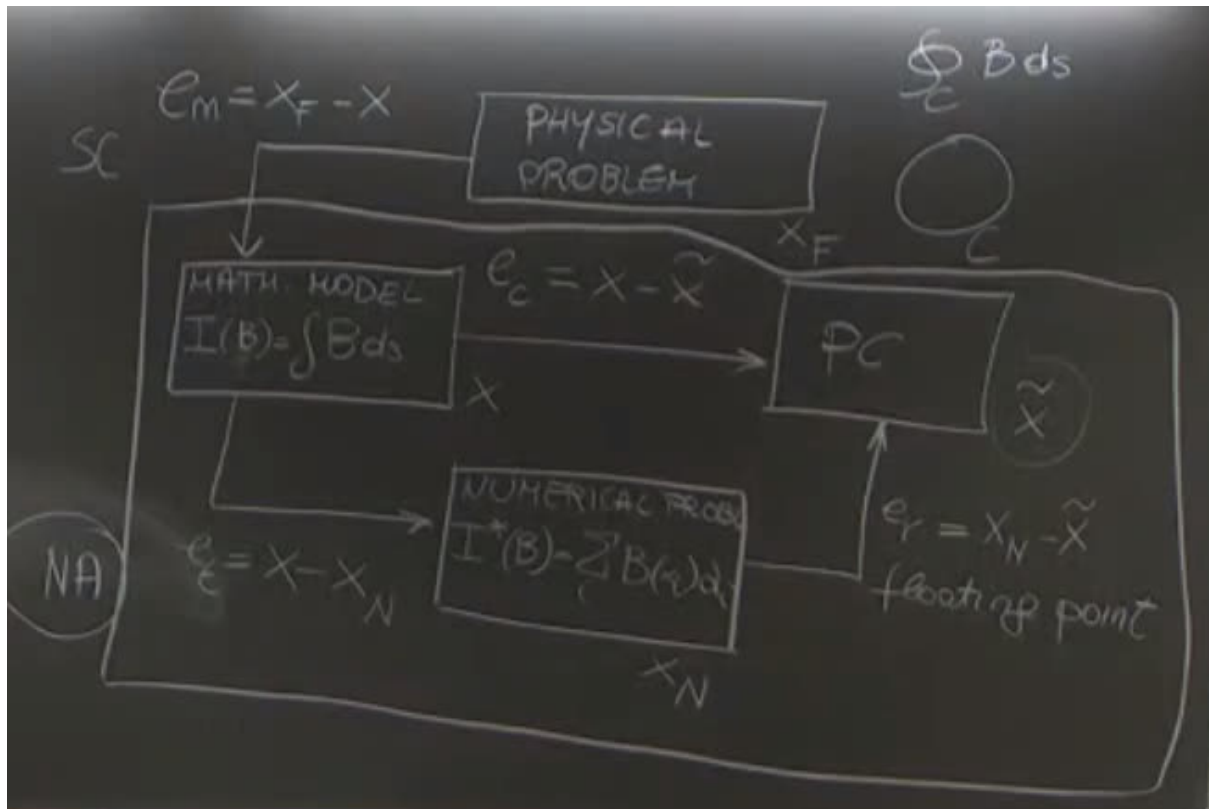
CS = Computer Science

E = Engineering

Numerical Analysis =  $M \cap CS$

Basic methods to approach math problems

Scientific Computing =  $M \cap CS \cap E$  Take a problem and replicate it on a digital device to understand better the situation



Where:

$$\begin{cases} x_F \text{ solution of physical model} \\ x \text{ solution of mathematical model} \\ \tilde{x} \text{ does not substitute reality simulation} \end{cases}$$

We replaced the integral with a summation, ia PC there is no concept of infinity.

What to do after we observe what's going on? Use a better model or a better  $x \leftrightarrow \tilde{x}$  mapping.

Errors:

$$\begin{cases} e_m = x_F - x \text{ modelling error between physical problem and mathematical model} \\ e_c = x - \tilde{x} \text{ computational error} = e_t + e_r \begin{cases} e_t = x - x_N \text{ truncation error} \\ e_r = x_N - \tilde{x} \text{ rounding error, floating point approximation} \end{cases} \end{cases}$$

Kinds of errors:

- $|x - \tilde{x}|$   
**Absolute error**,  $|e_c|$ . Consider absolute error based on model

- $\frac{|x-\tilde{x}|}{|x|}$   
**Relative error**, more meaningful, like percentage

## 1.2 Floating point representation

Not really important, skipping...

## 2 Nonlinear equations

$f$ , we want to find  $\alpha \in \mathbb{R}$  zero of  $f$  such that  $f(\alpha) = 0$

It is not easy to find the zero of a function when  $\mathbb{P}_n$   $n > 4$ . From here we need to approximate the zeros of a nonlinear function.

An example of physical phenomenon: the ideal gas equation:

$$pV = nRT$$

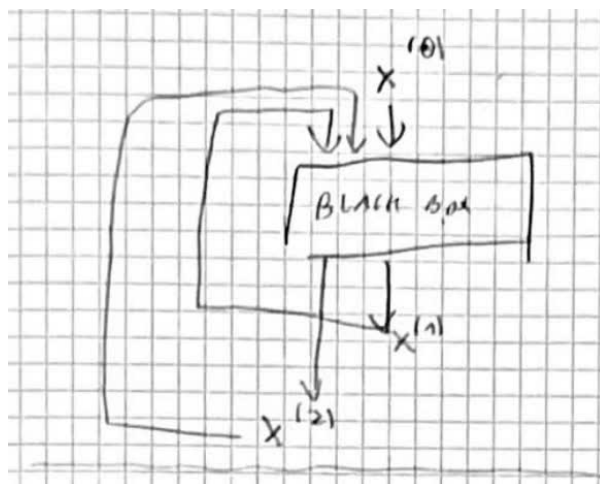
To find  $V$

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

A nonlinear equation, although we know the pressure, the temperature and the constants  $a$  and  $k$ , it is not easy to solve.

We use **iterative methods**

$x^{(0)}$   
Initial guess



$$x^{(k)} \simeq \alpha$$

Ideally I want

$$\lim_{k \rightarrow \infty} x^{(k)} = \alpha$$

**Convergence**, or equivalently

$$e^{(k)} = \alpha - x^{(k)}$$

$$\lim_{k \rightarrow \infty} e^k = 0$$

But we have to decide when to stop this approximation, **stopping criteria**: set a maximum number of iterations and a tolerance error.

We are looking at the approximation error

$$? \alpha \in \mathbb{R} \text{ s.t. } \underbrace{f(\alpha) = 0}_{\text{non linear}}$$

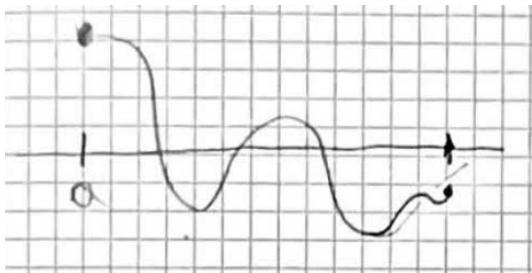
Two methods: bisection and Newton method

## 2.1 Bisection method

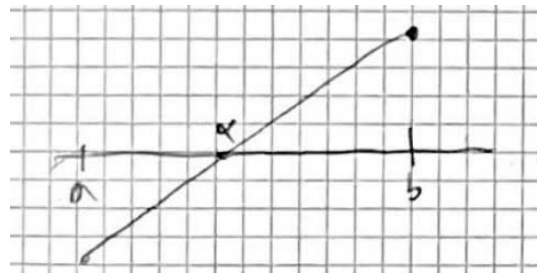
In mathematics, the bisection method is a root-finding method that applies to any continuous function for which one knows two values with opposite signs.

We first have to fix the hypothesis, which coincides on the hypothesis of the zero of nonlinear continuous functions:

1.  $f \in C^0([a, b])$   
Set of all continuous functions in  $[a, b] \subset \mathbb{R}$
2.  $f(a)f(b) < 0$



(a) The example has more zeros even, this is too much for us, even one is enough



(b) The function is taking values at opposite sign at endpoints, which means that the function  $f$  has at least one zero in the interval

Starting from an interval, we shrink it:

$$\begin{aligned}
 \alpha &\in I^{(0)} = [a, b] \\
 \alpha &\in I^{(1)} \subset I^{(0)} \\
 |I^{(0)}| &= \frac{|I^{(0)}|}{2} \\
 \alpha &\in I^{(2)} \subset I^{(1)} \\
 |I^{(2)}| &= \frac{|I^{(1)}|}{2} \\
 &\vdots
 \end{aligned}$$

Collection of intervals that are getting smaller and smaller, that all contain  $\alpha$  (the zero), so at the end we will get to  $\alpha$

Strength point of bisection: **always convergent, but has a lot of drawbacks too.**

But we are looking for  $x^{(k)}$  that approximates  $\alpha$

$$? x^{(k)} \simeq \alpha$$

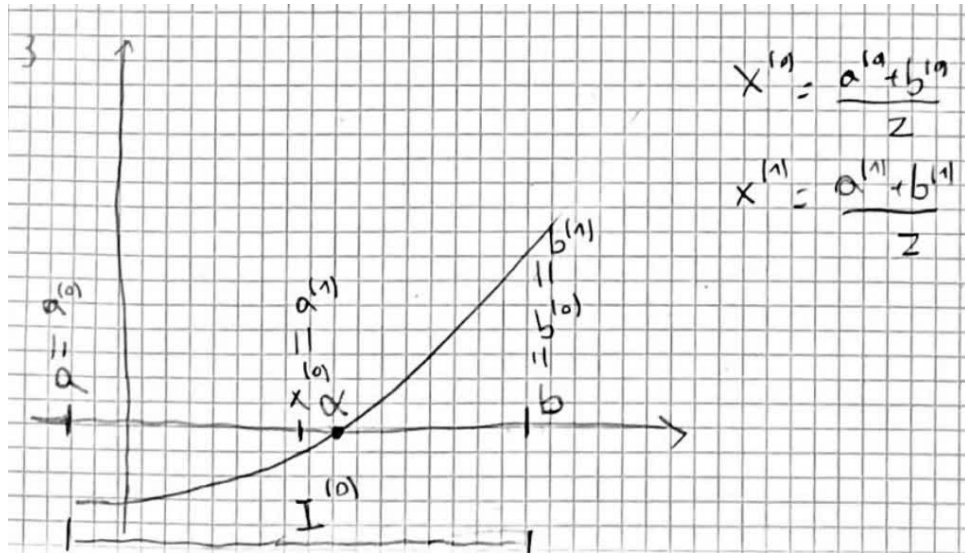
To do this we consider the midpoint of the interval

We see that in this example the first midpoint is already close to  $\alpha$ . Now we have to choose a new subinterval  $I^{(1)}$ , we will choose the right subinterval as we need to satisfy the hypothesis that the extremes have opposite signs. We then continue iteratively.

The algorithm is similar to binary search. Formally:

**Inputs:**  $a = a^{(0)}$ ,  $b^{(0)} = v$ ,  $f$ ;  $TOL$ ,  $Nmax$

Where  $TOL$  is the tolerance and  $Nmax$  the maximum number of iterations.



```

while (true)
    x(k) = (a(k) + b(k)) / 2
    if (f(x(k-1)) = 0) break;
    if f(a(k-1))f(x(k-1)) < 0
        a(k) = a(k-1), b(k) = x(k-1);
    else
        a(k) = x(k-1), b(k) = b(k-1);
    end

```

Regarding the intervals, we see that:

$$|I^{(k)}| = \frac{b-a}{2^k} = \frac{|I^{(0)}|}{2^k}$$

$$|e^{(k)}| = |\alpha - x^{(k)}| < \frac{|I^{(k)}|}{2} = \frac{b-a}{2^{k+1}}$$

$$\lim_{k \rightarrow \infty} |e^{(k)}| = (b-a) \lim_{k \rightarrow \infty} \left(\frac{1}{2}\right)^{(k+1)} = 0$$

Regarding the tolerance

$$? k \text{ s.t. } |e^{(k)}| \leq TOL = 10^{-9}$$

$$\frac{b-a}{2^{(k+1)}} \leq TOL$$

$$\frac{b-a}{TOL} \leq 2^{(k+1)}$$

$$\log_2 \left( \frac{b-a}{TOL} \right) \leq k+1$$

$$k \geq \log_2 \left( \frac{b-a}{TOL} \right) - 1$$

The right member will represent the  $N_{max}$

$$N_{max} = \left\lceil \log_2 \left( \frac{b-a}{TOL} \right) \right\rceil = \left\lceil \log_{10} \left( \frac{b-a}{TOL} \right) / \log_{10} 2 \right\rceil$$



### 2.1.1 Pros and Cons

Pros:

- Converges
- Can compute maximum number of iterations

Cons:

- We are losing the monotonicity of the method, a subsequent guess is not guaranteed to find a  $x$  closer to  $\alpha$  w.r.t. the previous  $x$ , **not monotonic w.r.t. to error, the error does not necessarily decreases at each step, we have no convergence order**
- Cannot find  $\alpha$  in a single step
- We are only exploiting the fact that the function is changing sign, function values of  $f$  are ignored

### 2.2 Newton method

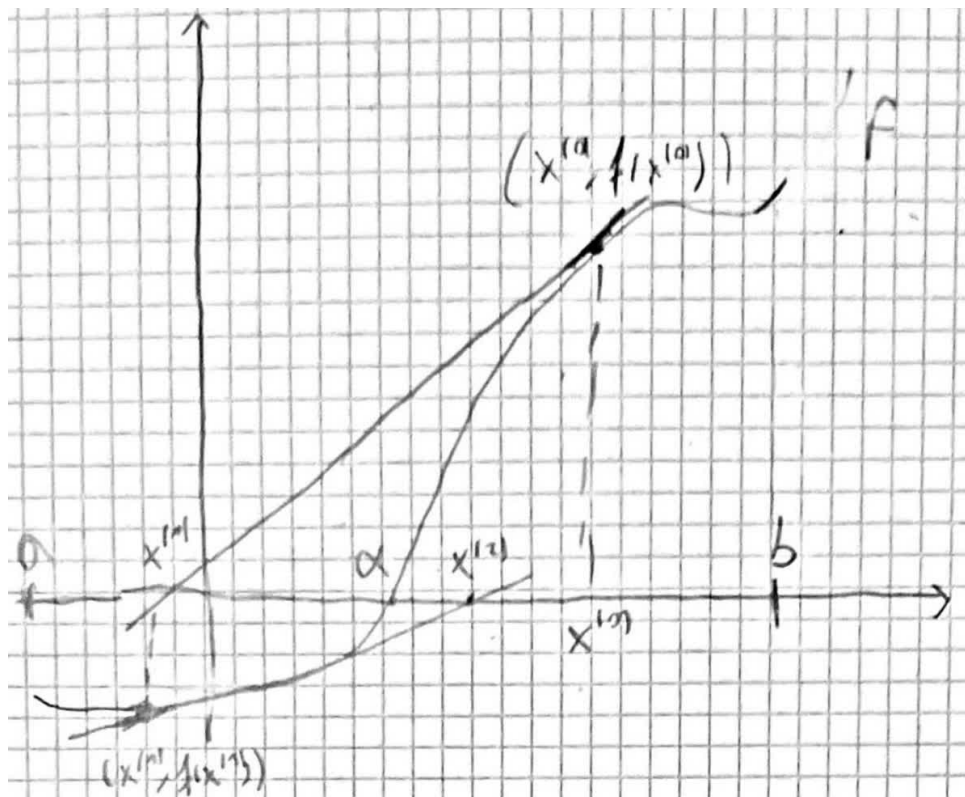
$N_{max}$  has to hold for all continuous functions in that interval  $[a, b]$ , so it has to be large.

The Newton we are not only exploiting the sign of  $f$ , but also the values. It is the most powerful method, and demands that:

$$f \in C^1([0, 0])$$

Set of functions continuous in their first derivative

Replace  $f$  with the tangent line at point of our initial guess  $x^{(0)}$ , then take the intersection of the tangent with the  $x$  axis, that will be our  $x^{(1)}$ . Then repeat:



tg to  $f$  at  $(x^{(k)}, f(x^{(k)}))$   
 $y(x) = f(x^{(k)}) + f'(x^{(k)})(x - x^{(k)})$   
 $x^{(k+1)}$  s.t.  $y(x^{(k+1)}) = 0$

$$f(x^{(k)}) + f'(x^{(k)})(x^{(k+1)} - x^{(k)})$$

From this equality we want to derive  $x^{(k+1)}$ :

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})} \quad k \geq 0$$

Assuming  $f'(x^{(k)}) \neq 0$

It can be proved that this algorithm is just a truncation of the Taylor expansion

### 2.2.1 Taylor expansion

We must first decide the center and where to evaluate the expansion. In our case the center is  $x^{(k)}$ , we evaluate at  $x^{(k+1)}$

$$f(x) = f(x^{(k)}) + f'(x^{(k)})(x - x^{(k)}) + O\left((x - x^{(k)})^2\right)$$

So if we evaluate at  $x^{(k+1)}$ , we simply replace  $x$ . We neglect the big  $O$  term and if  $k$  is sufficiently large, we can approximate to  $\alpha$ .

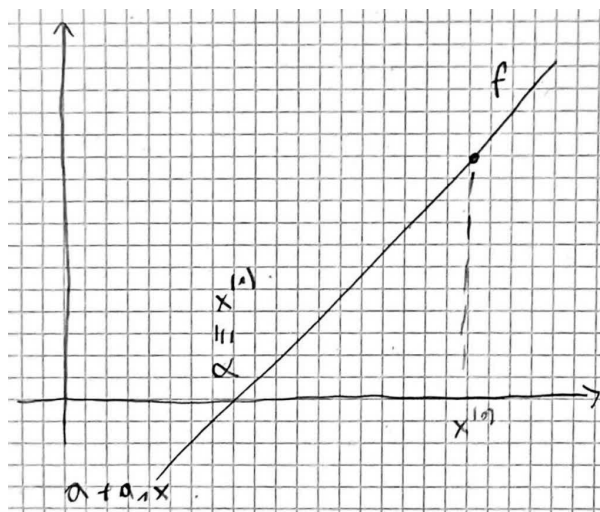
$$0 = f(\alpha) \simeq f(x^{(k+1)}) \cong f(x^{(k)}) + f'(x^{(k)})(x^{(k+1)} - x^{(k)})$$

Which is exactly the Newton method

### 2.2.2 Comparison with bisection

Newton can identify  $\alpha$  in a single step.

Geometrical proof:



Analytical proof:

$$x^{(1)} = x^{(0)} - \frac{f(x^{(0)})}{f'(x^{(0)})} = x^{(0)} - \frac{a_0 + a_1 x^{(0)}}{a_1} = -\frac{a_0}{a_1}$$

Bisection converges without initial guess, but Newton is better even though it needs an initial guess. **If initial guess not close to the zero, we do not converge**

### 2.2.3 Convergence

Two hypothesis

H1)  $x^{(0)}$  sufficiently close to  $\alpha$ . But we do not know  $\alpha$ , we could:

- Graphically plot it
- Do some steps of the bisection and use some outputs of it for Newton: **bisection-Newton: predictor-corrector**, use a weaker method then a stronger one. With bisection we know that we will converge, even if slowly, then after a desired steps (sufficiently close) we use Newton which is faster

H2)  $\alpha$  is a simple zero of  $f$

$$\begin{cases} f(\alpha) = 0 \\ f'(\alpha) \neq 0 \end{cases}$$

Reminder, an  $\alpha$  is a zero of order  $m$  if

$$\begin{cases} f(\alpha) = f'(\alpha) = f''(\alpha) = \dots = f^{(m-1)}(\alpha) = 0 \\ f^{(m)}(\alpha) \neq 0 \end{cases}$$

$\Rightarrow$  **Newton is convergent**

Adding a third hypothesis

H3)  $f \in C^2([a, b])$ , we can say that the following limit holds:

$$\lim_{k \rightarrow \infty} \frac{x^{(k+1)} - \alpha}{[x^{(k)} - \alpha]^2} = \underbrace{\frac{f''(\alpha)}{2f'(\alpha)}}_C$$

In general **convergence order** equal to  $P$  if  $\exists c$  independent from  $k$ , such that

$$\frac{|x^{k+1} - \alpha|}{|x^k - \alpha|^P} \leq C, \forall k \geq k_0$$

If  $P = 1$ , linear convergence. In our case  $P = 2$ , quadratic convergence. In general, with a higher convergence order the error reduces:

$$\begin{aligned} x^{(k)} - \alpha &= 10^{-2} \\ x^{(k)} - \alpha &\simeq 10^{-4} & P = 2 \\ x^{(k)} - \alpha &\simeq 10^{-6} & P = 3 \end{aligned}$$

The constant  $C$  does not have any requirement, but in some sense it is slowing the convergence (for  $C = 10, P = 3$ , the error  $\simeq 10^{-6} * 10 = 10^{-5}$ ), but for  $P = 1$  reduction of the error not guaranteed if  $C = 1$ , so:

$$P = 1 \rightarrow C < 1$$

In the Newton case:

$$C = \frac{f''(\alpha)}{2f'(\alpha)}$$

## 2.2.4 Modified Newton scheme

What if H2 does not hold, can we still use Newton? Yes, but we lose the quadratic order of the convergence.

If  $\alpha$  is a multiple zero of  $f$  (multiplicity  $m$ ) and if  $x^{(0)}$  is sufficiently close to  $\alpha \Rightarrow$  Newton converges linearly ( $P = 1$ , we lost an order of convergence)

We can use the modified Newton scheme:

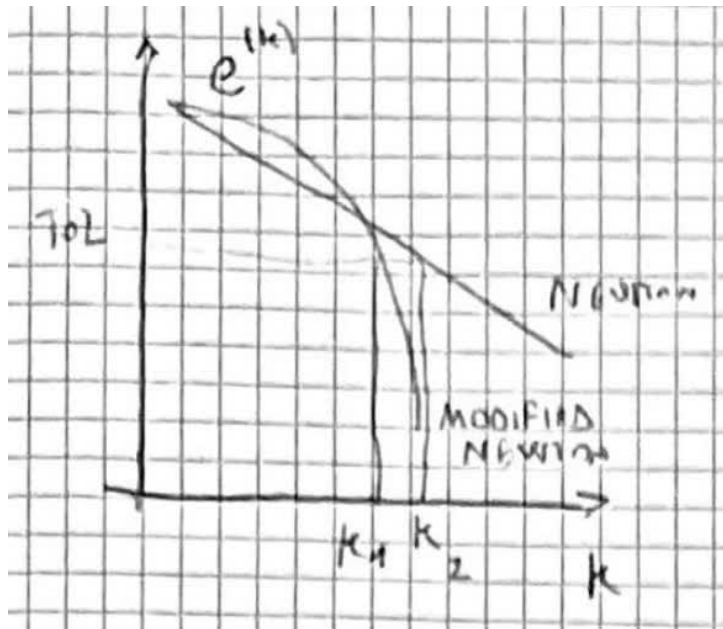
$$x^{(k+1)} = x^{(k)} - m \frac{f(x^{(k)})}{f'(x^{(k)})} \quad k \geq 0$$

Example:

$$f(x) = (x-1)\log(x)$$

$$\alpha = 1$$

$$m = 2$$



## 2.2.5 System of nonlinear equations, vector

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

Example

$$\begin{cases} f_1(x_1, x_2) = x_1^3 + \sin x_2 = 0 \\ f_2(x_1, x_2) = -x_1\sqrt{x_2} + \operatorname{tg}\left(\frac{x_2}{3x_1}\right) = 0 \end{cases}$$

We can rewrite this complex model to a more manageable form, a **vectorial way**

$$\vec{x} = [x_1, x_2, \dots, x_n]^T$$

$$\vec{f} = [f_1(\vec{x}), f_2(\vec{x}), \dots, f_n(\vec{x})]^T$$

$$\vec{0} = [0, \dots, 0]^T \in \mathbb{R}^n$$

So we can rewrite the system of nonlinear equations in:

$$\vec{f}(\vec{x}) = \vec{0}$$

Now applying the Newton method, we just introduce the vectors that contain the  $k$ -approximation error:

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}$$

Now with

$$\begin{aligned}\vec{x}^{(k)} &= [x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)}]^T \in \mathbb{R}^n \\ \vec{x}^{(k+1)} &= [x_1^{(k+1)}, x_2^{(k+1)}, \dots, x_n^{(k+1)}]^T \in \mathbb{R}^n \\ \vec{f}(\vec{x}^{(k)}) &= [f_1(\vec{x}^{(k)}), f_2(\vec{x}^{(k)}), \dots, f_n(\vec{x}^{(k)})]^T \in \mathbb{R}^n\end{aligned}$$

What about the derivative? Use Jacobian

$$x^{(k+1)} = x^{(k)} - \underbrace{\frac{f(x^{(k)})}{f'(x^{(k)})}}_{\delta x^{(k)}}$$

$$f'(x^{(k)})\delta x^{(k)} = -f(x^{(k)})$$

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} + \delta \vec{x}^{(k)}$$

$$? \delta \vec{x}^{(k)} = -\vec{f}'(\vec{x}^{(k)})$$

The Jacobian

$$(J_F)_y = \frac{\partial f_i}{\partial x_j} \quad i, j = 1, \dots, n$$

The question mark becomes

$$\underbrace{J_F(\vec{x}^{(k)})}_{\mathbb{R}^{n \times n}} \underbrace{\delta \vec{x}^{(k)}}_{\mathbb{R}^n} = -\underbrace{\vec{f}'(\vec{x}^{(k)})}_{\mathbb{R}^n}$$

## 2.2.6 Bisection - Newton method

They are predictor-corrector methods. Newton is conditional convergent (the initial guess must be close to the root), while bisection is unconditional convergent: **by using bisection we are sure that the initial guess is close to the root (predictor) then use Newton as corrector.**

## 2.3 Convergence order

$$\{x^{(k)}\} \simeq \alpha$$

Convergence order equal to  $p$  if  $\exists c$  independent from  $k$ , such that  $\frac{|x^{k+1}-\alpha|}{|x^k-\alpha|^p} \leq c, \forall k \geq k_0$

- **Bisection**, error no monotone, no convergence order
- **Newton**

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})} \quad k \geq 0$$

In Newton, under assumptions that:

- $f(\alpha) = 0$
- $f'(\alpha) \neq 0$

$\alpha$  simple zero, we have

- Convergence order of 1 if  $f \in C^1$
- Convergence order of 2 if  $f \in C^2$

If  $\alpha$  not simple,  $f \in C^2$ ,  $\alpha$  is a zero of order  $m$  ( $f(\alpha) = 0, \dots, f^{(m-1)}(\alpha) = 0, f^{(m)}(\alpha) \neq 0$ ), more simply when there is a power the zero is of the order of that power.

- **Modified Newton method**, unlike bisection, with newton we can compute the convergence order. Can find non-simple zeros, we have a convergence order of 2 for example. We use a different update rule:

$$x^{(k+1)} = x^{(k)} - m \frac{f(x^{(k)})}{f'(x^{(k)})} \quad k \geq 0$$

$$e = \{|x^k - \alpha|\}_k$$

$$e^k = |x^k - \alpha|$$

$$\text{Convergence order} = p = \frac{\log \left[ \frac{e^{k+2}}{e^{k+1}} \right]}{\log \left[ \frac{e^{k+1}}{e^k} \right]}$$

**Do not mix up exponential and error, the  $e^k$  there stands for the error vector considering index from  $k$  to  $end$**

## 2.4 Stopping criteria/point

$$\underbrace{|x^{(k)} - \alpha|}_{e^{(k)}} \leq \underbrace{CS}_{\text{Error estimator}} < TOL$$

$$\left| e^{(kmin)} \right| = \left| x^{(kmin)} - \alpha \right| \leq CS = C \left| x^{(kmin+1)} - x^{(kmin)} \right| \leq TOL \text{ (e.g. } = 10^{-9} \text{)}$$

A constant that multiplies our error estimator  $S$ . If large like  $10^4$  we lose 4 orders. It is called the **reliability**, if

- $C = O(1)$ , estimator reliable
- $C = O(10^s)$ , not reliable

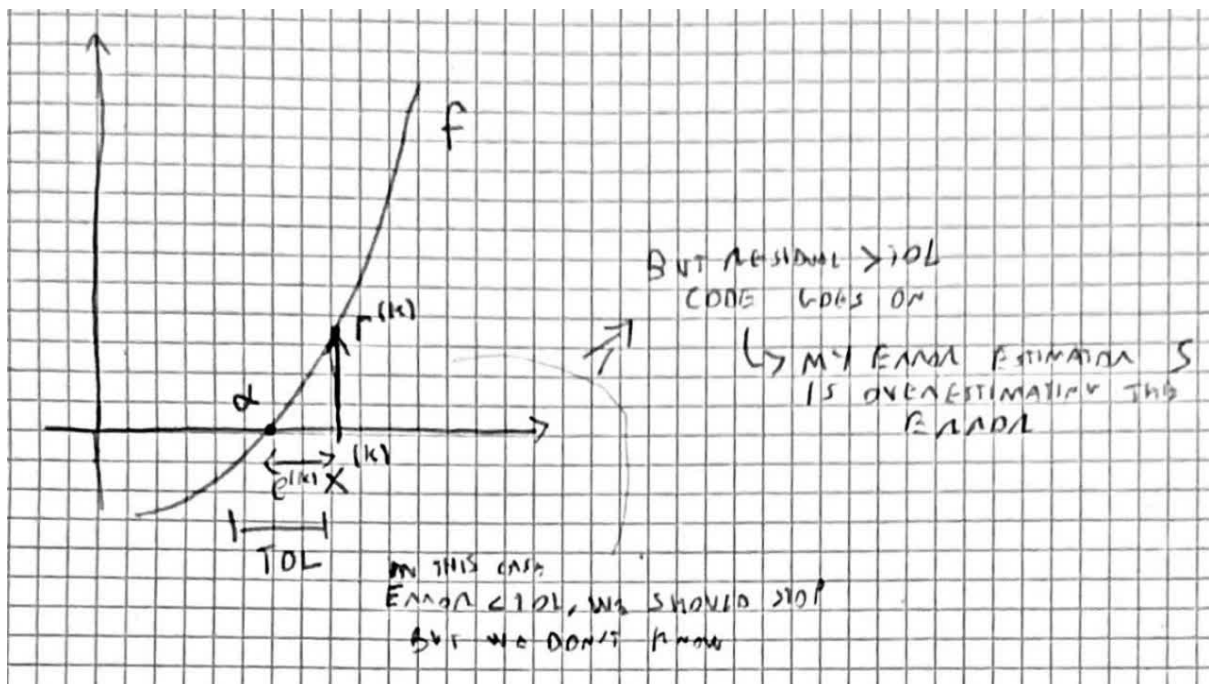
Two kinds of estimators for the error, when one not reliable, can rely on the other one

$$S = \begin{cases} x^{(k+1)} - x^{(k)} \text{ increment, if convergence this difference becomes smaller and smaller} \\ \textbf{Newton} \\ r^{(k)} = f(x^{(k)}) \text{ residual, huge if } x^{(k)} \text{ far from } \alpha, \text{ less it is, smaller is the error} \\ \textbf{Bisection} \end{cases}$$

The increment one, cycle till:

$$|x^{(k+1)} - x^{(k)}| > TOL \ \&\& \ i < TOL$$

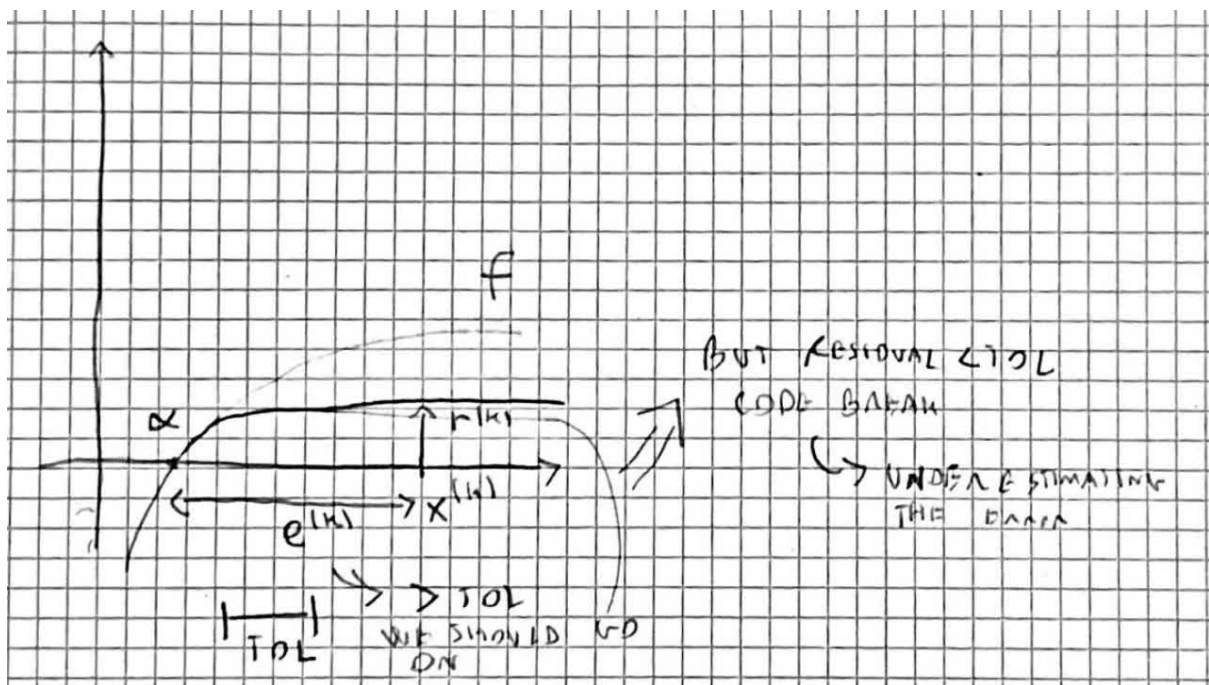
### 2.4.1 Reliability of the residual



**Figure 1:** Overestimating

My error estimator  $S$  is overestimating the error  $e^{(k)}$

$$|f'(\alpha)| \gg 1$$



**Figure 2:** Underestimating

My error estimator  $S$  is underestimating the error  $e^{(k)}$

$$|f'(\alpha)| \ll 1$$

What's better? Better when overestimating, we do a little more work, but at the end we get the better result

So this estimator is reliable when:

$$|f'(\alpha)| \simeq 1$$

## 2.5 Fixed Point Method

If we apply continuously for example the cos:

$$\begin{aligned} x^{(0)} &= 1 \\ x^{(1)} &= \cos(x^{(0)}) \\ x^{(2)} &= \cos(x^{(1)}) \\ &\vdots \\ x^{(k+1)} &= \cos(x^{(k)}) \end{aligned}$$

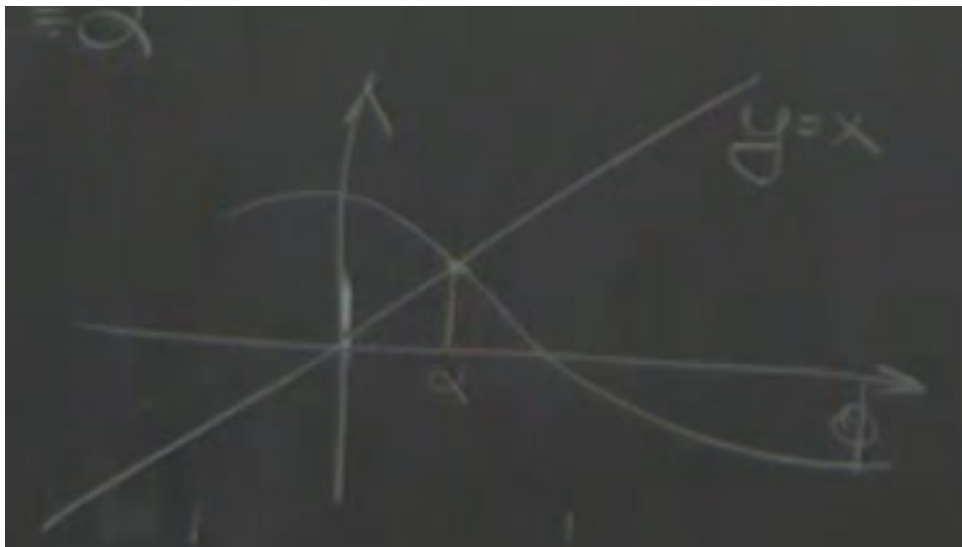
Which means  $\alpha = \cos(\alpha)$ , with  $\alpha$  known as **fixed point** of the function cos. In general:

$$\Phi : [a, b] \subset \mathbb{R} \rightarrow \mathbb{R}$$

$$\alpha \in \mathbb{R} \text{ fixed point } \phi(\alpha) = \alpha$$

Geometrically we are looking for the intersection of:

$$\begin{cases} y = \phi(x) \\ y = x \end{cases} \quad \text{Bisector line}$$



- Does any function has a fixed point? No, for example parallel line or exponential function, it grows fast and does not encounter the bisector line
- Does a function admit more fixed points? Yes, function can encounter the bisector line multiple times

How can a fixed point interest us? We can talk about a sort of duality problem:

$$\begin{cases} ? \alpha \in \mathbb{R} \text{ st } f(\alpha) = 0 & \text{Our original problem of zeros} \\ ? \alpha \in \mathbb{R} \text{ st } \phi(\alpha) = \alpha & \text{Fixed point problem} \end{cases}$$



### 2.5.1 Problems correlation

The transition:

$$\left[ \begin{array}{ccc} ? \alpha \in \mathbb{R} \text{ st } f(\alpha) = 0 & \Leftrightarrow & ? \alpha \in \mathbb{R} \text{ st } \phi(\alpha) = \alpha \\ f(x) = 0 & \Leftrightarrow & \underbrace{f(x) + x = x}_{\phi(x)} \\ \Rightarrow \text{Hp: } f(\alpha) = 0 & & \Leftarrow \text{Hp: } \phi(\alpha) = \alpha \\ \text{Th: } \phi(\alpha) = \alpha & & \text{Th: } f(\alpha) = 0 \\ \phi(\alpha) = \underbrace{f(\alpha)}_{=0} + \alpha = \alpha & & \phi(\alpha) = \underbrace{f(\alpha)}_{=\alpha} + \alpha = \alpha \end{array} \right]$$

But the  $\phi$  is not unique, we can build more and different fixed point functions (e.g. add constant multiplier), this is a **advantage**, we can pick a function that is for sure to have a fixed point

### 2.5.2 The method with Newton and Bisection

With an **iterative process**, now we try to solve

$$? \alpha \in \mathbb{R} \text{ st } \phi(\alpha) = \alpha$$

With iterative method

$$x^{(k+1)} = \phi(x^{(k)}) \quad k \geq 0$$

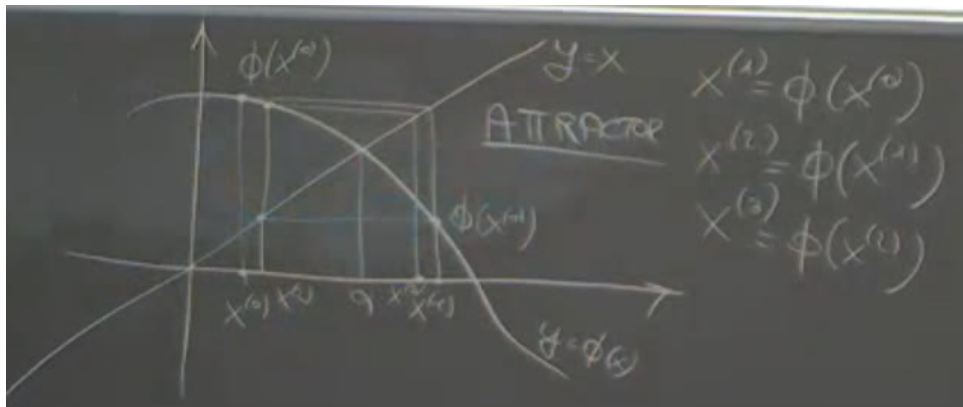
With initial guess  $x^{(0)}$ . Considering the previous methods bisection and Newton, can we rewrite them as a fixed point method?

- **Newton:** we can just pick a fixed point function as:

$$\phi_N(x) = x - \frac{f(x)}{f'(x)}$$

- **Bisection:** no, the midpoint depends on two variables, so bisection is not an example of fixed point method

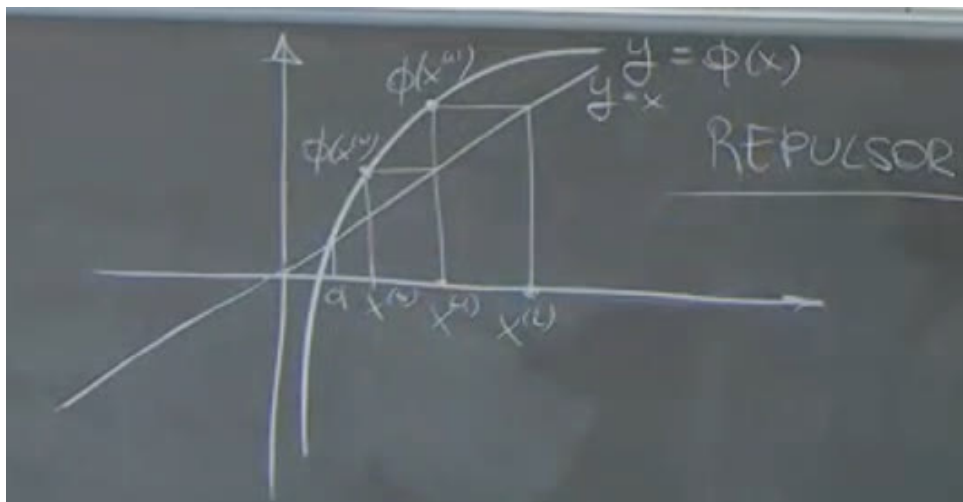
### 2.5.3 Convergence



**Figure 3:** Attractor

All even at left side, all even at right side.

Are we converging to  $\alpha$ ? Yes, the iterations are moving closer and closer to  $\alpha$ , which is an **attractor**.



**Figure 4:** Repulsor

In this case we are unlucky, we diverge, **repulsor**.

These are two particular cases, in total there are four: left-right convergent, left-right divergent, one direction convergent, one direction divergent.

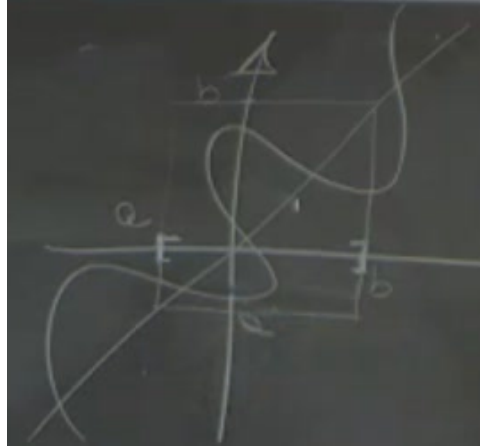
Can we identify some features that are responsible for the convergent and divergent trend? The **value of the derivative  $\phi'(\alpha)$  if less or greater than 1**

## 2.5.4 Global and Local Convergence Results

$$x^{(k+1)} = \phi(x^{(k)})$$

- **Global result**

1. Let  $\phi \in C^0([a, b])$  and s.t.  $\phi \in [a, b] \forall x \in [a, b]$



We are demanding that our function  $\phi$  has to take values inside that rectangle, must be limited in a certain portion of the plane.

Under this hypothesis, we can prove that exists at least a fixed point  $\alpha \in [a, b]$  for function  $\phi$  (the example has two), so **not uniqueness of fixed point**

2. If we in addition we assume that there exists an integer  $L < 1$  s.t.

$$|\phi(x_1) - \phi(x_2)| \leq L|x_1 - x_2| \forall x_1, x_2 \in [a, b]$$

(Lipschitz continuity, weaker demand w.r.t. derivability, weaker than  $C^1$ ) then  $\exists! \alpha \in [a, b]$  for  $\phi$  and

$$\{x^{(k)}\} \rightarrow \alpha \forall x^{(0)} \in \mathbb{R}$$

Collection of approximations, with this assumption we get **uniqueness of fixed point and convergence of method independently from initial guess.**

About the **rate of convergence of fixed point scheme with Lipschitz**, consider the error associated with  $k + 1$ :

$$|x^{(k+1)} - \alpha| = |\phi(x^{(k)}) - \phi(\alpha)| \leq L|x^{(k)} - \alpha|$$

$$\frac{|x^{(k+1)} - \alpha|}{|x^{(k)} - \alpha|} \leq L < 1$$

Which means that the fixed method is convergent with order 1 (power below is 1, but actually  $p \geq 1$ , so **order at least 1**)

But this practical is not practical

- **Local result or Ostrowski's theorem:** let  $\alpha$  be a fixed point for  $\phi$  in  $[a, b]$  (so we are already assuming the existence and uniqueness of  $\alpha$ ), with  $\phi \in C^1(I_\alpha)$  and  $I_\alpha$  neighborhood of  $\alpha$  (different from before, we here stronger assumptions than global which only required  $C^0$  and Lipschitz continuity, but  $C^1$  locally).

Under these hypotheses, if  $|\phi'(\alpha)| < 1$  then  $\exists \delta > 0$  s.t.  $\forall x^{(0)}$  with  $|x^{(0)} - \alpha| < \delta$ :

$$\{x^{(k)}\} \rightarrow \alpha \text{ and } \lim_{k \rightarrow \infty} \frac{x^{(k+1)} - \alpha}{x^{(k)} - \alpha} = \phi'(\alpha)$$

To summarize if:

$$\begin{cases} |\phi'(\alpha)| < 1 & \text{convergence} \\ |\phi'(\alpha)| > 1 & \text{divergence} \\ |\phi'(\alpha)| = 1 & \text{we cannot say anything} \end{cases}$$

### 2.5.5 Examples for the Local Result

1.  $\phi(x) = \cos(x)$  and  $\phi'(x) = -\sin(x)$

$$|\phi'(\alpha)| = |\sin(\alpha)| < 1 \quad \alpha \neq 0$$

2.  $\phi(x) = x^2 - 1$ , we want to know if fixed point method is convergent or not: first we derive the fixed point of  $\phi$ :

$$x = \phi(x) \rightarrow x^2 - x - 1 = 0$$

$$\alpha_{1,2} = \frac{1 \pm \sqrt{5}}{2}$$

So two fixed point, now compute derivative

$$\phi'(x) = 2x \rightarrow |\phi'(\alpha)| = |1 \pm \sqrt{5}|$$

In neither cases the module is less than 1, so divergent fixed point method

3.  $\log(x) = \gamma$  with  $\gamma \in \mathbb{R}$ , we are demanded to approximate this function with fixed point:

$$f(x) = 0 \rightarrow \underbrace{\log(x) - \gamma}_{f(x)} = 0$$

We can use whatever method/fixed point function we like:

- (a) With Newton

$$\phi_1(x) = \phi_N(x) = x - \frac{\log(x) - \gamma}{\frac{1}{x}} = x(1 - \log(x) + \gamma)$$

- (b)

$$\phi_2(x) = \log(x) - \gamma + x$$

- (c)

$$x \log(x) - \gamma x = 0 \rightarrow \phi_3(x) = \frac{x \log(x)}{\gamma}$$

For  $\gamma = -2$ , we can verify that  $\phi_1$  and  $\phi_3$  are ok, while for the  $\phi_2$  it does not work

### 2.5.6 Fixed point method with any order of convergence

**Proposition:** let us assume that the Ostrowski's theorem hypothesis are verified (**so we are in local setting**). If  $\phi \in C^p(I_\alpha)$  and  $\phi^{(i)}(\alpha) = 0$   $i = 1, \dots, p-1$  (derivatives till  $p-1$ ) and  $\phi^{(p)}(\alpha) \neq 0$  then  $\exists \delta > 0$  s.t.  $\forall x^{(0)}$  with  $|x^{(0)} - \alpha| < \delta$

$$\{x^{(k)}\} \rightarrow \alpha \text{ and } \lim_{k \rightarrow \infty} \frac{x^{(k+1)} - \alpha}{[x^{(k)} - \alpha]^p} = \underbrace{\frac{\phi^{(p)}}{p!}}_C$$

We can guarantee the convergence of our fixed point method and compute the convergence order.

## 2.5.7 Stopping criteria

Reminder

$$\left| e^{(kmin)} \right| = \left| x^{(kmin)} - \alpha \right| \leq CS = \left| x^{(kmin+1)} - x^{(kmin)} \right| \leq TOL \text{ (e.g. } = 10^{-9} \text{)}$$

Reminder:  $\exists \beta_k$  between  $\alpha$  and  $x^{(k)}$  (min value theorem, applicable since  $\phi$  is  $C^1$ )

$$\alpha - x^{(k+1)} = \phi(\alpha) - \phi(x^{(k)}) = \phi'(\beta_k)(\alpha - x^{(k)})$$

So, adding and subtracting  $x^{(k+1)}$ :

$$\alpha - x^{(k)} = \alpha - x^{(k+1)} + \underbrace{x^{(k+1)} - x^{(k)}}_{\delta^{(k)}} = \phi'(\beta_k)(\alpha - x^{(k)}) + \delta^{(k)}$$

$$\alpha - x^{(k)} = \frac{1}{1 - \phi'(\beta_k)}(x^{(k+1)} - x^{(k)})$$

For  $k$  sufficiently large we can identify  $x^{(k)}$  with  $\alpha$  and  $\beta_k$  with  $\alpha$ . So the error associated with the  $k$  iteration is:

$$\alpha - x^{(k)} \cong \underbrace{\frac{1}{1 - \phi'(\alpha)}}_C (x^{(k+1)} - x^{(k)})$$

For  $C$  as much as possible close to 1, it means that  $\phi'(\alpha)$  is very small, almost zero.

$$\phi'(\alpha) = 0 \rightarrow \text{convergence order of at least 2! Quadratic convergence}$$

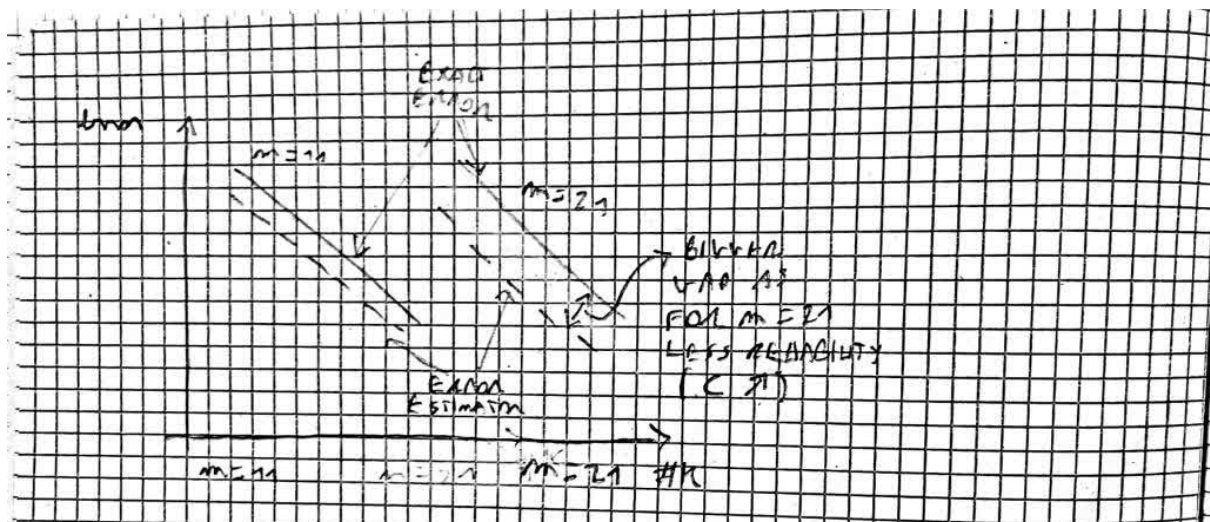
An example:

$$f(x) = (x-1)^{(n-1)} \log(x)$$

With one zero  $\alpha = 1$  with multiplicity  $m$ . We can prove that:

$$\phi'(\alpha) = 1 - \frac{1}{m}$$

Larger is  $m$ , closer that quantity is to 1, which means we are losing in terms of reliability ( $C$  is growing to infinity).



A stopping criterion for fixed point:  $m$  not too high.

## 2.6 Consistency

$$x^{(k+1)} = \phi(x^{(k)})$$

If  $\alpha$  is a fixed point of  $\phi$ , the method is consistent

### 3 Systems of Linear Equations: Direct Methods

We are talking about

$$Ax = b \Leftrightarrow \begin{cases} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n = b_2 \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n = b_n \end{cases}$$

$$A \in \mathbb{R}^{n \times n} \quad x, b \in \mathbb{R}^n$$

#### 3.1 Cramer Method

Reminder:

$$x_i = \frac{\det(A_i)}{\det(A)} \quad i = 1, \dots, n$$

Problem: finding the determinants computationally intensive, reminder of the Laplace rule:

$$\det(A) = \begin{cases} a_n & \text{if } n = 1 \\ \sum_{j=1}^n (-1)^{i+j} a_{ij} \Delta_{ij} & \forall i \quad \Delta_{ij} = \det(A_{ij}) \end{cases}$$

Where  $A_{ij}$  is the submatrix obtained by removing the  $i$ th row and  $j$ th column.

On average we perform a number of operations equal to  $3(n+1)!$ , absolutely unfeasible, we must move to numerical counterparts/approximations.

#### 3.2 Numerical Approximations: Direct Methods

We will talk about:

- 1) Direct methods (fixed number of steps to get the solution, no notion of convergence)
- 2) Iterative methods (non-ending number of steps, we will have to stop at a certain point, convergence)

There is no better alternative, both depend on the kind of the problem

At the basis of the direct methods we will use the factorization of the matrix: **LU factorization**.

##### 3.2.1 LU factorization of a matrix

Let  $A$  a matrix, it can be expressed as the product of two matrices  $A = LU$ , where  $L$  is a lower triangular matrix (elements different from zero on the main diagonal and on elements below) and  $U$  is an upper triangular matrix.

Consider

$$Ax = b \quad A \text{ nonsingular}$$

$$\underbrace{LU}_y x = b$$

The solution of the system then changes to

$$\begin{cases} Ly = b \\ Ux = y \end{cases}$$

Solve the upper system first, then the lower one

As the matrices are triangular (a lot of zeros, they are sparse), the systems are easier to solve! A matrix is defined as **sparse** if the number of entries different from zeros is a  $O(n)$  instead of  $O(n^2)$

A remark, if  $A$  is nonsingular, then it follows that  $L$  and  $U$  are nonsingular as

$$\det(A) = \det(LU) = \det(L) \det(U)$$

Also since the **determinant of a triangular matrix is the product of the entries in the diagonal**, all elements on the diagonal are non-zero.

### 3.2.2 Forward substitution method

Assume  $L$  and  $U$  are given, we start from:

$$Ly = b$$

$$\begin{bmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

$$l_{11}y_1 = b_1 \rightarrow y_1 = \frac{b_1}{l_{11}} \quad l_{11} \neq 0$$

$$l_{21}y_1 + l_{22}y_2 = b_2 \rightarrow y_2 = \frac{1}{l_{22}} [b_2 - l_{21}y_1]$$

$$l_{31}y_1 + l_{32}y_2 + l_{33}y_3 = b_3 \rightarrow y_3 = \frac{1}{l_{33}} (b_3 - l_{31}y_1 - l_{32}y_2)$$

So for a generic lower triangular matrix  $L$ :

$$y_1 = \frac{b_1}{l_{11}}$$

$$y_i = \frac{1}{l_{ii}} \left[ b_i - \sum_{j=1}^{i-1} l_{ij}y_j \right] \quad i = 2, \dots, n$$

With a number of operations for each  $i$ :

- 1 division for  $l_{ii}$
- In the squares  $i - 1$  subtractions
- Each subtractions is a product, we have  $i - 1$  multiplications

So the total number of operations is:

$$1 + \sum_{i=2}^n (1 + 2i - 2) = \sum_{i=1}^n (1 + 2i - 1) = \sum_{i=1}^n 1 + 2 \sum_{i=1}^n (i - 1) = n + 2 \frac{n(n-1)}{2} = \mathbf{n^2}$$

### 3.2.3 Backward substitution method

Consider now the upper triangular system:

$$Ux = b$$

$$\begin{bmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$u_{33}x_3 = y_3 \rightarrow x_3 = \frac{y_3}{u_{33}}$$

$$\dots \rightarrow x_2 = \frac{1}{u_{22}} [y_2 - u_{23}x_3]$$

...

So for a generic lower upper matrix  $U$ :

$$x_n = \frac{y_n}{u_{nn}}$$

$$x_i = \frac{1}{u_{ii}} \left[ y_i - \sum_{j=i+1}^n u_{ij}x_j \right] \quad i = n-1, \dots, 1$$

So the total number of operations is  $n^2$

### 3.2.4 Direct inspection

Suppose we know the matrix, we want to find the  $LU$  factorization:

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} = \underbrace{\begin{bmatrix} l_{11} & 0 \\ l_{21} & l_{22} \end{bmatrix}}_L \underbrace{\begin{bmatrix} u_{11} & u_{12} \\ u_{21} & 0 \end{bmatrix}}_U$$

In this case, we can write:

$$\begin{cases} l_{11}u_{11} = a_{11} \\ l_{11}u_{12} = a_{12} \\ l_{21}u_{11} = a_{21} \\ l_{21}u_{12} + l_{22}u_{22} = a_{22} \end{cases}$$

We can see that we have 6 unknowns and 4 equations. By convention we can assign "1" to the diagonal of the matrix  $L$

$$\underbrace{\begin{bmatrix} 1 & 0 \\ l_{21} & 1 \end{bmatrix}}_L \underbrace{\begin{bmatrix} u_{11} & u_{12} \\ u_{21} & 0 \end{bmatrix}}_U$$

So our system becomes:

$$\begin{cases} u_{11} = a_{11} \\ u_{12} = a_{12} \\ l_{21}u_{11} = a_{21} \\ l_{21}u_{12} + u_{22} = a_{22} \end{cases}$$

From which we easily find the unknowns.



If we consider a generic matrix  $N$ :

$$\underbrace{\begin{bmatrix} \ddots & & \ddots \\ & A & \\ \ddots & & \ddots \end{bmatrix}}_{n^2 \text{ equations}} = \underbrace{\begin{bmatrix} \ddots & & \\ & L & \\ & & \ddots \end{bmatrix}}_{\text{With } \frac{n(n+1)}{2} \text{ unknowns}} \underbrace{\begin{bmatrix} \ddots & & \\ & U & \\ & & \ddots \end{bmatrix}}_{\text{With } \frac{n(n+1)}{2} \text{ unknowns}}$$

Globally with  $n^2 + n$  unknowns. Just like before, we assign "1" to the diagonal of  $L$  so #equations=#unknowns

### 3.2.5 Gaussian elimination method

**GEM**, less computationally intensive than direct inspection. Consider the matrix

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

We assign for each entry a superindex to underline that it's the original matrix:

$$A^{(1)} = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ a_{21}^{(1)} & a_{22}^{(1)} & a_{23}^{(1)} \\ a_{31}^{(1)} & a_{32}^{(1)} & a_{33}^{(1)} \end{bmatrix}$$

1. First step is to move to the matrix  $A^{(2)}$ . Consider the coefficients:

$$l_{21} = \frac{a_{21}^{(1)}}{a_{11}^{(1)}} \quad \text{pivot } a_{11}^{(1)} \neq 0$$

$$l_{31} = \frac{a_{31}^{(1)}}{a_{11}^{(1)}} \quad \text{pivot } a_{11}^{(1)} \neq 0$$

$$A^{(2)} = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} \\ 0 & a_{32}^{(2)} & a_{33}^{(2)} \end{bmatrix}$$

Where

- $R2_{new} = R2_{old} - l_{21}R1_{old}$  and

$$a_{21}^{(2)} = a_{21}^{(1)} - l_{21}a_{11}^{(1)} = a_{21}^{(1)} - \frac{a_{21}^{(1)}}{a_{11}^{(1)}}a_{11}^{(1)} = 0$$

$$a_{22}^{(2)} = a_{22}^{(1)} - l_{21}a_{12}^{(1)}$$

$$a_{23}^{(2)} = a_{23}^{(1)} - l_{21}a_{13}^{(1)}$$

- $R3_{new} = R3_{old} - l_{31}R1_{old}$  and

$$a_{31}^{(2)} = a_{31}^{(1)} - l_{31}a_{11}^{(1)} = a_{31}^{(1)} - \frac{a_{31}^{(1)}}{a_{11}^{(1)}}a_{11}^{(1)} = 0$$

$$a_{32}^{(2)} = a_{32}^{(1)} - l_{31}a_{12}^{(1)}$$

$$a_{33}^{(2)} = a_{33}^{(1)} - l_{31}a_{13}^{(1)}$$

2. Similarly find  $A^{(3)}$ . Consider the coefficient:

$$l_{32} = \frac{a_{32}^{(2)}}{a_{22}^{(2)}} \quad \text{pivot } a_{22}^{(2)} \neq 0$$

$$A^{(3)} = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & a_{13}^{(1)} \\ 0 & a_{22}^{(2)} & a_{23}^{(2)} \\ 0 & 0 & a_{33}^{(3)} \end{bmatrix} = U$$

Where

$$\bullet R3_{new} = R3_{old} - l_{32}R2_{old} \text{ and}$$

$$a_{32}^{(3)} = a_{32}^{(2)} - l_{32}a_{22}^{(2)} = a_{32}^{(2)} - \frac{a_{32}^{(2)}}{a_{22}^{(2)}}a_{22}^{(2)} = 0$$

$$a_{33}^{(3)} = a_{33}^{(2)} - l_{32}a_{23}^{(2)}$$

Hence we found  $U$ . What about  $L$ ? It is:

$$\begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{bmatrix}$$

Let's see if  $LU = A$ , for example:

$$A = A^{(1)} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 0 & -1 \\ -1 & 1 & 5 \end{bmatrix}$$

1. Step 1:

$$A^{(2)} = \begin{bmatrix} 1 & 2 & 1 \\ 0 & -4 & -3 \\ 0 & 3 & 6 \end{bmatrix}$$

With

$$l_{21} = 2 \\ l_{31} = -1$$

2. Step 2:

$$A^{(3)} = \begin{bmatrix} 1 & 2 & 1 \\ 0 & -4 & -3 \\ 0 & 0 & \frac{15}{4} \end{bmatrix} = U$$

With

$$l_{32} = -\frac{3}{4}$$

And  $L$ :

$$L = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & -\frac{3}{4} & 1 \end{bmatrix}$$

We check that:

$$L \cdot U = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & -\frac{3}{4} & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 1 \\ 0 & -4 & -3 \\ 0 & 0 & \frac{15}{4} \end{bmatrix} = \begin{bmatrix} 1 & 2 & 1 \\ 2 & 0 & -1 \\ -1 & 1 & 5 \end{bmatrix} = A$$

**What about the cost? It is CUBIC:  $\frac{2}{3}n^3$ .**

When we want to solve  $Ax = b$  we have multiple ways:

1. Using GEM and LU factorization

- GEM  $\frac{2}{3}n^3$ .
- Resolution of

$$\begin{cases} Ly = b \\ Ux = y \end{cases}$$

Paying  $2n^2$

Paying in total  $\frac{2}{3}n^3 + 2n^2 \sim O(n^3)$

2. GEM on  $[A|b]$  paying  $(> \frac{2}{3}n^3) + 2n^2 \sim O(n^3)$

3. Using  $x = A^{-1}b$  paying an higher cost than the previous two methods. The inverse can be found applying GEM till:

$$[A|I_n] \Rightarrow [I_n|A^{-1}]$$

Consider methods 1 and 2, in which case one is the less costly than the other? When we have to solve multiple systems of the kind  $Ax = c_{i \rightarrow q}$  with right side different every time: in this case the first method is better as we pay the cost of  $LU$  factorization only once (so we pay  $\frac{2}{3}n^3 + q(2n^2)$ ) while with the second one we pay the whole cost every time (paying  $q(\frac{2}{3}n^3 + 2n^2)$ ).

**This is the case of method 3**, we are solving  $Ax = c_{i \rightarrow q}$  multiple times and in total we pay  $\frac{2}{3}n^3 + n(2n^2) = \frac{8}{3}n^3$

### 3.2.6 Necessary and sufficient criteria for GEM

$$A = A^{(1)} \begin{bmatrix} 1 & 1 & 3 \\ 2 & 2 & 2 \\ 3 & 6 & 4 \end{bmatrix} \rightarrow A^{(2)} = \begin{bmatrix} 1 & 1 & 3 \\ 0 & 0 & -4 \\ 0 & 3 & -5 \end{bmatrix}$$

With  $[l_{21} \ l_{31} \ l_{32}] = [2 \ 3 \ \frac{3}{0}]$ . The last pivot error! Some conditions:

- **Necessary and sufficient condition:** let  $A \in \mathbb{R}^{n \times n}$ , then  $LU$  factorization  $\exists$ ! **if and only if** the **principal submatrices**  $A_i$  of  $A$  for  $i = 1, \dots, n-1$  are nonsingular. The principal submatrices are:

$$A = \begin{bmatrix} A_1 & \vdots & \vdots & \vdots & \vdots \\ \cdots & A_2 & \vdots & \vdots & \vdots \\ \cdots & \cdots & A_3 & \vdots & \vdots \\ \cdots & \cdots & \cdots & \ddots & \vdots \\ \cdots & \cdots & \cdots & \cdots & A_{n-1} \\ \cdots & & & & \end{bmatrix}$$

As  $i$  goes to  $n-1$ , the matrix  $A_n = A$  can be singular, a  $LU$  exists anyway (but solution might not). This condition guarantess that a **LU factorization exists and it is unique**.

Some examples

- Singular matrix but  $LU$  factorization can be found

$$A_1 = \begin{bmatrix} 1 & 2 \\ 1 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 0 & 0 \end{bmatrix}$$

- We lose existence

$$A_2 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ l_{21} & 1 \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} \\ 0 & u_{22} \end{bmatrix}$$

As  $l_{21}u_{11} = 1$  is impossible to satisfy as  $u_{11} = 0$

- We lose uniqueness

$$A_3 = \begin{bmatrix} 0 & 1 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ l_{21} & 1 \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} \\ 0 & u_{22} \end{bmatrix}$$

As  $l_{21}u_{11} = 0$ , always satisfiable as  $u_{11} = 0$

- **Sufficient conditions for LU factorization existence and uniqueness**, 3 conditions that correspond to 3 families of lucky matrices

- 1)  $A$  belongs to the family of the strictly diagonally dominant by rows matrices, which means that the element of the diagonal is dominant w.r.t. the other elements in the same row: the absolute value of diagonal element is strictly greater than sum of absolute value of other elements:

$$|a_{ii}| > \sum_{j=1, j \neq i}^n |a_{ij}| \quad i = 1, \dots, n$$

An example that is not strictly diagonal by rows:

$$\begin{bmatrix} -4 & 0 & 3 \\ 1 & 2 & -4 \\ 7 & -1 & 10 \end{bmatrix}$$

The second row does not satisfy the condition, something that satisfies the condition would have second row as  $[1 \ -4 \ 2]$

- 2) In some sense the dual family of the previous case:  $A$  belongs to the family of the strictly diagonally dominant by columns matrices:

$$|a_{jj}| > \sum_{i=1, i \neq j}^n |a_{ij}| \quad j = 1, \dots, n$$

Strictly dominant by columns does not mean it is also strictly dominant by rows, and viceversa

- 3)  $A$  belongs to the family of symmetric positive definite (spd) matrices.

- A matrix is symmetric if:

$$a_{ij} = a_{ji} \quad i, j = 1, \dots, n$$

$$A = A^T$$

- A matrix is positive definite if:

$\forall x \neq 0 \in \mathbb{R}$  for each vector, we compute the scalar and verify that:

$$x^T A x > 0$$

But we perform an infinite number of checks, we consider the pair eigenvalues and eigenvectors:

$$A \in \mathbb{R}^{n \times n}$$

$$(\lambda, v)$$

$$Av = \lambda v$$

- A matrix is symmetric if all eigenvalues are real numbers
- A matrix is positive definite if all eigenvalues are positive  $\lambda > 0$

### 3.2.7 LU strategies for particular matrices

- If matrix  $A$  is spd, we consider the  $A = LU$ , how does the two factors inherit the symmetric aspect of the original matrix? As

$$A = LU = A^T = (LU)^T = U^T L^T$$

One is the transpose of the other, so we compute only one factor. In matlab Cholesky,  $R = \text{chol}(A)$ , and the computational cost from  $n^3$  will be reduced to  $\frac{n^3}{3}$ .

$$A = R^T R$$

Also the diagonal entries of matrix  $R$   $r_{ii} > 0$  will all be positive.

The Cholesky decomposition:

$$H = \begin{cases} h_{11} = \sqrt{a_{11}} \\ h_{ij} = \frac{1}{h_{jj}} \left( a_{ij} - \sum_{k=1}^{j-1} h_{ik} h_{jk} \right) & j = 1, \dots, i-1 \\ h_{ii} = \sqrt{a_{ii} - \sum_{k=1}^{i-1} h_{ik}^2} \end{cases}$$

Where

$$H^T = R$$

And our system becomes

$$\begin{cases} Ax = b \\ A = R^T R = HH^T \end{cases} \rightarrow R^T R x = b \rightarrow \begin{cases} R^T y = b \\ R x = y \end{cases} = \begin{cases} H y = b \\ H^T x = y \end{cases}$$

- Consider a 3-diagonal matrix, which has diagonal, first upper and first lower diagonals the only non-null overall, while all the others zero. The entries of those 3 diagonals are either 0,-1,+1

$$A = \begin{bmatrix} \ddots & \ddots & 0 & 0 & 0 \\ \ddots & \ddots & \ddots & 0 & 0 \\ 0 & \ddots & \ddots & \ddots & 0 \\ 0 & 0 & \ddots & \ddots & \ddots \\ 0 & 0 & 0 & \ddots & \ddots \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ \ddots & 1 & 0 & 0 & 0 \\ 0 & \ddots & 1 & 0 & 0 \\ 0 & 0 & \ddots & 1 & 0 \\ 0 & 0 & 0 & \ddots & 1 \end{bmatrix}}_{L \text{ lower bi-diagonal, with main diagonal all 1's}} \underbrace{\begin{bmatrix} \ddots & \ddots & 0 & 0 & 0 \\ 0 & \ddots & \ddots & 0 & 0 \\ 0 & 0 & \ddots & \ddots & 0 \\ 0 & 0 & 0 & \ddots & \ddots \\ 0 & 0 & 0 & 0 & \ddots \end{bmatrix}}_{U \text{ upper bi-diagonal}}$$

We to find the  $LU$  factorization, consider

$$A = \begin{bmatrix} a_1 & c_1 & 0 \\ e_2 & a_2 & c_2 \\ 0 & e_3 & a_3 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 0 & 0 \\ \beta_2 & 1 & 0 \\ 0 & \beta_3 & 1 \end{bmatrix}}_L \underbrace{\begin{bmatrix} \alpha_1 & \gamma_1 & 0 \\ 0 & \alpha_2 & \gamma_2 \\ 0 & 0 & \alpha_3 \end{bmatrix}}_U$$

It can be easily verified that

$$\begin{aligned} \gamma_1 &= c_1 & R1 * C2 \text{ first row of } L \text{ by second column of } U \\ \gamma_2 &= c_2 & R2 * C3 \end{aligned}$$

So we have 5 unknowns in  $\beta_i$  and  $\alpha_i$ :

$$\begin{aligned}
R1 * C1 & \quad \alpha_1 = a_1 \\
R2 * C1 & \quad \beta_2 \alpha_1 = e_2 \rightarrow \beta_2 = \frac{e_2}{\alpha_1} \\
R2 * C2 & \quad \beta_2 c_1 + \alpha_2 = a_2 \rightarrow \alpha_2 = a_2 - \beta_2 c_1 \\
R3 * C2 & \quad \beta_3 \alpha_2 = e_3 \rightarrow \beta_3 = \frac{e_3}{\alpha_2} \\
R3 * C3 & \quad \beta_3 c_2 + \alpha_3 = a_3 \rightarrow \alpha_3 = a_3 - \beta_3 c_2
\end{aligned}$$

So with a general  $n \times n$  matrix:

$$\begin{cases} \alpha_1 = a_1 \\ \beta_i = \frac{e_i}{\alpha_{i-1}} \\ \alpha_i = a_i - \beta_i c_{i-1} \\ i = 2, \dots, n \end{cases}$$

For every  $\beta_i$  we pay 1 division, for every  $\alpha_i$  we pay 1 multiplication and 1 difference, so in total we pay:

$$3(n-1)$$

Our problem is

$$Ax = b \rightarrow LUx = b \rightarrow \begin{cases} Ly = b \\ Ux = y \end{cases}$$

– The first equation:

$$\begin{bmatrix} 1 & 0 & 0 \\ \beta_2 & 1 & 0 \\ 0 & \beta_3 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

From which:

$$\begin{aligned}
y_1 &= b_1 \\
y_2 &= b_2 - \beta_2 y_1 \\
y_3 &= b_3 - \beta_3 y_2
\end{aligned}$$

Differently from forward substitution method:

- \* We don't have to subtract by a sum of elements
- \* There is no division term

And for a general  $n \times n$  matrix:

$$\begin{aligned}
y_1 &= b_1 \\
y_i &= b_i - \beta_i y_{i-1} \quad i = 2, \dots, n
\end{aligned}$$

With cost of  $2(n-1)$

– The second equation:

$$\begin{bmatrix} \alpha_1 & c_1 & 0 \\ 0 & \alpha_2 & c_2 \\ 0 & 0 & \alpha_3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

From which:

$$\begin{aligned}
x_3 &= \frac{y_3}{\alpha_3} \\
x_2 &= \frac{1}{\alpha_2} [y_2 - c_2 x_3] \\
x_1 &= \frac{1}{\alpha_1} [y_1 - c_1 x_2]
\end{aligned}$$

Differently from backward substitution method we don't have to subtract a summation term.  
For a general matrix:

$$x_n = \frac{y_n}{\alpha_n}$$

$$x_i = \frac{1}{\alpha_i} [y_i - c_i x_{i+1}]$$

With cost of  $3(n-1) + 1$

So in total we pay:

$$\underbrace{3(n-1)}_{\text{factorization}} + \underbrace{2(n-1)}_{Ly=b} + \underbrace{3(n-1)+1}_{Ux=y} = 8n-7$$

All this process is known as **Thomas algorithm**

### 3.2.8 Rows swapping: pivoting

Consider this matrix:

$$A = \begin{bmatrix} 1 & 1 & 3 \\ 2 & 2 & 2 \\ 3 & 6 & 4 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 1 & 3 \\ 0 & 0 & -4 \\ 0 & 3 & -5 \end{bmatrix}$$

We have  $a_{22}^{(2)} = 0$ , a pivot is nullable. Also if we verify the necessary and sufficient condition we see that the second principal matrix is singular. A solution is to swap the rows:

$$A = \begin{bmatrix} 1 & 1 & 3 \\ 3 & 6 & 4 \\ 2 & 2 & 2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 1 & 3 \\ 0 & 3 & -5 \\ 0 & 0 & -4 \end{bmatrix}$$

The necessary and sufficient condition is satisfied now. Remember, **the matrix must be nonsingular**.

**We perform the permutation when we meet the problem, otherwise we perform GEM all over again.** We exchange rows  $i$  with  $i+1$  if:

$$a_{ii}^{(i)} = 0 \ \&\& \ a_{i+1,i}^{(i)} \neq 0$$

To do such swap we need to define a **permutation matrix**  $P$ , an orthogonal matrix ( $PP^T = P^T P = I$ ), obtained from identity matrix by swapping some rows. Premultiplying it by our original matrix we perform the exchange, for example:

$$PA = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{31} & a_{32} & a_{33} \\ a_{21} & a_{22} & a_{23} \end{bmatrix}$$

**To exchange two rows, we just need to premultiply the original matrix by a permutation matrix. If we postmultiply  $AP$  we swap columns instead.**

This process is called **pivoting by rows**. So during GEM we keep track of swappings on the matrix  $P$ , so at the end we obtain  $LU$  factorization not for  $A$  but for  $PA$

$$PA = LU$$

And the solution

$$Ax = b \rightarrow PAx = Pb \rightarrow LUx = Pb$$

$$\begin{cases} Ly = Pb \\ Ux = y \end{cases}$$

Only on the top equation we have effect of the permutation matrix.

But what if the pivot is not null but very small?

$$a_{ii}^{(i)} = 2.5 \cdot 10^{-4}$$

At a certain point we will divide by it, amplifying the floating point error. Consider:

$$A = \begin{bmatrix} 1 & 1 + 0.5 \cdot 10^{-6} & 3 \\ 2 & 2 & 20 \\ 3 & 6 & 4 \end{bmatrix}$$

The necessary and sufficient condition holds, the determinant of the second principal matrix is not zero but very close to it, we can apply GEM. But the  $LU$  factorization might not be accurate. The idea is to keep the pivot as large as possible in the GEM.

In matlab we will use

$$[L, U, P] = \text{lu}(A);$$

### 3.3 Accuracy

This concerns iterative methods as well. Till now we saw that

$$LU + pivoting \rightarrow \text{accurate } LU \text{ factorization } A = LU$$

With  $A$  nonsingular. If this factorization is accurate, is the solution to the system  $Ax = b$  accurate as well? **No, an accurate  $LU$  factorization does not necessarily ensure to have an accurate solution.**

Consider the Hilbert matrix:

$$a_{ij} = \frac{1}{i+j-1}$$

$$\begin{bmatrix} 1 & 1/2 & 1/3 & 1/4 & \dots \\ 1/2 & 1/3 & 1/4 & \dots & \\ 1/3 & 1/4 & \dots & & \\ \dots & & & & \end{bmatrix}$$

It is a symmetric matrix and positive definite, so spd. Let  $A_n$  the Hilbert matrix of order  $n$ , consider the family of systems:

$$A_n x_n = b_n \quad b_n, x_n \in \mathbb{R}^n, A_n \in \mathbb{R}^{n \times n}$$

The exact solution is

$$b_n x_n = [1, 1, \dots, 1]^T$$

Consider the error:

$$R_n = P_n A_n - L_n U_n \quad (A)$$

With each  $n$  we associate:

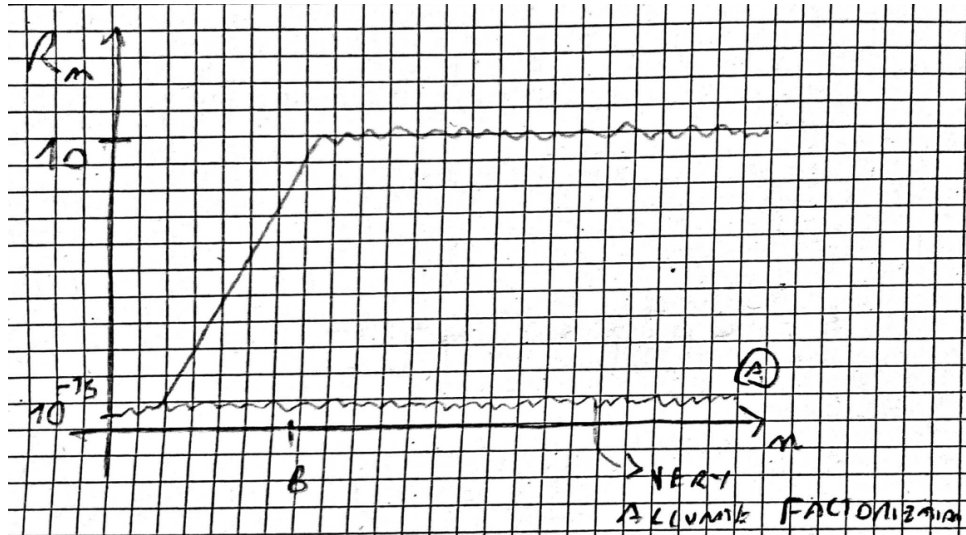
$$\max_{i,j} |r_{ij}|$$

Consider the approximation  $\tilde{x}_n$ , the relative error

$$E_n = \frac{\|x_n - \tilde{x}_n\|}{\|x_n\|} \quad (B)$$

With (B) is an accurate  $LU$  factorization, but the result is not accurate





The reason is the problem, the Hilbert matrix is **ill positioned problem**. The method does not solve  $Ax = b$ , but retrieves the solution of the **perturbed system**:

$$(A + \underbrace{\delta A}_{\mathbb{R}^{n \times n}})(x + \delta x) = b + \underbrace{\delta b}_{\mathbb{R}^n}$$

$\delta b$  and  $\delta A$  are perturbation on the data, matlab is computing the exact solution, but due to point floating point we will have a solution perturbation  $\delta x$ . We would like that for small perturbation on the data we get small perturbation on the problem (**well positioned problems**). Let the perturbed solution be

$$\tilde{x} = x + \delta x$$

Starting from  $\delta A = 0$ , we want to relate the two perturbations:

$$\frac{\|\delta A\|}{\|A\|} \quad \frac{\|\delta b\|}{\|b\|}$$

We must firstly introduce some lemmas.

### 3.3.1 Lemmas

- 1) Let  $B \in \mathbb{R}^{n \times n}$  be a spd matrix. Then it holds

$$\lambda_{\min}(B) \leq \frac{x^T B x}{x^T x} \leq \lambda_{\max}(B) \quad \forall x \neq 0 \in \mathbb{R}^n$$

There  $\lambda$ 's are the eigenvalues.

- 2) Let  $A$  be a nonsingular matrix. Then  $A^T A$  is spd

### 3.3.2 Perturbations relation

Starting from  $\delta A = 0$ , we want to relate the two perturbations:

$$\frac{\|\delta A\|}{\|A\|} \quad \frac{\|\delta b\|}{\|b\|}$$

- First step, we subtract term by term the exact system  $Ax = b$  from the perturbed one  $(A + 0)(x + \delta x) = b + \delta b$ :

$$A\delta x = \delta b$$

Consider the L2-norm (euclidean norm  $\|w\|^2 = w^T w$ ):

$$\|A\delta x\|^2 = \|\delta b\|^2 \rightarrow (A\delta x)^T (A\delta x) = \|\delta b\|^2 \rightarrow \delta x^T A^T A \delta x = \|\delta b\|^2$$

From the lemma 2  $A^T A$  is spd, so we can apply lemma 1 with  $B = A^T A$

$$\lambda_{\min}(A^T A) \delta x^T \delta x \leq \delta x^T A^T A \delta x = \|\delta b\|^2 \leq \lambda_{\max}(A^T A) \delta x^T \delta x$$

With  $\delta x^T \delta x = \|\delta x\|^2$ . Consider only the left part of the inequality:

$$\lambda_{\min}(A^T A) \|\delta x\|^2 \leq \|\delta b\|^2$$

$$\|\delta x\| \leq \frac{\|\delta b\|}{\sqrt{\lambda_{\min}(A^T A)}}$$

- As the next step consider

$$\|Ax\|^2 = \|b\|^2$$

Similar reasoning:

$$\lambda_{\min}(A^T A) x^T x \leq x^T A^T A x = (Ax)^T Ax = \|b\|^2 \leq \lambda_{\max}(A^T A) x^T x$$

With  $x^T x = \|x\|^2$ . Consider only the right part of the inequality:

$$\|b\| \leq \sqrt{\lambda_{\max}(A^T A)} \|x\|$$

$$\frac{1}{\|x\|} \leq \frac{\sqrt{\lambda_{\max}(A^T A)}}{\|b\|}$$

- We put everything together. As both inequalities are all positive quantities, we can combine them considering:

$$\begin{cases} a \leq b \\ c \leq d \end{cases} \rightarrow ab \leq cd$$

$$\frac{\|\delta x\|}{\|x\|} \leq \underbrace{\sqrt{\frac{\lambda_{\max}(A^T A)}{\lambda_{\min}(A^T A)}}}_{K(A) \geq 1 \text{ condition number}} \frac{\|\delta b\|}{\|b\|}$$

$K(A) \geq 1$  condition number

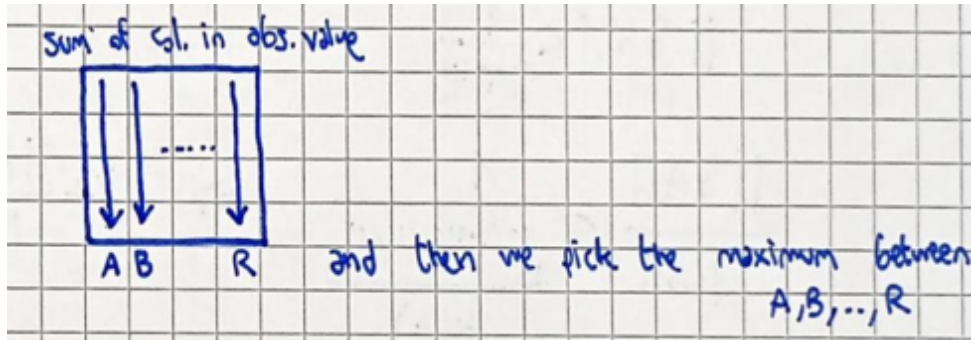
The condition number can only amplify, a small perturbation on the data will be small as well if the condition number is almost 1 (well conditioned problem), but if the condition number is very large this relationship does not hold (ill conditioned problem, just like for Hilbert)

### 3.3.3 Condition number

$$K(A) = \|A\| \cdot \|A^{-1}\|$$

A norm of a matrix has 3 possible definitions

- **Norm one**, maximum of the sum by columns



$$\|A\|_1 = \max_j \sum_{i=1}^n |a_{ij}|$$

$$K(A) = \|A\|_1 \cdot \|A^{-1}\|_1$$

- **Infinity norm**, first compress horizontally, then compress vertically, like dual of before

$$\|A\|_\infty = \max_i \sum_{j=1}^n |a_{ij}|$$

$$K(A) = \|A\|_\infty \cdot \|A^{-1}\|_\infty$$

- **2-norm or Spectral norm**, spectrum of a matrix is the collection of the eigenvalues

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)}$$

From before, we found that:

$$K(A) = \sqrt{\frac{\lambda_{\max}(A^T A)}{\lambda_{\min}(A^T A)}} = \|A\|_2 \|A^{-1}\|_2$$

So before we found a particular condition number:

$$K_2(A) = \|A\|_2 \|A^{-1}\|_2$$

- A special case, if  $A$  is symmetric ( $A = A^T$ ), so  $A^T A = A^2$

$$\lambda_{\max}(A^2) = [\lambda_{\max}(A)]^2$$

And

$$\|A\|_2 = \lambda_{\max}(A)$$

$$K_2(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$$

In matlab  $\text{cond}(A)$  for 2-norm,  $\text{condest}(A, \text{'parse'})$  for 1-norm

Now we can rewrite:

$$\delta b = A \underbrace{(x + \delta x)}_{\text{perturbation } \tilde{x}} - b = A\tilde{x} - b = -\tilde{r}$$

We get the residual

### 3.3.4 Perturbation on the matrix A

What if  $\delta A \neq 0$ ? If  $\|\delta A\| \|A^{-1}\| < 1$ , we can:

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{K(A)}{1 - K(A) \frac{\|\delta A\|}{\|A\|}} \left[ \frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right]$$

From the hypothesis

$$\|\delta A\| \cdot \|A^{-1}\| < 1 \rightarrow \|\delta A\| < \frac{1}{\|A^{-1}\|} \rightarrow \frac{\|\delta A\|}{\|A\|} < \frac{1}{\|A\| \cdot \|A^{-1}\|} \rightarrow K(A) \frac{\|\delta A\|}{\|A\|} < 1$$

Which means

$$1 - K(A) \frac{\|\delta A\|}{\|A\|} > 0$$

**In conclusion, before solving a problem check the condition number to see if it is ill conditioned or not. If condition number too huge problem**

## 4 Systems of Linear Equations: Iterative Methods

### 4.1 Iterative methods

We start from an initial guess  $x^{(0)} \in \mathbb{R}^n$ , that enters a black box that gives a  $x^{(1)}$ , and so on

$$\{x^{(k)}\} \quad x^{(k)} \simeq x \quad x^{(k)} \in \mathbb{R}^n$$

#### 4.1.1 Convergence

we will talk about convergence again

$$\lim_{k \rightarrow \infty} x^{(k)} = x$$

Or alternatively express convergence in terms of error

$$\begin{cases} e^{(k)} = x - x^{(k)} \\ \lim_{k \rightarrow \infty} e^{(k)} = 0 \end{cases}$$

#### 4.1.2 Consistency

In iterative methods we will always talk about **convergence, consistency and stability**.

Consistency:

$$x = Bx + g$$

We follow the recursive definition to create a new iteration:

$$x^{(k+1)} = Bx^{(k)} + g \quad B \in \mathbb{R}^{n \times n} \quad g \in \mathbb{R}^n$$

Where we have dependency of:

$$B = B(A) \quad g = g(A, b)$$

Consistency means that the equality has to hold w.r.t. the exact solution.

Notice that  $g$  depends on  $A$  as:

$$\begin{aligned} x &= Bx + g \\ (x - Bx) &= g \\ (I - B)x &= g \\ (I - B)A^{-1}b &= g \end{aligned}$$

#### 4.1.3 Convergence analysis

Subtract term by term

$$\begin{aligned} x^{(k+1)} &= Bx^{(k)} + g & - & \quad x = Bx + g \\ \underbrace{x - x^{(k)}}_{e^{(k+1)}} &= \underbrace{B(x - x^{(k)})}_{Be^{(k)}} \\ e^{(k+1)} &= Be^{(k)} \end{aligned}$$

Compatibility of vector norms says that:

$$\|e^{(k+1)}\| = \|Be^{(k)}\| \leq \|B\|_2 \|e^{(k)}\|$$

Use this inequality for convergence analysis

$$\|e^{(k)}\| \leq \|B\|_2 \|e^{(k-1)}\| \leq \|B\|_2^2 \|e^{(k-2)}\| \leq \dots \leq \|B\|_2^k \|e^{(0)}\|$$

Where the following quantity is called **spectral radius**

$$||B||_2^k = [\rho(B)]^k$$

Spectral radius is the maximum of the eigenvalue of the matrix taken in absolute value (in matlab  $\max(\text{abs}(\text{eig}(A)))$ ). A reminder, the 2-norm:

$$||B||_2 = \sqrt{\lambda_{\max}(B^T B)}$$

And if  $B$  is spd (eigenvalues are real positive)

$$||B||_2 = \lambda_{\max}(B)$$

Which means we are comparing  $\lambda_{\max}$  with itself, as in this case it is both spectral radius and 2-norm of the matrix! (**only if spd**)

To summarize:

$$||e^{(k)}|| \leq [\rho(B)]^k ||e^{(0)}||$$

And for convergence, the spectral radius must be less to 1

$$\begin{cases} \rho(B) < 1 \\ \lim_{k \rightarrow \infty} [\rho(B)]^k ||e^{(0)}|| = 0 \end{cases}$$

Attention, we **assumed that the method is consistent, make sure that the method is consistent**

This inequality represents a **necessary and sufficient condition for convergence**: let use consider the iterative scheme given by  $x^{(k+1)} = Bx^{(k)} + g$  and let us assume that such a scheme is **consistent**, then the scheme turns out to be convergent independently of the initial guess,  $\forall x^{(0)} \in \mathbb{R}^n$  **if and only if**  $\rho(B) < 1$

A scheme to be convergent, we must check consistency and spectral radius. If one of these does not hold, the method is not convergent.

Additionally, the smaller  $\rho(B)$ , the quicker the convergence.

We can also rewrite the inequality and get

$$\frac{||e^{(k)}||}{||e^{(0)}||} \leq [\rho(B)]^k \leq TOL$$

To find  $k$ =minimum number of iterations

## 4.2 Richardson schemes

We have to introduce

$$P \in \mathbb{R}^{n \times n} \quad \text{nonsingular matrix (preconditioner)}$$

$$\alpha_k \neq 0 \in \mathbb{R}$$

We rewrite the original system ( $Ax = b$ ) to

$$\alpha_k Ax = \alpha_k b$$

And rewrite matrix  $A$  with splitting

$$\alpha_k A = P - (P - \alpha_k A)$$

So:

$$\begin{aligned} P - (P - \alpha_k A)x &= \alpha_k b \\ Px &= \alpha_k b + (P - \alpha_k A)x \end{aligned}$$

We arbitrary choose the left side associated to  $k + 1$  and the right side associated to  $k$ :

$$Px^{(k+1)} = \alpha_k b + (P - \alpha_k A)x^{(k)}$$

Notice, by repliacing  $x^{(k)}$  with exact solution, we get consistency: **all the Richardson schemes are consistent by construction, so no need to check it, but we have to point it out!** We want to reach:

$$x^{(k+1)} = Bx^{(k)} + g$$

Multiplying by  $P^{-1}$  we will obtain:

$$x^{(k+1)} = \underbrace{\alpha_k P^{-1} b}_{g_{\alpha_k}} + \underbrace{(I - \alpha_k P^{-1} A)x^{(k)}}_{B_{\alpha_k}}$$

At this point we distinguish Richardson schemes to:

- Stationary, if  $\alpha_k = \alpha$ , which means it never changes at every iteration, a bad choice of  $\alpha$  may hinders convergence, it must be in suitable range
- Dynamic if  $\alpha_k$  changes at each iteration. The parameter can be tuned, it is better, but it's more computational intensive/demanding

$\alpha$  is the learning rate, accelerates convergence. Develop it:

$$x^{(k+1)} = \alpha_k P^{-1} ( \underbrace{b - Ax^{(k)}}_{r^{(k)} \text{ kth-residual}} ) + x^{(k)}$$

We can interpret it as "next iteration=previous+correction"

$$x^{(k+1)} = x^{(k)} + \alpha_k z^{(k)}$$

Where this value is the **preconditioned residual**

$$z^{(k)} = P^{-1} r^{(k)}$$

To compute the  $z^{(k)}$  just solve

$$Pz^{(k)} = r^{(k)} \leftrightarrow Ax = b$$

Notice that  $P$  is nonsingular, and it is **arbitrarily chosen**, so we choose a diagonal, 3 diagonal or triangular matrix. Also as  $P$  is always the same, use  $LU$  factorization just once! We save a lot of computations!

## 4.3 Stationary Richardson Schemes

### 4.3.1 Jacobi and Gauss-Seidel methods

- **Jacobi**, we start from an example

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 = b_1 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 = b_2 \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 = b_3 \end{cases}$$

For  $a_{11} \neq 0$

$$\begin{cases} x_1 = \frac{1}{a_{11}} [b_1 - a_{12}x_2 - a_{13}x_3] \\ x_2 = \frac{1}{a_{22}} [b_2 - a_{21}x_1 - a_{23}x_3] \\ x_3 = \frac{1}{a_{33}} [b_3 - a_{31}x_1 - a_{32}x_2] \end{cases}$$

So far no approximation, consider

$$\begin{cases} x_1^{(k+1)} = \frac{1}{a_{11}} [b_1 - a_{12}x_2^{(k)} - a_{13}x_3^{(k)}] \\ x_2^{(k+1)} = \frac{1}{a_{22}} [b_2 - a_{21}x_1^{(k)} - a_{23}x_3^{(k)}] \\ x_3^{(k+1)} = \frac{1}{a_{33}} [b_3 - a_{31}x_1^{(k)} - a_{32}x_2^{(k)}] \end{cases}$$

Then the generic case

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left[ b_i - \sum_{j=1, j \neq i}^n a_{ij}x_j^{(k)} \right]$$

$$i = 1, \dots, n \quad a_{ii} \neq 0$$

This method **can be parallelized**, each variable  $k+1$  depends only on variables at step  $k$ , strength point

- **Gauss-Seidel**, says to use the variables at the same step, it should converge/perform faster, but not guaranteed

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left[ b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)} \right]$$

$$i = 1, \dots, n \quad a_{ii} \neq 0$$

**Cannot be parallelized**

We now rewrite the schemes so we can reach Richardson-like form

$$x^{(k+1)} = \underbrace{\alpha_k P^{-1} b}_{g_{\alpha_k}} + \underbrace{(I - \alpha_k P^{-1} A)}_{B_{\alpha_k}} x^{(k)}$$

- **Jacobi**

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left[ b_i - \sum_{j=1, j \neq i}^n a_{ij}x_j^{(k)} \right]$$

$$D = \begin{bmatrix} a_{11} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & a_{nn} \end{bmatrix}$$

$$\vec{x}^{(k)} = [x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)}]^T$$



$$\vec{x}^{(k+1)} = [x_1^{(k+1)}, x_2^{(k+1)}, \dots, x_n^{(k+1)}]^T$$

$$b = [b_1, \dots, b_n]^T$$

So:

$$Dx^{(k+1)} = b - (A - D)x^{(k)}$$

$$x^{(k+1)} = x^{(k)} + D^{-1}(b - Ax^{(k)})$$

$$\begin{cases} x^{(k+1)} = x^{(k)} + 1 \cdot D^{-1}r^{(k)} \\ \alpha_{k_J} = 1, P_J = D \end{cases}$$

$$x^{(k+1)} = \underbrace{(I - D^{-1}A)}_{B_J} x^{(k)} + \underbrace{D^{-1}b}_{g_J}$$

• **Gauss-Seidel**

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left[ b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)} \right]$$

$$A = \begin{bmatrix} \ddots & & A - D + E \\ & D & \\ -E & & \ddots \end{bmatrix}$$

Where matrix  $E$  is strictly lower-triangular. Similar manipulations from before, we will get

$$x^{(k+1)} = x^{(k)} + (D - E)^{-1}(b - Ax^{(k)})$$

$$\begin{cases} x^{(k+1)} = x^{(k)} + (D - E)^{-1}r^{(k)} \\ \alpha_{k_{GS}} = 1, P_{GS} = (D - E) \end{cases}$$

$$x^{(k+1)} = \underbrace{(I - (D - E)^{-1}A)}_{B_{GS}} x^{(k)} + \underbrace{(D - E)^{-1}b}_{g_{GS}}$$

**Both stationary, so parameter  $\alpha$  not used to accelerate convergence, as it is constant.**

#### 4.3.2 Convergence for Jacobi and Gauss-Seidel

For Richardson schemes  $x^{(k+1)} = Bx^{(k)} + g$ , the requirements were for consistency and spectral radius  $\rho(B) < 1$  (necessary and sufficient). The consistency property for Richardson schemes is automatically guaranteed.

Additionally, we have other sufficient conditions

• **Jacobi**

- Necessary condition:  $\rho(B_J) < 1$
- Sufficient conditions (attention, referred to matrix **A**):

1. If  $A$  is strictly diagonally dominant by rows
2. If  $A$  is strictly diagonally dominant by columns

• **Gauss-Seidel**

- Necessary condition:  $\rho(B_{GS}) < 1$
- Sufficient conditions (attention, referred to matrix  $A$ ):
  1. If  $A$  is strictly diagonally dominant by rows
  2. If  $A$  is strictly diagonally dominant by columns
  3. If  $A$  is spd

We said that Gauss-Seidel does not necessarily outperforms Jacobi. A particular configuration

Let  $A \in \mathbb{R}^{n \times n}$  nonsingular, tridiagonal, with all  $a_{ii} \neq 0$ : both Jacobi and Gauss-Seidel are convergent or are divergent (they cannot be "discordant") and if they both are convergent

$$\rho(B_{GS}) = [\rho(B_J)]^2$$

The rate of convergent is determined by the spectral radius, which means GS converges faster, more specifically

$$\begin{aligned} \rho(B_J) = \frac{1}{4} \rightarrow \left(\frac{1}{4}\right) \leq \varepsilon \rightarrow k \geq \log_4\left(\frac{1}{\varepsilon}\right) \\ \rho(B_{GS}) = \frac{1}{4^2} \rightarrow \left(\frac{1}{4^2}\right) \leq \varepsilon \rightarrow 2k \geq \log_4\left(\frac{1}{\varepsilon}\right) \rightarrow k \geq \frac{1}{2} \log_4\left(\frac{1}{\varepsilon}\right) \\ \#iterations_{GS} \simeq \frac{\#iterations_J}{2} \end{aligned}$$

#### 4.3.3 Optimal acceleration parameter for stationary schemes

Let  $A$  and  $P$  be spd matrices. Then, the stationary Richardson scheme is convergent  $\forall x^{(0)} \in \mathbb{R}^{n \times n}$  if and only if (necessary and sufficient)

$$0 < \alpha < \frac{2}{\lambda_{\max}(P^{-1}A)}$$

Moreover

$$\alpha_{opt} = \frac{2}{\lambda_{\max}(P^{-1}A) + \lambda_{\min}(P^{-1}A)}$$

Finally

$$\underbrace{\|e^{(k)}\|_A}_{x-x^{(k)}} \leq \underbrace{\left(\frac{K(P^{-1}A) - 1}{K(P^{-1}A) + 1}\right)^k}_{< 1} \|e^{(0)}\|_A \quad k \geq 0$$

With  $K$  the condition number and  $\|e\|_A$  the energy norm:

$$\|w\|_A = \sqrt{w^T A w} \quad w \in \mathbb{R}^n$$

A proof for the necessary and sufficient condition:

- Prove that

$$0 < \alpha < \frac{2}{\lambda_{\max}(P^{-1}A)}$$

Consider  $B_\alpha = I - \alpha P^{-1}A$ , let  $\lambda_i$  eigenvalues of  $P^{-1}A$  such that

$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_n \geq 0$$

And let the generic eigenvalue of  $B_\alpha$ :

$$1 - \alpha\lambda_i$$

For convergence consistency holds since it is a Richardson scheme and  $\rho(B) < 1$  must hold as well, which means:

$$|1 - \alpha\lambda_i| < 1 \Rightarrow -1 < 1 - \alpha\lambda_i < 1 \Rightarrow \begin{cases} 1 - \alpha\lambda_i < 1 \\ 1 - \alpha\lambda_i > -1 \end{cases}$$

The first condition is true for  $\alpha > 0$ , about the second we will get  $\alpha < \frac{2}{\lambda_i} \forall i$  which will result in:

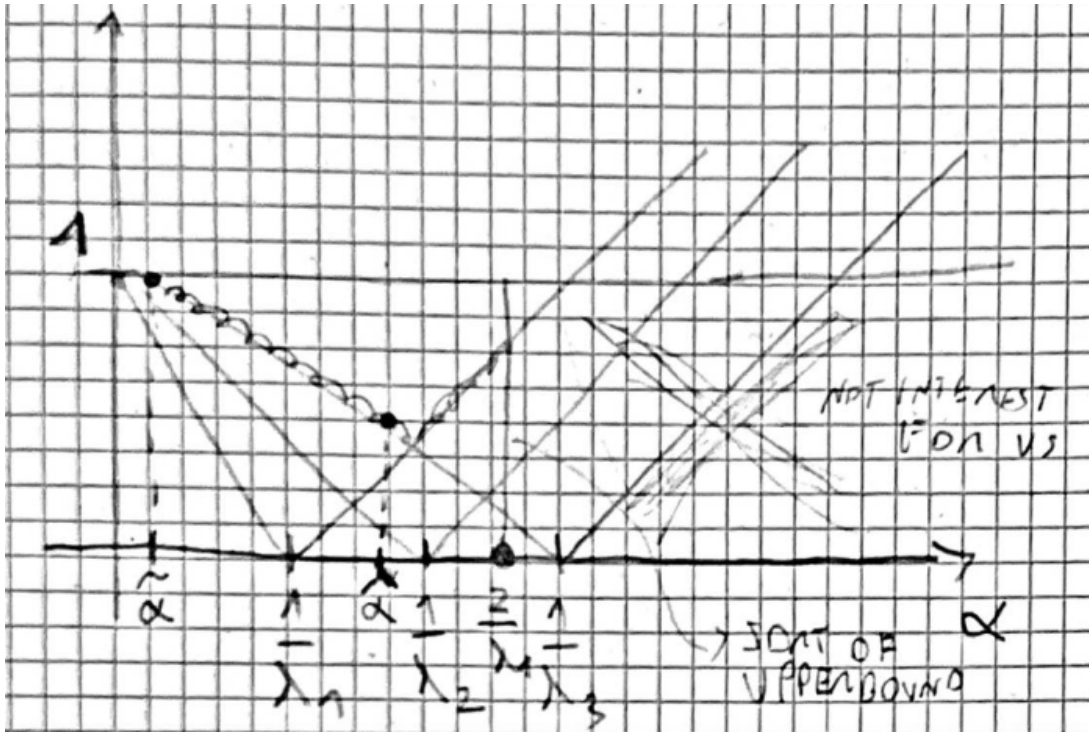
$$\frac{2}{\lambda_1} \leq \frac{2}{\lambda_2} \leq \dots \leq \frac{2}{\lambda_n} \Rightarrow \alpha < \frac{2}{\lambda_1} = \frac{2}{\lambda_{\max}(P^{-1}A)}$$

- Prove that  $\alpha_{opt}$  is that quantity. Consider a matrix  $A \in \mathbb{R}^{3 \times 3}$  with  $\lambda_1 > \lambda_2 > \lambda_3$ . Remind that

$$B_\alpha = I - \alpha P^{-1}A$$

$$|1 - \alpha\lambda_1| \quad |1 - \alpha\lambda_2| \quad |1 - \alpha\lambda_3|$$

Plotting them w.r.t.  $\alpha$ :



Analyzing those 3 plots we can understand why the optimal value is that: we are searching for  $\alpha_{opt}$  such that the method is as fast as possible. A method is quicker than the other if its spectral radius is smaller than the second, which means we are finding the value for  $\alpha$  such as that the maximum eigenvalue (which is related to the spectral radius) is as small as possible.

On the graph we draw a vertical line in a  $\tilde{\alpha}$  which will meet the 3 functions: we will get 3 eigenvalues, we consider the maximum.

Observing the plot, we can see that the maximum eigenvalue is on the branches highlighted, and the  $\alpha$  that minimizes the maximum eigenvalue is given by the intersection of the two branches

$$\underbrace{1 - \alpha\lambda_3}_{\text{Positive branch}} = \underbrace{\alpha\lambda_1 - 1}_{\text{Negative branch}} \Rightarrow \alpha_{opt} = \frac{2}{\lambda_1 + \lambda_3}$$

Which is exactly

$$\alpha_{opt} = \frac{2}{\lambda_{\max}(P^{-1}A) + \lambda_{\min}(P^{-1}A)}$$

About the maximum rate of convergence, the spectral radius (with  $\lambda_n$  the minimum):

$$\rho(B_{\alpha_{opt}}) = 1 - \alpha_{opt}\lambda_n = 1 - \frac{2\lambda_n}{\lambda_1 + \lambda_n} = \frac{\lambda_1 - \lambda_n}{\lambda_1 + \lambda_n}$$

- Prove that

$$\underbrace{\|e^{(k)}\|_A}_{x-x^{(k)}} \leq \underbrace{\left(\frac{K(P^{-1}A) - 1}{K(P^{-1}A) + 1}\right)^k}_{< 1} \|e^{(0)}\|_A \quad k \geq 0$$

With  $K$  the condition number and  $\|e\|_A$  the energy norm:

$$\|w\|_A = \sqrt{w^T A w} \quad w \in \mathbb{R}^n$$

We said that the preconditioner matrix  $P$  is:

- Nonsingular
- Easy to solve, as we have to solve  $Pz^{(k)} = r^{(k)}$
- **Additional condition**,  $K(P^{-1}A)$  small

#### 4.4 Dynamic Richardson Schemes

If  $A$  and  $P$  are spd matrices, the dynamic Richardson scheme is convergent if (sufficient)

$$\alpha_{k,opt} = \frac{[z^{(k)}]^T r^{(k)}}{[z^{(k)}]^T A z^{(k)}} \quad \forall k \geq 0, z^{(k)} = P^{-1}r^{(k)}$$

**We directly have optimal value for  $\alpha$ , this method is called preconditioned gradient method.** Moreover, we can exactly prove the same inequality

$$\underbrace{\|e^{(k)}\|_A}_{x-x^{(k)}} \leq \underbrace{\left(\frac{K(P^{-1}A) - 1}{K(P^{-1}A) + 1}\right)^k}_{< 1} \|e^{(0)}\|_A \quad k \geq 0$$

Remark: for  $P = I$ , our optimal recipe for  $\alpha$  becomes:

$$\alpha_{k,opt} = \frac{[r^{(k)}]^T r^{(k)}}{[r^{(k)}]^T A r^{(k)}} \quad \forall k \geq 0, z^{(k)} = P^{-1}r^{(k)}$$

#### Gradient method.

##### 4.4.1 The algorithm

The algorithm has 4 steps

```
// set initial guess, we can also define associated initial residual
 $x^{(\phi)} \in \mathbb{R}^n$      $r^{(\phi)} = b - Ax^{(\phi)}$ 
for k=0,1,...
    if  $P=I$ 
        // got to step 2, 3 steps algorithm in this case
    else
        1)  $Pz^{(k)} = r^{(k)}$ 
```

- 2)  $\alpha_{k,opt} = \frac{[z^{(k)}]^T r^{(k)}}{[z^{(k)}]^T A z^{(k)}}$
- 3)  $x^{(k+1)} = x^{(k)} + \alpha_{opt,k} z^{(k)}$
- 4)  $r^{(k+1)} = r^{(k)} - \alpha_{opt,k} A z^{(k)}$

If stationary, step 2) outside the for cycle

```
// set initial guess, we can also define associated initial residual
 $x^{(\phi)} \in \mathbb{R}^n$      $r^{(\phi)} = b - Ax^{(\phi)}$ 
 $\alpha_{opt} = \dots$ 
for k=0,1,...
    if P=I
        // got to step 3, 2 steps algorithm in this case
    else
        1)  $Pz^{(k)} = r^{(k)}$ 
        3)  $x^{(k+1)} = x^{(k)} + \alpha_{opt} z^{(k)}$ 
        4)  $r^{(k+1)} = r^{(k)} - \alpha_{opt} A z^{(k)}$ 
```

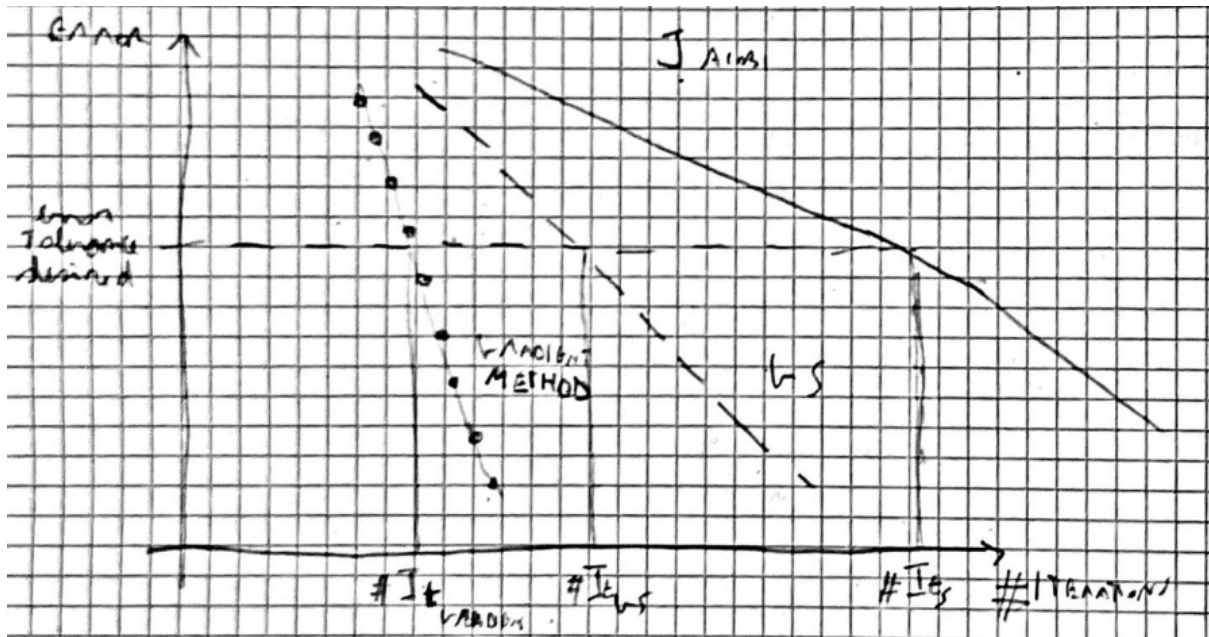
#### 4.4.2 Stationary and dynamic: which one is the best

Now that we have found the optimal acceleration parameters  $\alpha_{opt}$  for both stationary and dynamic cases, which one is better? For stationary we have to find the eigenvalues, but if the matrix is very big this is too demanding, though there are methods to cope with this.

It is a good idea to prefer dynamic scheme, as in the definition of the acceleration parameter we are using all parameters that are required to be computed for the solution of the system (in the stationary to find the eigenvalues we have to use other parameters unrelated and useless for the solution of our original problem).

An example

$$\begin{cases} 2x_1 + x_2 = 1 \\ x_1 = 3x_2 = 0 \end{cases}$$



#### 4.4.3 Proof for the optimal acceleration parameter

If  $A$  is spd, the solution to our system is:  $Ax = b \Leftrightarrow$  equivalent to minimize the quadratic form

$$Q(x) = \frac{1}{2}x^T Ax - x^T b$$

as minimizing this quadratic form is solving the gradient w.r.t. 0:

$$\nabla Q(x) = Ax - b = 0$$

Note that the quadratic  $Q(x)$  is a paraboloid and we are finding its minimum. How to proceed:

$$x^{(k+1)} = x^{(k)} + \gamma_k d^{(k)}$$

with  $d^{(k)}$  the direction, **the steepest descent** and  $\gamma_k$  the step size. We choose the steepest direction, so the gradient:

$$d^{(k)} = -\nabla Q(x^{(k)}) = b - Ax^{(k)} = r^{(k)}$$

The steepest direction corresponds to the residual. About the step

$$Q(x^{(k)} + \gamma_k r^{(k)}) = \tilde{Q}(\gamma_k)$$

Which is the value of  $\gamma_k$  that minimizes  $\tilde{Q}$ ? Compute the derivative and impose it to 0

$$\frac{d\tilde{Q}}{d\gamma_k} = 0$$

$$\tilde{Q}(\gamma_k) = Q(x^{(k)} + \gamma_k r^{(k)}) = \frac{1}{2} (x^{(k)} + \gamma_k r^{(k)})^T A (x^{(k)} + \gamma_k r^{(k)}) - (x^{(k)} + \gamma_k r^{(k)})^T b$$

Compute the derivative

$$\frac{d\tilde{Q}}{d\gamma_k} = [r^{(k)}]^T Ax^{(k)} + \gamma_k [r^{(k)}]^T Ar^{(k)} - [r^{(k)}]^T b = 0$$

$$\gamma_k = \frac{[r^{(k)}]^T (b - Ax^{(k)})}{[r^{(k)}]^T Ar^{(k)}} = \frac{[r^{(k)}]^T r^{(k)}}{[r^{(k)}]^T Ar^{(k)}}$$

We found  $\alpha_k = P\gamma_k$

$$Pz^{(k)} = r^{(k)}$$

$$\alpha^{(k)} = \frac{[z^{(k)}]^T r^{(k)}}{[z^{(k)}]^T Az^{(k)}}$$

#### 4.5 Conjugate gradient method

Gradient method can work with Hilbert system: conjugate gradient method, 5 steps algorithm that selects a new direction  $p^{(k)}$  instead of  $d^{(k)}$ :

$$[p^{(j)}]^T Ap^{(k+1)} = 0 \quad j = 0, \dots, k$$

New direction  $A$ -orthogonal (or  $A$  conjugate) w.r.t. the previous direction.

```
// set initial guess, we can also define associated initial residual
x(0) ∈ ℝ^n      r(0) = b - Ax(0)
for k=0, 1, ...
    1) α_k = ...
    2) x^(k+1) = x^(k) + α_k p^(k)
    3) r^(k+1) = ...
    4) β_k = ... // another constant for computing new direction
    5) p^(k+1) = r^(k+1) - β_k p^(k)
```

This method wants  $A$  spd, also the errors:

$$\|e^{(k)}\|_A \leq C^k \|e^{(0)}\|_A$$

$$C := C \left( \sqrt{K(P^{-1}A)} \right) \quad \text{Depends on sqrt root of condition number of } A$$

Consider the Hilbert problem/system, which we remind has a very bad condition number:

$$H_n x_n = b_n$$

n	$K(A_n)$	\	PG	P=D	PCG	P=D
4		$O(10^{-13})$	$O(10^{-3})$	995	$O(10^{-2})$	3
6					$O(10^{-2})$	4
8						4
10						5
12						5
14		$O(10)$	$O(10^{-3})$	1379	$O(10^{-3})$	5

If we can work in an exact arithmetic, this method becomes a direct method.

## 4.6 Stopping Criteria

Consider the error estimators: increment and residual

### 4.6.1 Residual

$$Ax = b$$

$$S = r^{(k)} = b - Ax^{(k)}$$

$$e^{(k)} = x - x^{(k)}$$

We want to relate the estimator

$$\|e^{(k)}\| \quad S$$

Normalize residual

$$\frac{\|e^{(k)}\|}{\|x\|} \leq C \frac{\|r^{(k)}\|}{\|b\|} \leq TOL \quad x \neq 0, b \neq 0$$

The constant  $C$  discriminates the reliability or not of our estimator: if it's huge, our estimator is not reliable.

We want to stop at the minimum iteration  $kmin$

$$\frac{\|e^{(kmin)}\|}{\|x\|} \leq C \underbrace{\frac{\|r^{(kmin)}\|}{\|b\|}}_{\text{Prove this}} \leq TOL$$

Remind that

$$\frac{\|\delta x = x - \tilde{x}\|}{\|x\|} \leq K(A) \frac{\|r = b - A\tilde{x}\|}{\|b\|}$$

If we consider  $\tilde{x} = x^{(k)}$ , we have that:

$$\delta x = e^{(k)} \quad r = r^{(k)}$$

Which means

$$C = K(A)$$

So regarding the reliability, we have to look at the condition number

#### 4.6.2 Increment

$$Ax = b$$

$$S = \delta^{(k)} = x^{(k+1)} - x^{(k)}$$

$$kmin \in \mathbb{N} \text{ s.t. } \|e^{(kmin)}\| \leq C \underbrace{\|\delta x^{(k)}\|}_{\text{Prove this}} \leq TOL$$

We start from the below and using triangular inequality:

$$\|e^{(k)}\| = \|x - x^{(k)}\| = \|\underbrace{x - x^{(k+1)}}_{e^{(k+1)}} + \underbrace{x^{(k+1)} - x^{(k)}}_{\delta^{(k)}}\| \leq \|e^{(k+1)}\| + \|\delta^{(k)}\| \leq \rho(B)\|e^{(k)}\| + \|\delta^{(k)}\|$$

$$\|e^{(k)}\| \leq \underbrace{\frac{1}{1 - \rho(B)}}_C \|\delta^{(k)}\|$$

As we want  $C$  as close as possible to 1, we want the spectral radius as close as possible to 0. So increment is a reliable estimator if the spectral radius is very close to 0.

About  $kmin$  (minimum number of iterations):

$$e^k = x - x^k$$

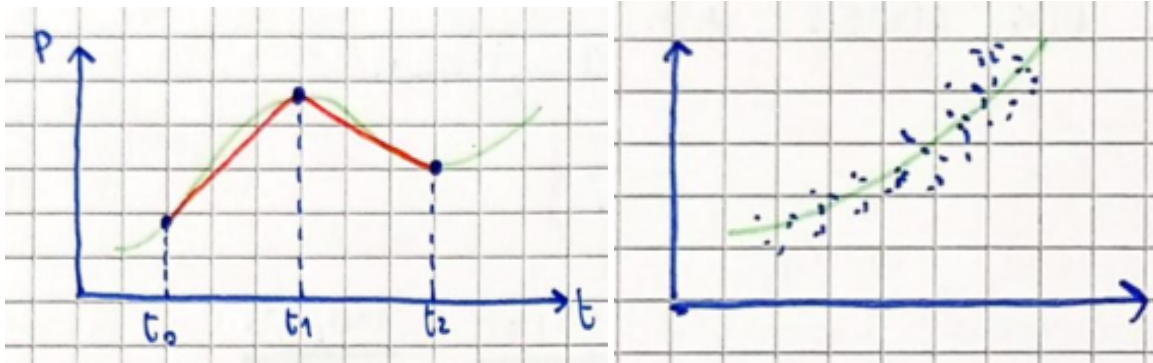
$$kmin = \frac{\log \left[ \frac{\varepsilon(1 - \|B\|_2)}{\log \|x^1 - x^0\|} \right]}{\log \|B\|_2}$$

With  $\varepsilon = TOL$



## 5 Approximation of Functions and Data

Approximation of data:



Left interpolation, right least square.

Approximation of functions like integral: approximate it to polynomial. The analytical tool to approximate a function is the Taylor expansion, but it suffers some problems

- Need higher order (more derivatives)
- It works well when we consider a neighborhood of the center of the expansion, but moving out of the neighborhood it will start to perform badly

### 5.1 Interpolation

Instead of Taylor, consider:

$$\left[ \begin{array}{cc} \text{Function} & \text{Data} \\ f & \{(x_i, y_i)\} \quad i = 0, \dots, n \\ y_i = f(x) & x_i \text{ distinct} \end{array} \right]$$

Identify a function  $\tilde{f}$  s.t.  $\tilde{f}(x_i) = y_i$  for  $i = 0, \dots, n$ , **interpolation conditions**,  $n + 1$  conditions. The function  $\tilde{f}$  can be:

- Polynomial
- Trigonometric expansion (Fourier expansion)
- Rational

$$\frac{a_0 + a_1x + \dots + a_px^p}{b_0 + b_1x + \dots + b_qx^q}$$

We consider polynomial interpolation

#### 5.1.1 Polynomial interpolation

Theorem: let  $\{x_i, y_i\}_{i=0}^n$  ( $x_i$  are interpolation nodes,  $y_i$  are values to be interpolated) with  $x_i$  all distinct. Then  $\exists!$  polynomial degree of  $\leq n$  (interpolating/interpolation polynomial) s.t. (we guarantee interpolation condition):

$$\underbrace{\Pi_n(x_i) = y_i}_{n+1 \text{ conditions}} \quad i = 0, \dots, n$$

Proof for uniqueness: by contradiction assume that we have 2 interpolating polynomials of order  $n$

$$\Pi_n \in \mathbb{P}_n \quad \Pi_n(x_i) = y_i \quad i = 0, \dots, n$$

$$\Pi_n^* \in \mathbb{P}_n \quad \Pi_n^*(x_i) = y_i \quad i = 0, \dots, n$$

Consider the difference

$$D(x) = \Pi_n(x) - \Pi_n^*(x) \in \mathbb{P}_n$$

And

$$D(x_i) = \Pi_n(x_i) - \Pi_n^*(x_i) = y_i - y_i = 0 \quad i = 0, \dots, n$$

A polynomial of degree  $n$  has at most  $n$  intersections with the  $x$ -axis, but in this case  $D(x)$  of degree  $n$  has  $n + 1$  zeros: the only way that we can satisfy the  $n + 1$  conditions is that  $D(x)$  is identically equal to 0. which means

$$D(x) = 0 \Rightarrow \Pi_n(x) - \Pi_n^*(x) = 0 \Rightarrow \Pi_n(x) = \Pi_n^*(x)$$

Which contradicts our initial assumption.

Finding the characteristic polynomial: assume that values to be interpolated are all null except one

$$y_i = 0 \quad \forall i \neq k \quad y_k = 1$$

$$\left\{ \begin{array}{ccc} x_0 = 0 & x_1 = 0.5 & x_2 = 1 \\ y_0 = 0 & y_1 = 1 & y_2 = 0 \end{array} \right\}$$

Let change notation, instead of  $\Pi_2$  we use  $\phi_k$

$$\phi_1 \in \mathbb{P}_2 \quad \underbrace{\phi_1(x_0 = 0) = 0}_A \quad \underbrace{\phi_1(x_1 = 1/2) = 1}_B \quad \underbrace{\phi_1(x_2 = 1) = 0}_C$$

We want to build a polynomial of degree 2 that is zero in those two points:

$$\left\{ \begin{array}{ll} (x-0)(x-1) & \text{satisfies A and B} \\ \frac{(x-0)(x-1)}{(0.5-0)(0.5-1)} & \text{satisfies C} \end{array} \right.$$

$$\phi_1(x) = \dots = -4x(x-1)$$

With a generic case

$$\phi_k(x_i) = \delta_{ik} = \begin{cases} 0 & i \neq k \\ 1 & i = k \end{cases}$$

Kronecker delta, to express it as a polynomial with degree  $n$ , the characteristic polynomial:

$$\phi_k(x) = \prod_{j=0, j \neq k}^n \frac{x - x_j}{x_k - x_j}$$

Moving to a more general case: instead of a set of arbitrary values

$$\left\{ \begin{array}{ccc} x_0 & x_1 & x_2 \\ y_0 & y_1 & y_2 \\ \Pi(x_0) = y_0 & \Pi(x_1) = y_1 & \Pi(x_2) = y_2 \end{array} \right\}$$

Expressing  $\Pi_2$  as linear combination of  $\phi_0$ ,  $\phi_1$  and  $\phi_2$

$$\left\{ \begin{array}{cccc} \phi_0 & \phi_0(x_0) = 1 & \phi_0(x_1) = 0 & \phi_0(x_2) = 0 \\ \phi_1 & \phi_1(x_0) = 0 & \phi_1(x_1) = 1 & \phi_1(x_2) = 0 \\ \phi_2 & \phi_2(x_0) = 0 & \phi_2(x_1) = 0 & \phi_2(x_2) = 1 \end{array} \right\}$$

$$\Pi_2(x) = a\phi_0(x) + b\phi_1(x) + c\phi_2(x)$$

Solving we will get:

$$a = y_0 \quad b = y_1 \quad c = y_2$$

The **Lagrange form**:

$$\Pi_n(x) = \sum_{k=0}^n y_k \phi_k(x)$$

$$\Pi_n(x) = \sum_{k=0}^n y_k \prod_{j=0, j \neq k}^n \frac{x - x_j}{x_k - x_j}$$

In matlab:

- **c = polyfit(x,y,n)**, to build interpolating polynomial
  - x vector collecting interpolation nodes,  $x_i$
  - y is  $y_i$
  - n is the degree of the polynomial, but is redundant as there is a strict relation between number of data and degree of polynomial (for  $n$  data, the polynomial will have degree of  $n - 1$ ), in least squares it will have a meaning

It returns the coefficients of our interpolating polynomial

$$p_n(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$$

With c(1) coefficient of  $x^n$ , c(2) of  $x(n-1)$

- **d = polyval(c,z)**, to evaluate interpolating polynomial at point  $z$ 
  - If  $z$  single number  $\mathbb{R}$ , d will be a number
  - If  $z$  is  $\mathbb{R}^q$ , vector

### 5.1.2 Interpolation error

Error at nodes is null, at points that are not nodes? Consider a function continuous in a certain interval  $I$ :

$$f \in C^0(\bar{I}) \quad I(x_0, x_n)$$

$$\{(x_i, y_i = f(x_i))\}_{i=0}^n \quad x_i \text{ distinct}$$

Assuming

$$f \in C^{n+1}(\bar{I})$$

We define the interpolation error:

$$\forall x \in I$$

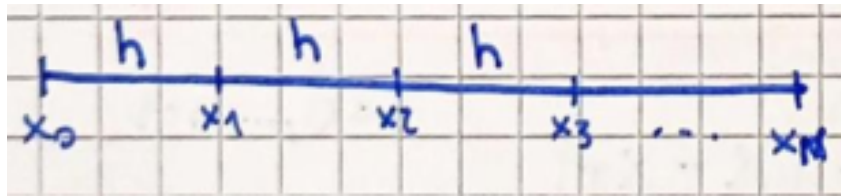
$$E_n f(x) = f(x) - \Pi_n f(x) = \frac{f^{(n+1)}(\alpha(x))}{(n+1)!} \prod_{k=0}^n (x - x_k)$$

And as expected  $E_n f(x_i) = 0$ . The weak point is that we assume regularity in the function, which depends on the number of nodes: such regularity uncommon. Another drawback is that the  $\alpha(x)$  depends on  $x$  but we don't know the exact value of  $\alpha(x) \in I$ . In practice this result useless, so we consider the maximum value.

If we have more and more information, more samples, the degree of polynomial increases and we have more zeros (the function meets the x axis more times, like sinusoid) but the quality of the approximation

improves.

Consider an uniform partition of the interval



$$h = \frac{x_n - x_0}{n}$$

$$x_k = x_{k+1} + h \quad k = 1, \dots, n$$

$$x_j = x_0 + jh \quad j = 0, \dots, n$$

In this case we can prove:

$$\left| \prod_{k=0}^n (x - x_k) \right| \leq n! \frac{h^{n+1}}{4}$$

Therefore:

$$\max_{x \in I} |E_n f(x)| \leq \underbrace{\frac{h^{n+1}}{4(n+1)}}_A \cdot \underbrace{\max_{x \in I} |f^{(n+1)}(x)|}_B$$

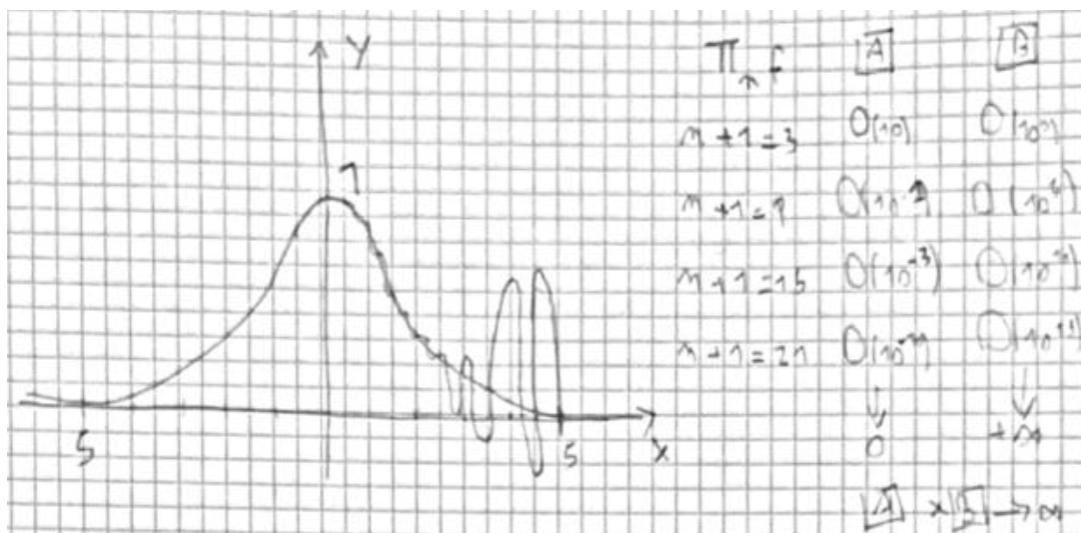
The two blocks:

- $A \rightarrow 0$  for  $n \rightarrow \infty$
- $B$  for  $n \rightarrow \infty$  depends on  $f$ :

$$B \rightarrow \begin{cases} 0 \\ \text{constant} \\ +\infty \begin{cases} \text{If A goes to 0 quicker, OK} \\ \text{If B goes to } \infty \text{ quicker, not OK} \end{cases} \end{cases}$$

Consider infact the function

$$f(x) = \frac{1}{1+x^2} \quad x \in I = [-5, 5]$$



We are in presence of spurious oscillations, in the case this error goes to  $\infty$ , we have the **Runge phenomenon: more samples, possibility of spurious oscillations.**

To deal with this we can:

- Particular choice of nodes, Chebyshev nodes, as we chose uniform sampling, uniform distribution of nodes, not good idea: with Chebyshev we choose more nodes near the endpoints where the Runge phenomenon occurs, while at the center of the interval less nodes. To divide into  $n$  parts:

1.  $n$  equal parts first

$$\frac{\Pi_i}{n} \quad i = 0, \dots, n$$

2. Compute the projection on the x-axis, the minus sign is in order to fix the order

$$-f\left(\frac{\Pi_i}{n}\right)$$

3. To use Chebyshev, use the following function that maps the interval  $I$  to the generic interval  $[a, b]$ , maps the nodes of previous step

$$x_i = \frac{a+b}{2} + \frac{b-a}{2} \hat{x}_i$$

- When we increase  $n$ , a lot of x-axis crosses, so we can use low degree polynomials and work interval by interval: **piecewise linear interpolation**

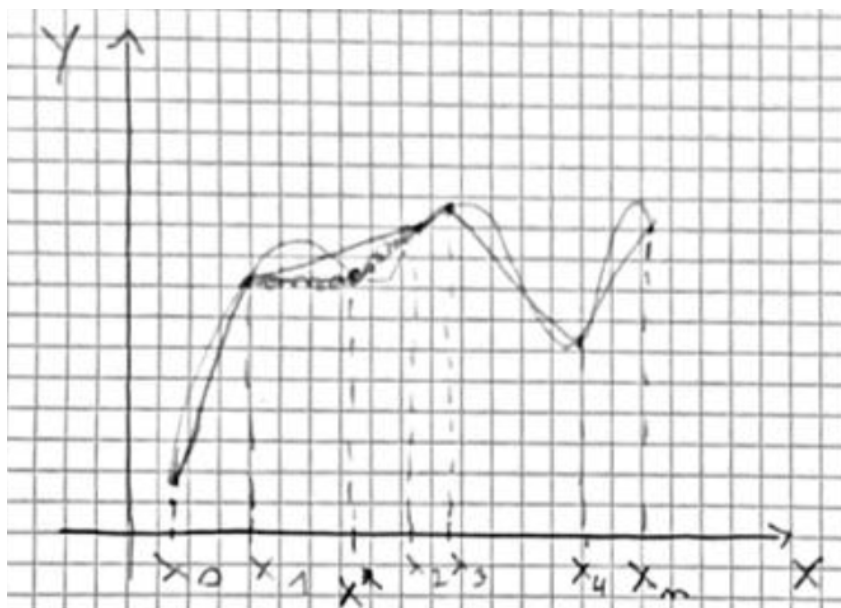
### 5.1.3 Piecewise linear interpolation

We do not have to select an uniform distribution of nodes

$$I_j = [x_j, x_{j+1}] \quad h_j = x_{j+1} - x_j$$

$$H = \max_j h_j$$

The idea is consider each subinterval and replace the function  $f$  with conjunction of endpoints (polynomial with low degree, so we avoid oscillations). By increasing the number of nodes, the approximation improves



We indicate linear interpolant as:

$$\begin{aligned}\Pi_1^H f &\in C^0(\bar{I}) \\ \Pi_1^H f|_{I_j} &\in \mathbb{P}_1(I_j) \\ \Pi_i^H f(x_k) &= f(x_k) \quad k = 0, \dots, n\end{aligned}$$

Express it as conjunction of endpoints that includes those 3 requirements:

$$\Pi_1^H f(x) = f(x_j) + \frac{f(x_{j+1}) - f(x_j)}{x_{j+1} - x_j}(x - x_j) \quad x \in I_j$$

How to prove that the error becomes 0 for  $n \rightarrow \infty$ ? We consider the local error on the subintervals  $I_j$ , we have (for  $n = 1$ ):

$$\max_{x \in I_j} |E_1^H f(x)| \leq \frac{h_j^2}{4 \cdot (1 + 1)} \cdot \max_{x \in I_j} |f''(x)| \quad f \in C^2(\bar{I}_j)$$

Instead of  $I_j$  and  $h_j$ , we want inequality on all  $I$ :

$$\max_{x \in I} |E_1^H f(x)| \leq \frac{H^2}{8} \cdot \max_{x \in I} |f''(x)| \quad f \in C^2(\bar{I})$$

Now we have the error that is decreasing:

- $A = \frac{H^2}{8}$ , for  $n \rightarrow \infty$ , the maximum length  $H$  will become smaller and smaller as we are increasing samples!
- $B = \max_{x \in I} |f''(x)|$  does not change, it's not dependent on  $n$

We can also use piecewise parabola or cubic interpolation, when we join pieces of parabola is the approximation more regular? Is it a  $C^1$  or still  $C^0$ ? Not  $C^1$ , so by increasing local degree of polynomial we do not gain regularity, but we are improving accuracy of approximation. For example  $n = 2$

$$\max_{x \in I_j} |E_2^H f(x)| \leq \frac{h_j^3}{4 \cdot (2 + 1)} \cdot \max_{x \in I_j} |f^{(3)}(x)| \quad f \in C^3(\bar{I}_j)$$

$$\max_{x \in I} |E_1^H f(x)| \leq \frac{H^3}{12} \cdot \max_{x \in I} |f^{(3)}(x)| \quad f \in C^3(\bar{I})$$

This is a problem, we would like something that is globally very smooth.

In matlab:

- **d = interp1(x,y,z)**, in some sense merges both polyfit and polyval, output will have same dimension as  $z$
- **d = interp2(x,y,z)**, cubic interpolation

A remark: the matlab function plot is doing a sampling of the function, which is finer in the gradient of the function. This is known as **adaptive sampling**

### 5.1.4 Cubic spline interpolation

Again, a piecewise interpolation, but we join endpoints with a cubic function

1.  $S_3|_{I_j} \in \mathbb{P}$
2. Spline means function **smooth globally**, the pieces are joined so that the function is globally,  $S_3 \in C^2(\bar{I})$
3.  $S_3(x_i) = f(x_i) \quad i = 0, \dots, n$

In matlab **d = spline(x,y,z)**, build and directly evaluate. For each interval we have a polynomial of degree 3:  $S_3$  (so  $a_i$  for  $i = 4$ ). We have 4 unknowns for each interval, let #intervals =  $n$ , so  $4n$  unknowns. The procedure:

- 1)  $S_3(x_i) = f(x_i) \quad i = 0, \dots, n$
- 2) We demand  $S_3$  continuous in the nodes,  $S_3 \in C^0([x_0, x_n])$

$$[S_3(x_i)]^- = [S_3(x_i)]^+ \quad i = 1, \dots, n-1$$

- 3) We demand  $S'_3$  continuous in the nodes,  $S'_3 \in C^0([x_0, x_n])$

$$[S'_3(x_i)]^- = [S'_3(x_i)]^+ \quad i = 1, \dots, n-1$$

- 4) We demand  $S''_3$  continuous in the nodes,  $S''_3 \in C^0([x_0, x_n])$

$$[S''_3(x_i)]^- = [S''_3(x_i)]^+ \quad i = 1, \dots, n-1$$

In total we have  $(n+1) + 3(n-1) = 4n-2$  conditions, but we need two more:

- $S''_3(x_0) = S''_3(x_n) = 0$ , natural cubic interpolating spline
- Not-a-knot-condition:  $S''_3$  continuous at  $x_1, x_{n-1}$

$$[S'''_3(x_1)]^- = [S'''_3(x_1)]^+ \quad [S'''_3(x_{n-1})]^- = [S'''_3(x_{n-1})]^+$$

The error:

$$\max_{x \in [x_0, x_n]} |f^{(r)}(x) - S_3^{(r)}(x)| \leq C_r \cdot H^{4-r} \cdot \max_{x \in [x_0, x_n]} |f^{(4)}(x)| \quad r = 0, 1, 2$$

## 5.2 Least Squared Approximation

Data  $\{(x_i, y_i)\}_i^n$ ,  $x_i$  distinct, we find a polynomial  $\tilde{f} \in \mathbb{P}_m$  of degree  $m \geq 1, m \ll n$  and  $\tilde{f}(x_i) \neq y_i$  such that:

$$\underbrace{\sum_{i=0}^m [\tilde{f}(x_i) - y_i]^2}_A \leq \underbrace{\sum_{i=0}^n [p_m(x_i) - y_i]^2}_B \quad \forall p_m \in \mathbb{P}_m$$

We want to minimize the right term.

### 5.2.1 Degree n

$m = n$  Lagrange interpolant,  $\tilde{f} = \Pi_n$

### 5.2.2 Degree 1

$m = 1$  regression line

$$p_1(x) = b_0 + b_1x \quad b_0, b_1 \in \mathbb{R}$$

$$\tilde{f}(x) = a_0 + a_1x \quad a_0, a_1 \in \mathbb{R}$$

We want to find the specific polynomial  $\tilde{f}$  (so the two coefficients). Consider the definition, the blocks:

$$\bullet A = \sum_{i=0}^n [a_0 + a_1x_i - y_i]^2 = \Phi(a_0, a_1)$$

$$\bullet B = \sum_{i=0}^n [b_0 + b_1x_i - y_i]^2 = \Phi(b_0, b_1)$$

So:

$$\Phi(a_0, a_1) \leq \Phi(b_0, b_1) \quad \forall b_0, b_1 \in \mathbb{R}$$

We compute the partial derivatives to find  $a_0$  and  $a_1$ :

$$\left. \frac{\partial \Phi}{\partial b_0} \right|_{(b_0, b_1) = (a_0, a_1)} = 0 \quad \left. \frac{\partial \Phi}{\partial b_1} \right|_{(b_0, b_1) = (a_0, a_1)} = 0$$

By developing  $\Phi(b_0, b_1)$

$$\Phi(b_0, b_1) = \sum_{i=0}^n [b_0^2 + b_1^2 x_i^2 + y_i^2 + 2b_0 b_1 x_i - 2b_0 y_i - 2b_1 x_i y_i]$$

And the partial derivatives:

$$\frac{\partial \Phi}{\partial b_0} = \sum_{i=0}^n [2b_0 + 2b_1 x_i - 2y_i] \quad \frac{\partial \Phi}{\partial b_1} = \sum_{i=0}^n [2b_1 x_i^2 + 2b_0 x_i - 2x_i y_i]$$

Putting  $a_0$  and  $a_1$

$$\begin{cases} \sum_{i=0}^n [2a_0 + 2a_1 x_i - 2y_i] = 0 \\ \sum_{i=0}^n [2a_1 x_i^2 + 2a_0 x_i - 2x_i y_i] = 0 \end{cases} \rightarrow B \vec{a} = \vec{f}$$

$$\begin{cases} \sum_{i=0}^n a_0 + \sum_{i=0}^n a_1 x_i = \sum_{i=0}^n y_i \\ \sum_{i=0}^n a_1 x_i^2 + \sum_{i=0}^n a_0 x_i = \sum_{i=0}^n x_i y_i \end{cases} \Rightarrow \begin{cases} a_0(n+1) + a_1 \sum_{i=0}^n x_i = \sum_{i=0}^n y_i \\ a_1 \sum_{i=0}^n x_i^2 + a_0 \sum_{i=0}^n x_i = \sum_{i=0}^n x_i y_i \end{cases}$$

$$B = \begin{bmatrix} (n+1) & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \quad \vec{f} = [\sum y_i, \sum x_i y_i]^T \quad \vec{a} = [a_0, a_1]^T$$

Where  $B$  is spd (we can use gradient method)

### 5.2.3 Degree m generic

$$\tilde{f}(x) = a_0 + a_1x + a_2x^2 + \dots + a_mx^m$$

$$p_m(x) = b_0 + b_1x + b_2x^2 + \dots + b_mx^m$$

$$\left. \frac{\partial \Phi}{\partial b_i} \right|_{(b_0, b_1, \dots, b_m) = (a_0, a_1, \dots, a_m)} = 0 \quad i = 0, \dots, m$$

Just like before, we want to find  $a_0, \dots, a_m$ , after the calculations:

$$B = \begin{bmatrix} (n+1) & \sum x_i & \sum x_i^2 & \dots & \sum x_i^m \\ \sum x_i & \sum x_i^2 & \sum x_i^3 & \dots & \sum x_i^{m+1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum x_i^m & \sum x_i^{m+1} & \sum x_i^{m+2} & \dots & \sum x_i^{2m} \end{bmatrix}$$

$$\vec{a} = \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_m \end{bmatrix} \quad \vec{f} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \\ \vdots \\ \sum x_i^m y_i \end{bmatrix}$$



### 5.3 Numerical Integration

$$I(f) = \int_a^b f(x)dx \quad f \in C^0([a, b])$$

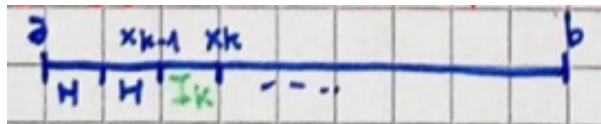
The quadrature rule: approximate a function:

$$\tilde{I}(f) \simeq I(f)$$

$$\tilde{I}(f) = \int_a^b \tilde{f}(x)dx \quad \tilde{f}(x) \simeq f(x)$$

The steps:

- Subdivide into  $M$  intervals, the points are known as quadrature nodes



$M + 1$  points, so  $M$  intervals,  $x_i$  are quadrature nodes. The amplitude:

$$H = \frac{b-a}{M}$$

Uniform partition of quadrature nodes:

$$x_{i+1} = x_i + H \quad i = 0, \dots, M-1$$

$$x_i = x_0 + (i)H \quad i = 1, \dots, M$$

- Expand the additivity and associativity of integral:

$$\int_a^b f(x)dx = \sum_{k=1}^M \int_{I_k} f(x)dx \cong \sum_{k=1}^M \int_{I_k} \tilde{f}(x)dx = \tilde{I}(f)$$

With  $\tilde{f}(x) \in \mathbb{P}_m$ , for different choices of  $m$  different quadrature rules

#### 5.3.1 Newton-Cotes

We have first to define:

- Order of accuracy**, associated only to the composite rule, "rate of convergence for the quadrature rule to zero- of the error" (it's the power of  $H$  in the composite error)
- Degree of exactness**, associated both to composite and simple, maximum degree of the polynomials which are exactly integrated by your quadrature rule. Suppose:

$$I(f) = \int_a^b f(x)dx \quad \tilde{I}(f)$$

We start from polynomial of degree 0  $p_0$ , but they are infinite: choose a representant

$$I(1) = ? = \tilde{I}(1)$$

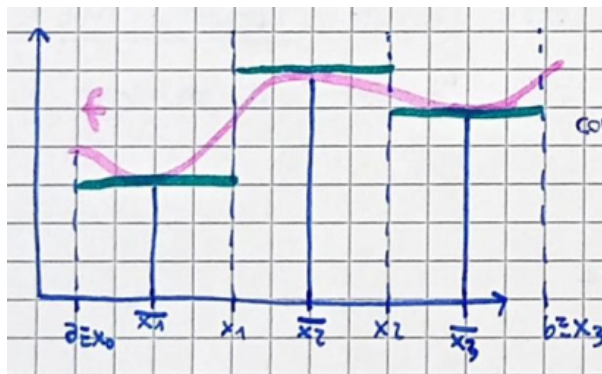
Then continue till this check does not hold. But in practice the following equality holds:

$$de = [\text{degree of derivative}] - 1$$

If we know the explicit expression of the error ( $E$ ) we try different degrees in order to make it zero

We have

- $m = 0$ , **midpoint quadrature rule**, which means for each subinterval  $\tilde{f}$  is a constant function



$$\begin{cases} \tilde{f} \in \mathbb{P}_0 \\ I_k = [x_{k-1}, x_k] \\ \bar{x}_k = \frac{x_{k-1} + x_k}{2} \end{cases}$$

The composite midpoint quadrature rule:

$$\tilde{I}_{MP}^c(f) = H \sum_{k=1}^M f(\bar{x}_k)$$

Where  $c$  stands for composite, the simple midpoint quadrature rule version:

$$\tilde{I}_{MP}(f) = (b-a)f\left(\frac{a+b}{2}\right)$$

**Errors:**

- Simple midpoint quadrature rule error

$$I(f) - \tilde{I}_{MP}(f) = \int_a^b [f(x) - f(\bar{x})]$$

Use the Taylor expansion centered at  $\bar{x}$ ,  $f \in C^2([a, b])$ , second order:

$$f(x) - f(\bar{x}) = f'(\bar{x})(x - \bar{x}) + \frac{f''(\alpha(x))}{2}(x - \bar{x})^2$$

Making the computations...

$$\tilde{E}_{MP} = I(f) - \tilde{I}_{MP}(f) = \frac{(b-a)^3}{24} f''(\beta)$$

$$f \in C^2([a, b]) \quad \beta \in [a, b]$$

$\beta$  (from min value theorem of integral) cannot be found, in practice upperbound for worst case

- Composite midpoint quadrature rule error

$$I(f) - \tilde{I}_{MP}^c(f) = \sum_{k=1}^M \left[ \int_{I_k} f(x) dx - \tilde{I}_{MP}(f|_{I_k}) \right] = \sum_{k=1}^M \frac{H^3}{24} f''(\beta_k)$$

Again from min value theorem of summation (dual of integral one):

$$= \frac{H^3}{24} f''(\gamma) \sum_{k=1}^M 1 = \frac{H^3}{24} f''(\gamma) M =$$

With  $H = \frac{b-a}{M} \rightarrow M = \frac{b-a}{H}$ , so:

$$\tilde{E}_{MP}^c = I(f) - \tilde{I}_{MP}^c(f) = \frac{(b-a)}{24} H^2 f''(\gamma)$$

$$f \in C^2([a, b]) \quad \gamma \in [a, b]$$

**Order of accuracy:**  $oa_{MP} = 2$

**Degree of exactness:**  $de_{MP} = 1$

–  $p_0$  degree 0

$$I(1) = ? = \tilde{I}(1)$$

$$(b-a) = ? = (b-a) f\left(\frac{a+b}{2}\right) = (b-a)$$

OK

–  $p_1$  degree 1

$$I(x) = ? = \tilde{I}(x)$$

$$\frac{x^2}{2} \Big|_a^b = ? = (b-a) \frac{a+b}{2}$$

OK, make the computations

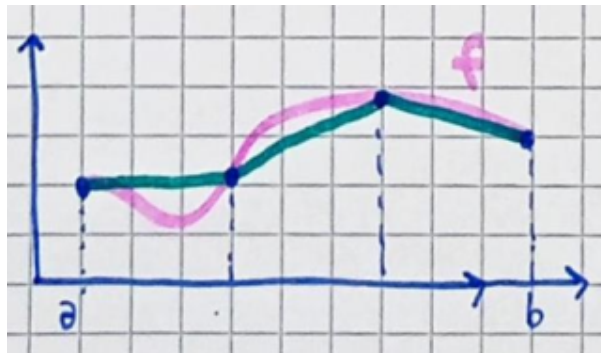
–  $p_1$  degree 2

$$I(x^2) = ? = \tilde{I}(x^2)$$

$$\frac{x^3}{3} \Big|_a^b = ? = (b-a) \left[ \frac{a+b^2}{2} \right]$$

KO, we stop here, so  $de_{MP} = 1$

- $m = 1$ , **trapezoidal quadrature rule**, for each subinterval a linear function



The composite trapezoidal quadrature rule (just basis times height):

$$\tilde{I}_T^c(f) = \frac{H}{2} \sum_{k=1}^M [f(x_{k-1}) + f(x_k)] = \frac{H}{2} [f(a) + f(b)] + H \sum_{k=1}^{M-1} f(x_k)$$

The first expression associated to intervals, the second to nodes. The simple trapezoidal quadrature rule version:

$$\tilde{I}_T(f) = \frac{(b-a)}{2} [f(a) + f(b)]$$

**Errors:**

- Simple trapezoidal quadrature rule

$$I(f) - \tilde{I}_T(f) = \int_a^b [d(x) - \Pi_1(f)(x)] dx$$

The term inside the integral is:

$$E_1 = \frac{f''(\eta(x))}{2} (x-a)(a-b)$$

From

$$E_n f(x) = \frac{f^{(n+1)}(\alpha(x))}{(n-1)!} \prod_{k=0}^n (x-x_k)$$

Making the computations...

$$\tilde{E}_T = I(f) - \tilde{I}_T(f) = -\frac{(b-a)^3}{12} f''(\delta)$$

$$f \in C^2([a, b]) \quad \delta \in [a, b]$$

- Composite trapezoidal quadrature rule

$$\tilde{E}_T = I(f) - \tilde{I}_T^c(f) = -\frac{(b-a)}{12} H^2 f''(\rho)$$

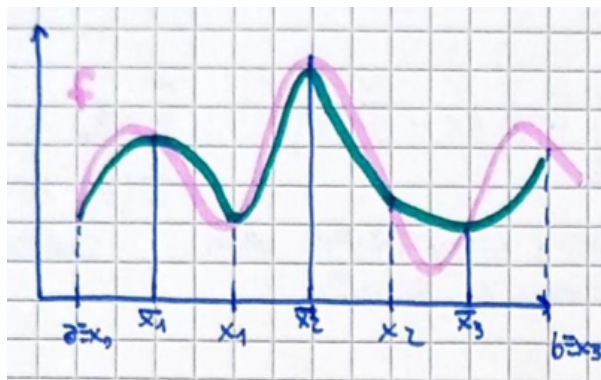
$$f \in C^2([a, b]) \quad \rho \in [a, b]$$

**Order of accuracy:**  $oa_T = 2$

**Degree of exactness:**  $de_T = 1$

Though  $oa$  and  $de$  same as midpoint, we see that the error is twice that of the midpoint: midpoint is easier and has lower error

- $m = 2$ , **Simpson quadrature rule**

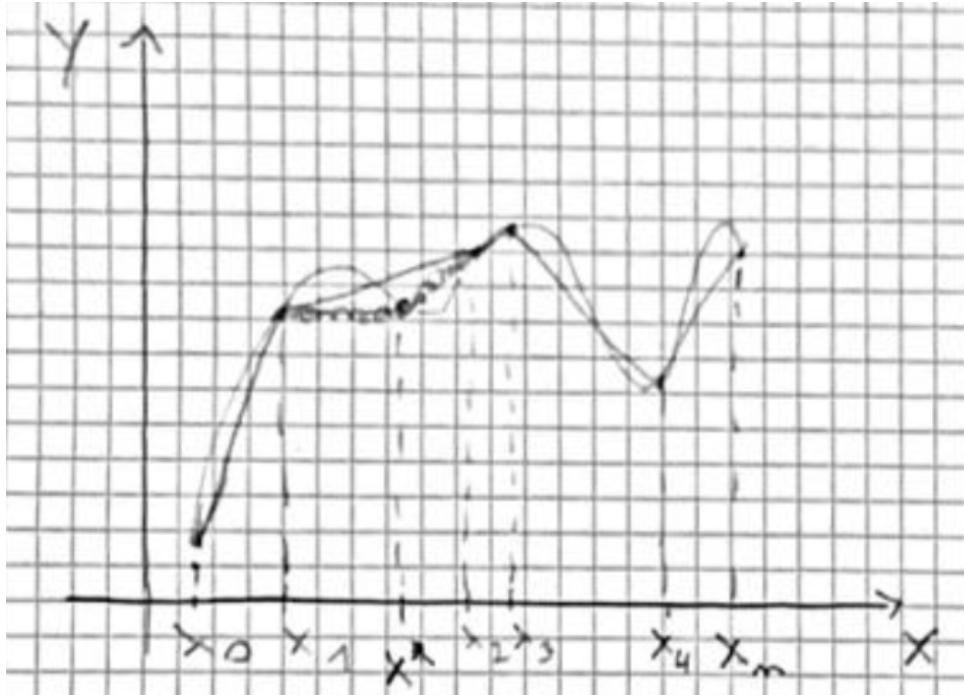


$$\tilde{I}_S^c(f) = \frac{H}{6} \sum_{k=1}^M [f(x_{k-1}) + 4f(\bar{x}_k) + f(x_k)] =$$

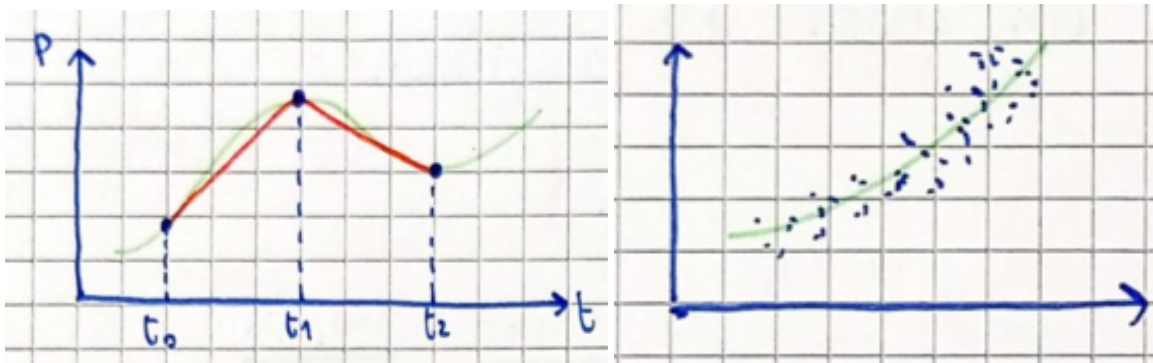
$$= \frac{H}{6} [f(a) + f(b)] + \frac{H}{3} \sum_{k=1}^{M-1} f(x_k) + \frac{2}{3} H \sum_{k=1}^M f(\bar{x}_k)$$

The simple Simpson quadrature rule version:

$$\tilde{I}_S(f) = \frac{(b-a)}{6} \left[ f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right]$$



**Figure 5:** Modified Newton graph



```

while (true)
     $x^{(k)} = (a^{(k)} + b^{(k)})/2$ 
    if ( $f(x^{(k-1)}) = 0$ ) break;
    if ( $f(a^{(k-1)})f(x^{(k-1)}) < 0$ )
         $a^{(k)} = a^{(k-1)}, b^{(k)} = x^{(k-1)}$ ;
    else
         $a^{(k)} = x^{(k-1)}, b^{(k)} = b^{(k-1)}$ ;
    end

```

