

# Numerical Analysis ?

Subject of writing numerical approx to problem  
 Idea: Numerical is the heart of studying math which is close to reality but not the same as the desired quantity  
 (approximation)

- Algorithms
- Analysis (errors)
- Implementation (python)

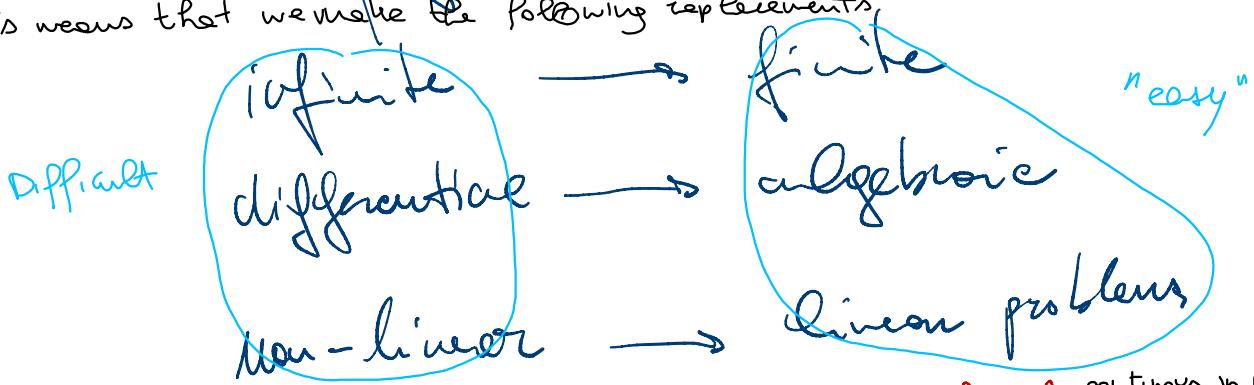
And we want to make sure this approx works within certain errors  
 Idea: get as close as possible to the original problem without a good approx

We can only approximate **Well posed problems**. i.e.

- ||| 1) there exist a unique solution  
 2) it depends continuously on the data |||

idea behind well posed is that we simplify in a controlled manner

This means that we make the following replacements:



Example:

Given a function "f", a point  $x_0$ ,  
 can get "f'(x\_0)" (derivative of f in  $x_0$ )

Is it well posed? Yes, because exist a UNIQUE JT ← derivative of f at  $x_0$   
 If f is not derivable  $\Rightarrow$  problem is not well posed  $\Rightarrow$  WE MUST CHECK DOMAIN!  
 → use finite difference

I'm not claiming to know anything about f except the evaluation of its body

$$f'(x_0)$$

$$\approx$$

$$\frac{f(x_0 + h) - f(x_0)}{h}$$

anything you choose here will be affected by what you can do and what you can't

with small h  
 $\uparrow$

(s.t.  $x_0 + h \in (a, b)$ )

Error: Taylor expansion around  $x_0$ :

Typical ways to estimate error is

$$f(x_0 + h) = f(x_0) + h \cdot f'(x_0) + \frac{h^2}{2} f''(x_0) + \dots$$

When you compute powers of  $2$  ( $2^{**i}, i \dots$ ),  
The accuracy of computer is perfect upto  
 $\sim 10^{-4}$ , and in reality is not perfect upto  
 $10^{14} - 10^{-15}$

In theory, as  $h \rightarrow 0$ , the finite approx should  
get closer to the exact result (exact derivative)

In practice much else happens

$$f(x_0 + h) \approx \sum_{i=0}^n \frac{f^{(i)}(x_0) h^i}{i!}$$

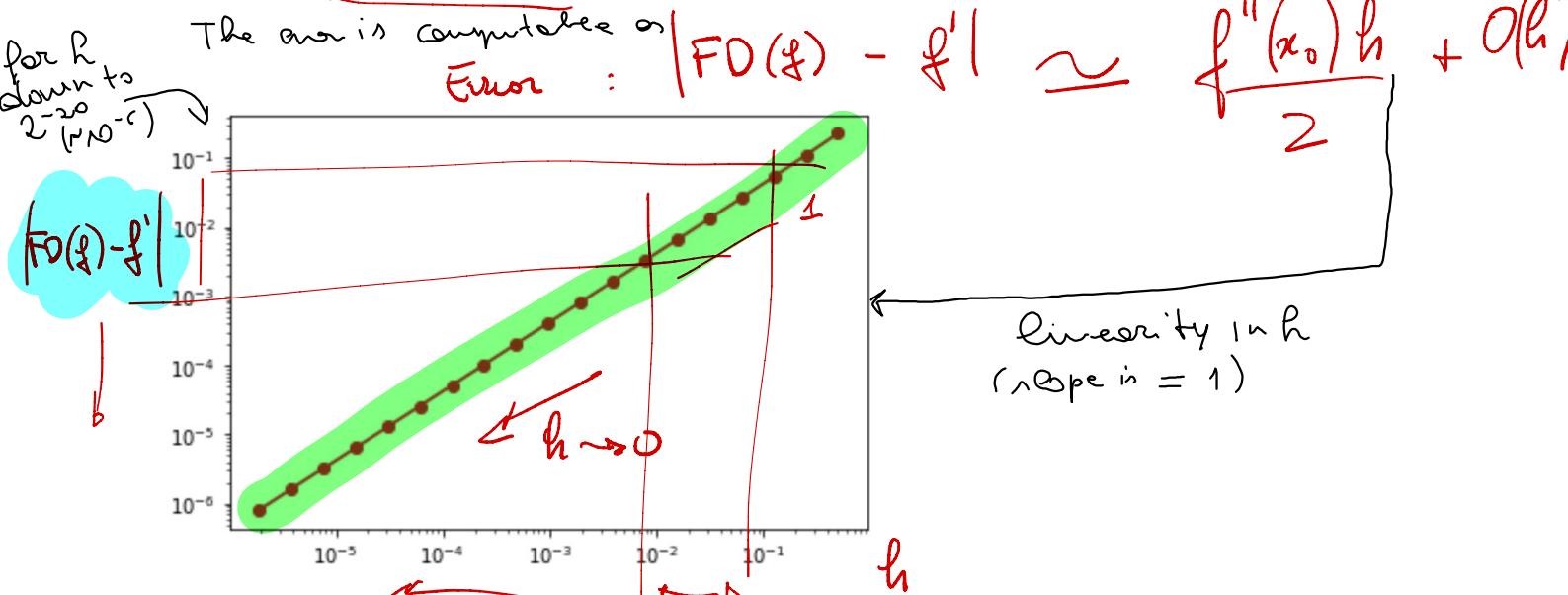
$$FD := \frac{f(x_0 + h) - f(x_0)}{h}$$

for  $n \rightarrow \infty$ ,  $\approx \equiv$   
 provided that function  
 is nicely behaved, so that  
 I can do it for  $n$  functions  
 at each step

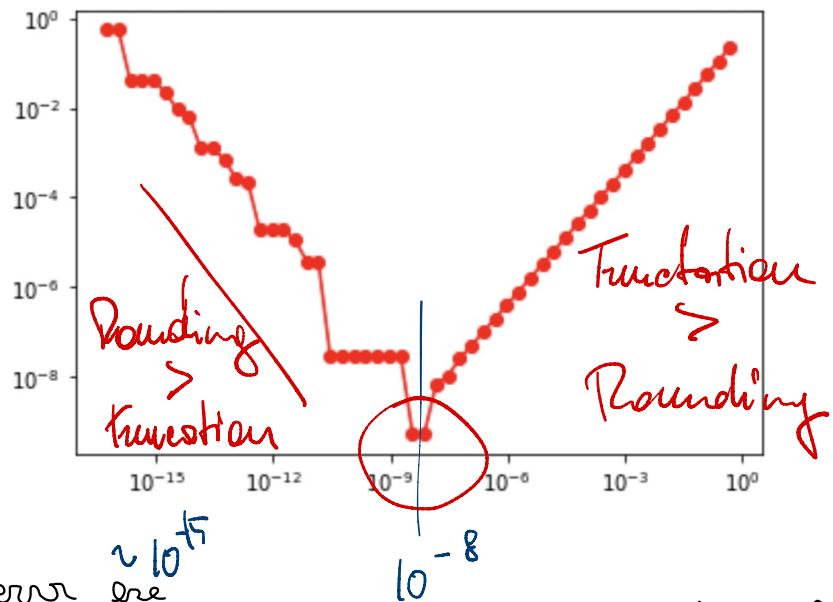
$$f(x_0) + \frac{1}{2} f''(x_0) h + O(h^2)$$

exact derivative  $=$

+ much smaller than  $O(h^2)$  means that  
 when  $h \rightarrow 0$  it's even smaller compared  
 to  $f''(x_0) h$



interesting things happens when  $h$  keeps getting smaller, but the algorithm above doesn't tell us that: it tells us that the error, being  $\sim h$ , will go to zero in a linear fashion as  $h$  goes down



$\sqrt{\epsilon} \sim \text{optimal } h$

$\epsilon$ : precision of machine

(Rounding precision)

Types of error are

$\sim 10^{-15}$

Data error  $\rightarrow$  always present

(Data)

If you have error in measurements of data i.e. the  $x_0$  above, and depends on used algorithms

Computational Error  $\rightarrow$  (Truncation error)

(Algorithm)

Truncation error

(Algorithm error in Exact Arithmetic)

Rounding Error  $\rightarrow$  due to algorithms of computers and numbers we deal with every day

(Finite Arithmetic)

the one of PC

$\sim$  error you make on the algorithm level by using approximating exact of the code

If doesn't take account of the fact that machines operate in a different arithmetic

arithmetic in pc is of power of 2

diff between what you expect from algorithm by putting exact st, and what you get when you run your algorithm (the approximation)

### FINE DIFF:

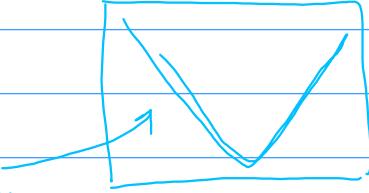
Infinite series - Truncated series

linear  
the graph above is  
**TRUNCATION**  
**Error ONLY**

because trunc. error dominates over rounding error

but, as h drops, rounding error starts becoming more relevant

here it dominates over the trunc error



Pb: Recompute value of a function "F"

$$F : \mathbb{R} \longrightarrow \mathbb{R}$$

at a argument  $x$

So we have

•  $x$  : true value of input

what to compute

•  $F(x)$  : desired value of output

affect by rounding,  
data error, ...

Error on  
input data

•  $\hat{x}$  : approximate input (may be "off" for rounding)

Error on  
outcome  
due to  
approximation

•  $\hat{F}(\hat{x})$  : computed output (may be "off" for both  
rounding and truncation)

$$\underbrace{F(\hat{x}) - F(x)}_{\text{effect of rounding on input data}} =$$

ROUNDING OR COMPUTATIONAL  
ERROR DUE TO NUMBER OF DIGITS  
IN INPUT (NUMBER OF DIGITS  
IN APPROXIMATED PROBLEM)

Truncation error is  
exist on floating

①

②

or

sensitivity of algorithm with  
respect to the original problem

sensitivity of original problem  
with respect to perturbations in  
the data

$$\underbrace{\hat{F}(\hat{x}) - F(\hat{x})}_{\substack{\text{You are } \uparrow \\ \text{controlling the } \downarrow \\ \text{algorithm, and } \\ \text{only know the } \uparrow \\ \text{condition is true, but } \\ \text{you can't remove it} \\ \text{You don't do slope}}}$$

$$+ \underbrace{F(\hat{x}) - F(x)}_{\substack{\text{sensitivity of } \uparrow \\ \text{original problem} \\ \text{with respect to } \uparrow \\ \text{perturbations in } \uparrow \\ \text{the data}}}$$

③

COMPUTATIONAL ERROR  
(diff between expected output  
and exact output)

the longest number S.f.  
in exact arithmetic

④ in exact arithmetic

③ in rounding arithmetic

Rounding precision :  $\epsilon$   
(machine precision)

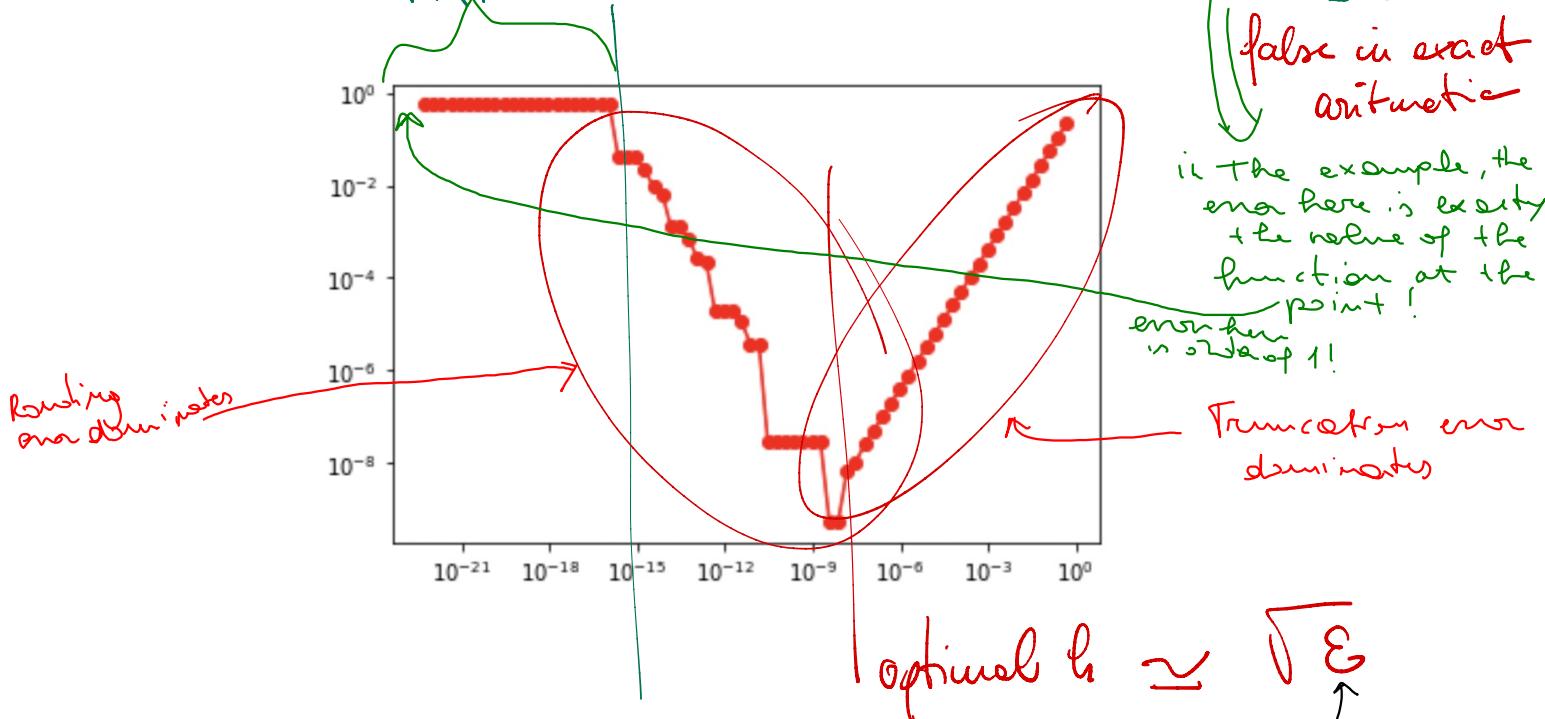
$$fl(1 + \epsilon) = fl(1)$$

gives optimal  
precision

$\mu < \epsilon$  including  
error dominates

floating point representation of a number

$$f(f(x+h)) = f(f(x)) \xrightarrow{\text{if } f'(x) = 0} f(x+h) - f(x) = 0$$



$$\sin(1) = 0.8414709848078965$$

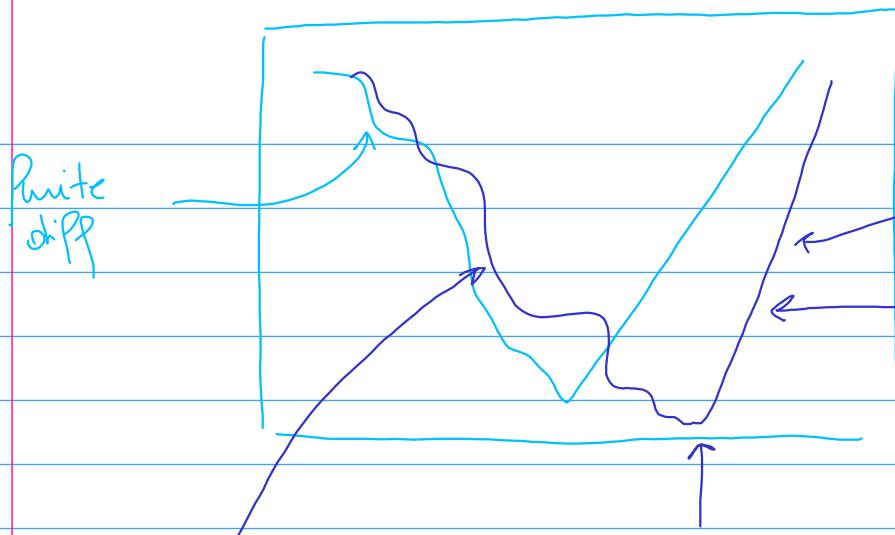
for G4 S.2  
doubles  
(python)  
 $\epsilon \sim 10^{-15}$

$$\frac{f(x+h) - f(x)}{h}$$

This  $h$  must be chosen correctly to prevent overflow or underflow, otherwise you'll get good results for some  $f$  and bad for others.

All of this also depends on where we evaluate  $f$  ( $x_n$ ) and function eval at the same time

! Don't use abs error as a measure of error, but use the RELATIVE Error



Not  
available,  
this gets  
bad  
further  
 $\sim \varepsilon^{1/3}$

optimal  
 $h$  in row  
 $\sim \varepsilon^{1/3}$

constant finite  
difference.  
BETTER, because  
here we have  
another power  
of  $n$ , which  
is power of 2  
(slope  $\ell n n$  is 2<sup>-1</sup>)

# Numerical Analysis - sec. 2

problem can be written  
in a functional approach

8.10.2020

$$y = F(x)$$

where

$$F: \mathcal{X} \longrightarrow \mathcal{Y}$$

$$\text{data} \rightarrow x \longrightarrow y = F(x) \quad \begin{matrix} \text{F is a given function} \\ \mathcal{X} \text{ output/result} \end{matrix}$$

## Functional Evaluation

EXAMPLES

$$1) \mathcal{X} : \mathbb{R}, \mathcal{Y} : \mathbb{R}$$

F: continuous func.

$$y = F(x)$$

### 1) Sum of two numbers:

$$\mathcal{X} : \mathbb{R}^2, \mathcal{Y} : \mathbb{R} \quad : F: \text{sum of numbers.}$$

$\downarrow$  pair of numbers

$$x = (a, b)$$

$$F(x) = F((a, b)) = a+b$$

### 3) Computation of derivative

$$\mathcal{X} : C^1([a, b]) \quad \mathcal{Y} : \mathbb{R}$$

F: evaluation of first derivative in  $x_0 \in [a, b]$

$$y = F(x) := \underbrace{x'}_{\text{function}}(x_0)$$

Only well posed problems are characterized by 2 properties

$$1) \forall x \in \mathcal{X}, \exists! y \in \mathcal{Y} \text{ s.t. } F(x) = y$$

$$2) \exists k \text{ (condition number)} \text{ s.t. } \forall x, \hat{x} \in \mathcal{X}$$

$$\|F(x) - F(\hat{x})\| \leq k \|x - \hat{x}\| \quad \begin{matrix} \text{basically} \\ F \text{ is} \\ \text{continuous} \\ (\text{bounded} \\ \text{for every}) \end{matrix}$$

$\overbrace{\text{each of these}}^{\text{are unique}}$

$\mathcal{Y}$

$\mathcal{X}$

We assume that

## Real Vector Spaces (Normed)

Def (real vector spaces)

RVS: collection of objects for which it makes sense  
"V" to "sum" and "scale"

$$+ : V \times V \longrightarrow V$$

$$(u, v) \longrightarrow w = u + v$$

$$\cdot : V \times \underline{\mathbb{R}} \longrightarrow V$$

such that

$$\forall u, v \in V, \forall \alpha, \beta \in \mathbb{R}$$

$$\alpha u + \beta v = w \in V$$

$\mathbb{R}^2$

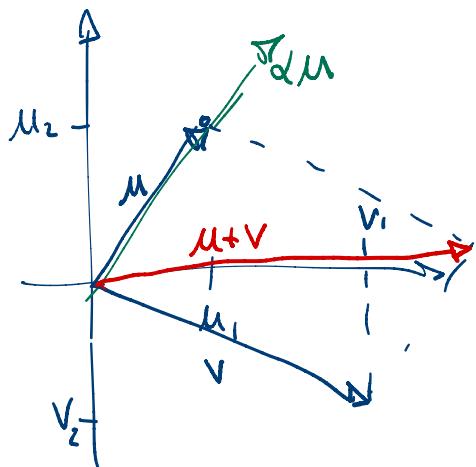
$C^0([a, b])$

$$u = (u_1, u_2)$$

$$v = (v_1, v_2)$$

$$u+v = (u_1+v_1, u_2+v_2)$$

$$\alpha u = (\alpha u_1, \alpha u_2)$$



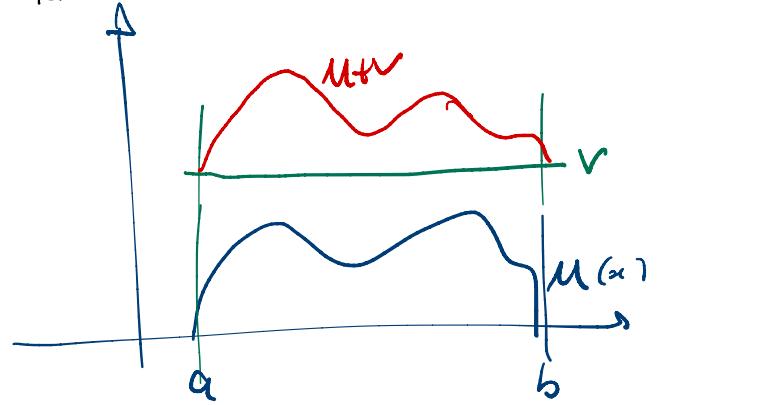
$$u : [a, b] \longrightarrow \mathbb{R}$$

$$v : [a, b] \longrightarrow \mathbb{R}$$

$$w = u+v : [a, b] \longrightarrow \mathbb{R}$$

because  $u, v$   
are defined  
in  $[a, b]$

$$\alpha \longrightarrow u(\alpha) + v(\alpha)$$



Way to measure things

Introduce  $\downarrow$  NORMS

# Norms on Real Vector Spaces

a function  $\|\cdot\| : V \rightarrow \mathbb{R}_0^+$  positive or zero

$$1) \|u\| > 0 \quad \forall u \in V$$

$$2) \|u+v\| \leq \|u\| + \|v\| \quad \text{Triangle Inequality}$$

$$3) \|\alpha u\| = |\alpha| \|u\|$$

$$\text{optional } 4) \|u\| = 0 \iff u = 0$$

without ④ it is a semi-norm

We will use

$L_p$  norms on  $\mathbb{R}^n$

$L_p$  norms on  $S \subset \mathbb{R}^d$

$\|\cdot\|_*$  operational norms induced by  $\|\cdot\|_V$

Def.  $\mathbb{R}^n$ :  $L_p$  norm  $\|u\|_p := \left( \sum_{i=1}^n |u_i|^p \right)^{\frac{1}{p}}$   $\downarrow p \rightarrow \infty$

$\|u\|_\infty := \max_i |u_i|$   $\ell_\infty$  norm

$S \subset \mathbb{R}^d$   $L_p$  norms  $\|u\|_p := \left( \int_S |u|^p \right)^{\frac{1}{p}}$   $\downarrow p \rightarrow \infty$  L<sub>p</sub>-norm

$\|u\|_\infty := (\text{ess}) \sup_{x \in S} |u(x)|$

$\|\cdot\|_*$  induced by vector space norm  $\|\cdot\|_V$  on  $V$  and  $\|\cdot\|_W$  on  $W$

$\uparrow$   
used to measure  
how functional  
functions work on  
what is their measure  
in a space

$A: V \rightarrow W$

$\|A\|_* := \sup_{0 \neq x \in V} \frac{\|A(x)\|_W}{\|x\|_V}$

You def the norms of domain and codomain  
of  $A$ , and then you can def the operat norm

$\ell_p$  norm of Matrices in  $\mathbb{R}^{n \times m}$

which is an obj of this kind

$A^{n \times m}$

$x \in \mathbb{R}^m$

$\cdot x$

$A \cdot x$

$\underline{y}$

$\in \mathbb{R}^n$

$$A : \mathbb{R}^m \rightarrow \mathbb{R}^n$$

$$A \cdot x = \underline{y} \in \mathbb{R}^n$$

$\ell_p$  norm  
of the matrix  
functional  
norm because linear  
operator norm of a function

$$A \in \mathbb{R}^{h \times m}$$

$$\|A\|_p := \sup_{0 \neq x \in \mathbb{R}^n} \frac{\|Ax\|_p}{\|x\|_p} = \|A\| \text{ induced by } \|.\|_p \text{ on } \mathbb{R}^n$$

$$A(x) = y$$

$$\sum_{j=1}^m A_{ij} x_j = y_i$$

$$A \in L(\mathbb{R}^m, \mathbb{R}^n)$$

$$\| \cdot \|_* \text{ norm on } L(\mathbb{R}^m, \mathbb{R}^n) \equiv \mathbb{R}^{n \times m}$$

Why we need them? Every quantity related to a well posed problem can be put in a way between two spaces and their norms.

## Well posed problems

$$1) \quad \forall x \in \mathcal{X}, \exists! y \in \mathcal{Y} \text{ s.t. } F(x) = y$$

Normed Vector Space

Normed Vector Space

$$2) \quad \exists k \text{ (condition number)} \text{ s.t. } \forall x, \hat{x} \in \mathcal{X}$$

$$\|F(x) - F(\hat{x})\|_y \leq (k_{\text{Abs}}) \|x - \hat{x}\|_x$$

Here we have norms

$$3) \quad \text{Relative cond. number (only if } x \neq 0, F(x) \neq 0)$$

$$\text{Krel st. } \forall x, \hat{x}$$

INDEPENDENT ON  $x$

$$\frac{\|F(x) - F(\tilde{x})\|_Y}{\|F(x)\|_Y} \leq \kappa_{rel} \frac{\|x - \tilde{x}\|_X}{\|x\|_X}$$

Prop

A problem is well posed if  $\kappa_{rel/abs}$  is

"small"  $\leftrightarrow$  ~~Krel/abs must be controllable~~

~~with respect to what  $\rightarrow$  relative depends on problem handled~~

Example : sum of two numbers: ( $L_1$  norm)

$$X: \mathbb{R}^2 \quad Y: \mathbb{R}$$

$$x \in \mathbb{R}^2 \rightarrow \|x\|_1 := |x_1| + |x_2| \quad \|y\|_1 = |y|$$

We have that

$$\frac{\|F(x) - F(\tilde{x})\|_1}{\|F(x)\|_1} = \frac{|x_1 - \hat{x}_1 + x_2 - \hat{x}_2|}{|x_1 + x_2|} \leftarrow \| \cdot \|_{Y=1}$$

$$\frac{\|x - \tilde{x}\|_1}{\|x\|_1} = \frac{|x_1 - \hat{x}_1| + |x_2 - \hat{x}_2|}{|x_1| + |x_2|}$$

$$\Delta y, \quad \Delta x \quad \Delta x = x - \hat{x} \quad \Delta y = F(x) - F(\tilde{x})$$

We now find  $\kappa_{abs}$ :  $\exists \kappa_{abs} \mid \forall x_1, x_2, |\Delta y| \leq \kappa \|\Delta x\| ?$

$$|\Delta y| \leq \kappa \|\Delta x\| \leq \kappa(|\Delta x_1| + |\Delta x_2|)$$

Assume that  $\Delta x \neq 0$

$$\frac{|\Delta g|}{|\Delta x|} \leq \kappa$$

$$\frac{|\Delta x_1 + \Delta x_2|}{|\Delta x_1| + |\Delta x_2|} \leq \frac{|\Delta x_1| + |\Delta x_2|}{|\Delta x_1| + |\Delta x_2|} \leq 1$$

From pt 5 of absolute  $\kappa$ , this is perfect: you perturb the data, and the perturbation of the result follows perfectly the pt b of the inputs ( $\kappa_{abs} = 1$ )

$$\kappa_{abs} = 1$$

but ~~not~~ local st. Indeed:

$$\frac{|\Delta y|}{|\Delta x|}, \frac{|x|}{|y|} \leq \frac{|x|}{|y|} \leq \kappa_{rel}$$

$\kappa_{abs} \leq 1$

we cannot control this number!

because

$\Rightarrow$  if you try to sum 2 numbers with closing values and opposite signs, you're dealing with an ill-posed problem, but the sum of 2 numbers is not guaranteed to be well-conditioned.

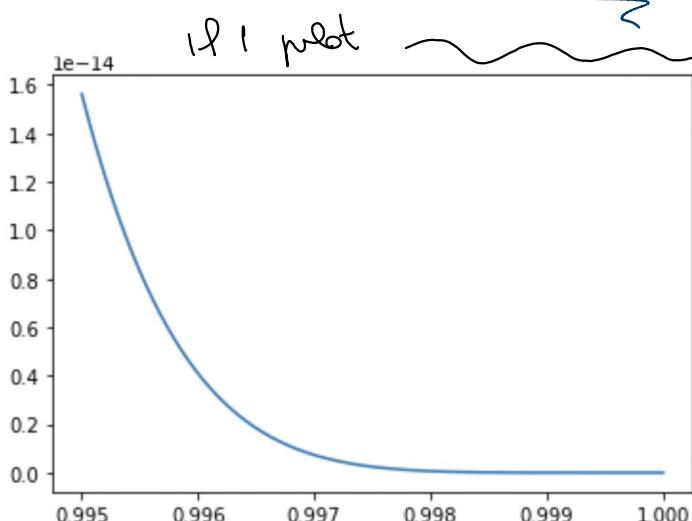
Example:

$$\frac{|x|}{|y|} = \frac{|x_1| + |x_2|}{|x_1 + x_2|}$$

can be as large as we want

$$y = (1-x)^6$$

$x$  is small

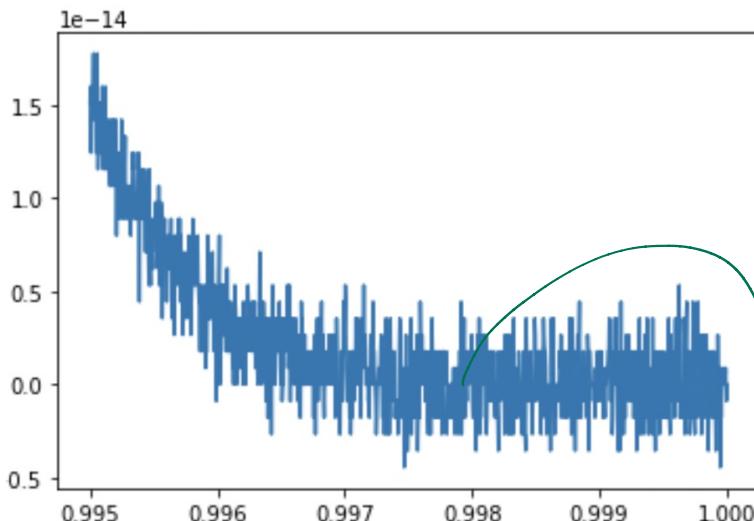


```
x = linspace(.995, 1, 1029)
```

```
y = (x-1)^**6
```

but if I round the power I get a bad result

$$y = x^{**6} - 6*x^{**5} + 15*x^{**4} - 20*x^{**3} + 15*x^{**2} - 6*x + 1$$



due to addition of very close numbers (in value) but with  $\neq$  signs  $\rightarrow$  we say that this is an ill posed problem  
 $x - x = \text{floating point errors}$   
(rounding errors)

Consider:

If  $x_1 = -(x_2 + \varepsilon)$

so what we obtain from machine epsilon

$$\frac{|x_1|}{|y|} = \frac{|x_2 + \varepsilon| + |x_2|}{|-x_2 - \varepsilon| + |x_2|} \leq \frac{2|x_2| + |\varepsilon|}{|\varepsilon|}$$

if  $x_2$  not close to zero then

$$\frac{|x_1|}{|y|} \leq \frac{2|x_2|}{|\varepsilon|} + 1$$

so, relating this to the Krel quantity  $\frac{|\Delta y|}{|y|}$  if  $|\varepsilon|$  is too small we get a problem --

$$\frac{|\Delta y|}{|y|} \cdot \frac{|x|}{|\Delta x|} \leq K_{\text{rel}} \Rightarrow$$

which means  $K_{\text{rel}} = \infty$  to agree with its def

$$\frac{|\Delta y|}{|\Delta x|} \leq K_{\text{rel}}$$

It has to be greater for every chosen  $\varepsilon$ )

In this one, it may also depend on  $x$

Sum of 2 numbers is not a "RELATIVELY"  
well posed problem ( $\leftrightarrow$  in relative terms  
 $\leftrightarrow$  k rel) relative change in result is caused by  
relative change in input, what you get is  
not controllable

But it is an "ABSOLUTELY" well posed problem

THIS DEPENDS ON  
THE EVALUATION POINT

$\uparrow$  difference!  
 $\downarrow$   
if my pt will be 0  
for non zero  $x \leftarrow$  difference!  
 $\Rightarrow$  k rel  $\rightarrow \infty$   
 $\Rightarrow$  ill posed problem

! Rounding errors are  
MAGNIFIED by  
big old train  
numbers

1) 2! solution

2) pb. is stable  $\Leftrightarrow$  well posed  $\Leftrightarrow$  finite cond. Num.

## $X, Y$ Normed Vector Spaces

$$F: X \rightarrow Y$$

]! solution  $\Rightarrow F$  is an INVERTIBLE FUNCTION  
positive and possibly infinite  
used whenever is required to study stability of a problem and implies that there is no more redundancy in data

We have

1)  $k_{abs} \in \mathbb{R}^+ \cup \{+\infty\}$

s.t.  $x, \hat{x} \in X$  LIPSCHITZ PROPERTY,  $\|F(x) - F(\hat{x})\|_Y \leq k_{abs} \|x - \hat{x}\|_X$

which is  $\Leftrightarrow$  to  $k_{abs} \geq \frac{\|F(x) - F(\hat{x})\|_Y}{\|x - \hat{x}\|_X}$  choose the smallest number greater than or equal to  $k_{abs}$

If  $k_{abs} \neq \infty \Rightarrow$  problem is stable (well conditioned)

1)  $k_{rel} \in \mathbb{R}^+ \cup \{+\infty\}$

s.t.  $x, \hat{x} \in X$

Similarly as before, the smallest number greater than or equal to

$k_{rel} : \sup_{x, \hat{x} \in X}$

$$\frac{\|F(x) - F(\hat{x})\|_Y}{\|F(x)\|_Y} \leq k_{rel} \frac{\|x - \hat{x}\|_X}{\|x\|_X}$$

denominates as the relativity

I split the norm because I have inequalities but then I can obviously the definition and put  $k_{rel}$  even to this

$$\frac{\|F(x) - F(\hat{x})\|_Y}{\|x - \hat{x}\|_X} \frac{\|x\|_X}{\|F(x)\|_Y}$$

If  $F$  is not linear  $\Rightarrow$  has in necessarily the Lipschitz constant of  $F$

$$k_{abs} \sup_{x, \hat{x} \in X} \|F'(F(x))\|_X$$

use def. of operator norm applied to  $F^{-1}$ : can do this because  $F$  is invertible

$$\frac{\|x - \hat{x}\|_X}{\|F(x)\|_Y}$$

$F^{-1}$  is

If operator is linear

Example:

$$Ax = y \Rightarrow y = F(x) \quad y = F(x)$$

$$\delta x = x - \hat{x} \quad \delta y = y - \hat{y}$$

Absolute stability:

$$A\hat{x} = \hat{y} \Rightarrow \hat{y} = F(\hat{x})$$

$$A\delta x = \delta y$$

$$\|Ax - A\hat{x}\| = \|A(x - \hat{x})\| = \|\delta y\|$$

I can collect the  $A$  thanks to the linearity of the problem

Exhibit property of operator norm of a matrix (follows from its definition)

$$\|Sg\| \leq \|A\| \|Sx\|$$

end we get  $K_{abs}$

Finding  $K_{abs}$  gives (by applying the def of  $K_{abs}$  above)

$$\frac{\|Sg\|}{\|y\|} \leq \|A\| \|A^{-1}\| \frac{\|Sx\|}{\|x\|}$$

Condition number of A

(this in def: condition number of a matrix is  $\|matrix\| \|matrix^{-1}\|$ )

it will depend on the norm you use

Prop

A problem is (absolutely) relatively well posed if  
( $K_{abs}$ )  $K_{abs} < +\infty$ .

What is an approximation?

Sequence of problems  $F_n$  ( $n \in \mathbb{N}$ )

applied to a sequence of "approximated data"  $x_n$   
(implied approximated sequence of domains  $X_n$ )

s.t.

$$\lim_{n \rightarrow \infty} |F_n(x_n) - F(x)| = 0 \quad \leftarrow$$

CONVERGENCE OF THE  
PROBLEM: this is really  
what we want to check

Note 1:

$$X_n \xrightarrow{\quad} X$$

Simple case:  $X_n \subset X$

2) Note 2: " $F_n \xrightarrow{\quad} F$ "

thin in the algorithm

consistency

Approximated  
method converges to  
the exact  
problem  
that we have

Assume that  $X_n = X$

$F_n \rightsquigarrow$  I can apply the  
algorithm to all pt  
of the "continuous  
space"

2)  $\xrightarrow{\text{can be rewritten}} \lim_{n \rightarrow \infty} |F_n(x) - F(x)| \rightarrow 0$

I'm asking that the algo  
applied to exact input  
minus the exact problem  
(in modulus) goes to zero  
for increasing  $n$

LAX-RICHTMEYER THEOREM:

For consistent pb ( $\lim_{n \rightarrow \infty} |F_n(x) - F(x)| = 0$ )

# CONVERGENCE

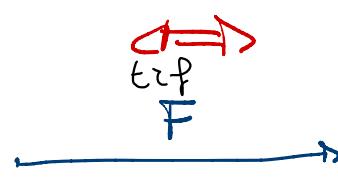
$\Leftrightarrow$

# STABILITY

Proof  $x^*$   $\xrightarrow{\text{input}}$

If  $F$  is well posed  $\Rightarrow F'(y_0) = \hat{x}$

$$F'(y_0) = \hat{x}$$



# STABILITY

output  $y$

$y$

$y_n$

can be  
expressed of  
 $x^*$  (e.g.  $x^* + \epsilon$ )  
or  $x^* + \delta x$

if you want to prove convergence you must show that

$$\underset{n \rightarrow \infty}{\lim} |F_n(x_n) - F(x)| \rightarrow 0$$

We have:

$$\left| F_n(x_n) - F(x_n) + F(x_n) - F(\hat{x}) + F(\hat{x}) - F(x) \right| = \dots \rightarrow 0$$

consistency  $\rightarrow 0$

stability  $\leq K|x_n - \hat{x}_n|$

stability  $\leq K|\hat{x} - x|$

All the errors falls into 3 categories

Stability errors

(of continuous pb)

(of discrete pb)

consistency error

if we control this, consistency is implied by stability of the problem.

If continuous problem is stable, and the discrete problem is stable and consistent  $\Rightarrow$  we have convergence

Example:

Interpolation

(Polynomial interpolation)

We consider the use of

One dimension.

and we study the following problem

Pb: Find a Polynomial of order  $n$  that coincides with the input function at (given) points

(Polynomial interpolation)

for simplicity but it is a general topic

$\mathcal{X} := C^0([a, b])$

continuous functions

They are not inputs!

They are part of the problem itself

CONSIDERATIONS

1. in  $\mathcal{X}$  we use the  $L^\infty$  norm:

$$\|u\|_\infty := \max_{x \in [a, b]} |u(x)|$$

We want to understand what  $F$  is and what's the condition number of the problem

$\mathbb{P}^n([a,b])$  is finite dimensional, while  $\mathcal{Y}$  is not because the outcome is a polynomial space of dimension  $n+1$ .

2. The space  $\mathcal{Y}$  is  $\mathbb{P}^n([a,b])$ : we use again  $\|\cdot\|_\infty$  as a norm.

Note 1:  $\mathcal{Y}$  is finite dimensional, and we can write linear combination

$$\text{it as } \mathcal{Y} = \mathbb{P}^n([a,b]) = \text{Span}\{v_i\}_{i=0}^n$$

$$\forall p \in \mathcal{Y} = \mathbb{P}^n([a,b]) \quad \exists! \{p^i\}_{i=0}^n \in \mathbb{R}^{n+1} \text{ s.t.}$$

$$p(x) = \sum_{i=0}^n p^i v_i(x)$$

↑ coefficients

linear independence of  $v_i, v_j \Leftrightarrow p^i = 0 \quad i=0, \dots, n \Leftrightarrow p(x) = 0 \quad \forall x \in [a,b]$

Polynomial interpolation: Given  $n+1$  points  $\{x_i\}_{i=0}^n$  (The algorithm takes a function (continuous) and returns a polynomial)

$$F: \mathcal{X} \equiv C^0([a,b]) \longrightarrow \mathcal{Y} \equiv \mathbb{P}^n([a,b])$$

Where  $p(x_i) = u(x_i) \quad i=0, \dots, n$

function  $\rightarrow$   $p$  polynomial

$\underbrace{\sum_{j=0}^n v_j p^j}_{\text{contravariant coeffs}} \Leftrightarrow \underbrace{p(x_i) = u(x_i)}_{n+1 \text{ conditions}}$  the given points

Can do this because I know that the space  $\mathbb{P}^n$  is finite dimensional, meaning that I can uniquely identify every element of that space with a vector of itself with respect to a given basis.

and invert order reading to Einstein notation  $\sum_{j=0}^n v_{ij} p^j = u_i \Leftrightarrow \sum v_{ij} = u_i$

↑ covariant coeffs

$$V_{ij} = v_j(x_i)$$

Vandermonde Matrix

what we doing: going from a vector of values of a pt in  $C^0([a,b])$  to a vector of coefficient of a basis in  $\mathbb{P}^n([a,b])$  (where  $v_i(x) = \text{pow}(x, i)$ )

Einstein notation:

$$V_{ij} p^j \equiv \sum_{j=0}^n V_{ij} p^j$$

we set  $(V_{ij})^{-1} = V^{ij}$  so that

$$V^{ij} V_{jk} = \delta_k^i = \begin{cases} 1 & \text{if } i=k \\ 0 & \text{if } i \neq k \end{cases}$$

so if needed

$$V_{ij} p^j = u_i$$

$$p^j = V^{ji} u_i \quad \Leftrightarrow \quad p = V^{-1} u$$

We define the "Reduced problem":  $X = \mathbb{R}^n$

where

$X$ : set of values of  $u$  in  $x_i$

with norm  $\|u\|_2 := \left( \sum_{i=0}^n u_i^2 \right)^{\frac{1}{2}}$  and same for  $y$  polynomial

The original problem can be cast in this problem with 2 additional functions

We want to find Condition number of  $F$ :  $u \rightarrow V u = p$

that is input  $\{u_i\}_{i=0}^n \rightarrow$  output  $\{p^j\}_{j=0}^n$

If I was suppose to perturb the input, starting from

$$\|p - \hat{p}\|_2 \leq K_{abs} \|u - \hat{u}\|_2$$

We know ALREADY that the error is given by the norm of the matrix (find out starting from  $\|p - \hat{p}\|_2$  and using)

$$\|p - \hat{p}\|_2 \leq \|V\| \|u - \hat{u}\|_2$$

L' As long as  $V$  is invertible, problem\* is absolutely well posed, and  $\|V^{-1}\|_*$  is not infinite

The op. norm is defined here as:  $\|V^{-1}\|_* := |\lambda(V^{-1})| = \max_i |\lambda_i(V^{-1})| = \|V^{-1}\|_2$  eigenvalues of  $V^{-1}$  op. norm induced by the 2-norm

To make the problem better conditioned, we need to minimize  $K_{abs} = \|V^{-1}\|_2$  that is  $V \equiv I$

then we have

$$\nabla = I \rightarrow v_j(x_i) = s_{ji}$$

which means we have a

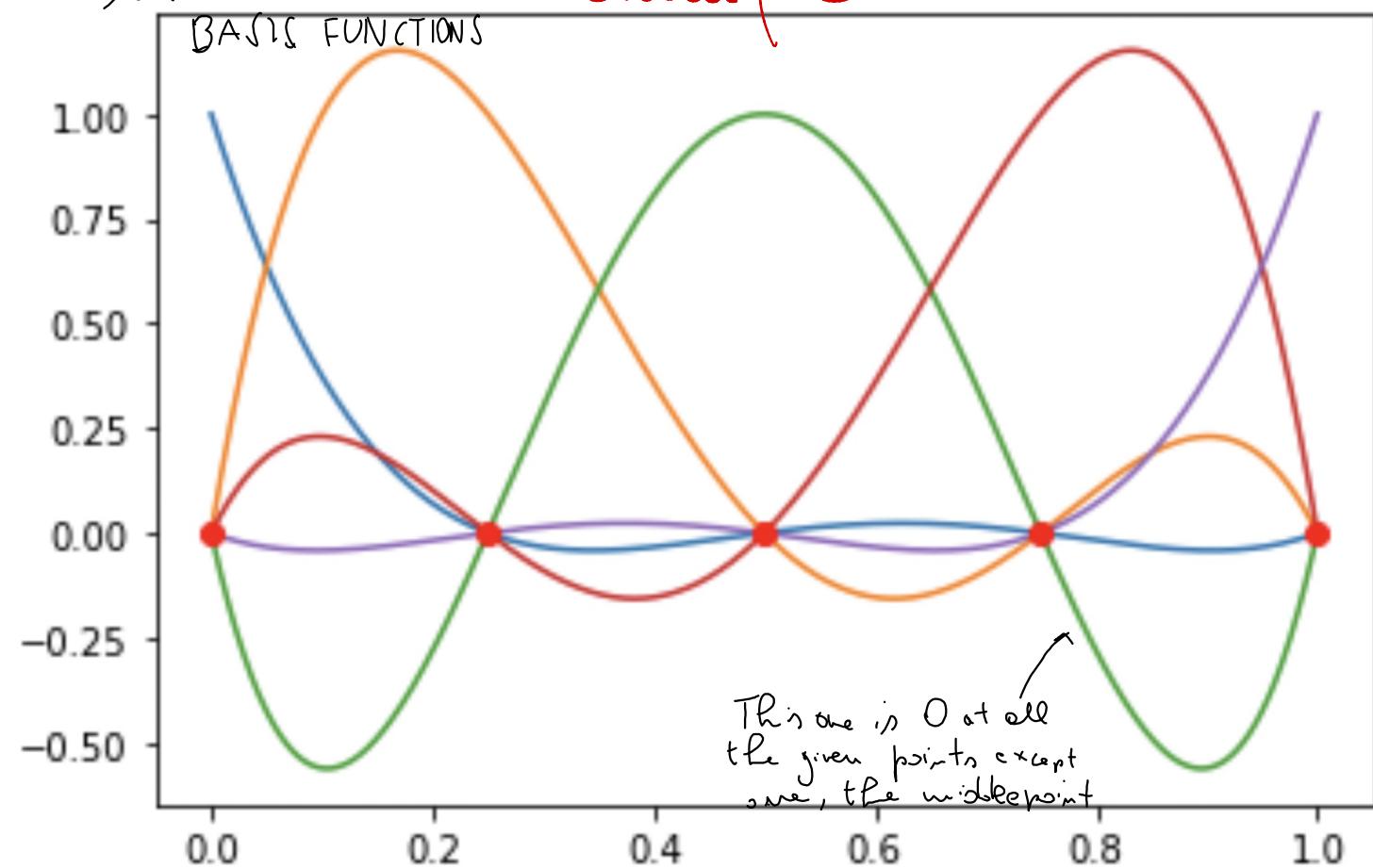
Polynomial that is zero at every point  $x_i$  except on  $x_5$  where it is 1. We call this polynomial

$$e_i := \prod_{\substack{j=0 \\ j \neq i}}^n \frac{(x - x_j)}{(x_i - x_j)}$$

poly of order n  
=>  $e_i(x_j \neq x_i) = 0, e_i(x_i) = 1$

Lagrange Basis

example with  $n = 5$



This example shows you how to choose the basis function such that

$$\nabla = I$$

And you can use directly these basis functions to construct the polynomial interpolation

$$p(x) = \sum_{j=0}^n u(x_j) v_j(x)$$

where the coefficient are exactly

$$p_j = u(x_j)$$

because  $\nabla^{-1} = I$

If you have columns of  $\nabla$  that are very similar to one another (f.e. because they get closer and closer to one another)  $\nabla$  will not be invertible (the condition number will be very bad!)

The condition number of Vandermonde matrix will deteriorate as we increase the number of points (simply because, the more points we have the similar all the basis functions will get!)

# Numerical Analysis — lec 4 L+H —

General interpolation problem from  $\mathcal{X} := C^0([a, b])$  to the space  $\mathcal{Y} := V := \text{span}\{v_i\}_{i=0}^n$  with fixed support points  $\{x_i\}_{i=0}^n$

$$\text{II: } C^0([a, b]) \longrightarrow V$$

$$u(x) \longrightarrow p(x) := \sum_{i=0}^n p^i v_i(x)$$

Such that  $p$  and  $u$  coincides in  $\{x_i\}_{i=0}^n$

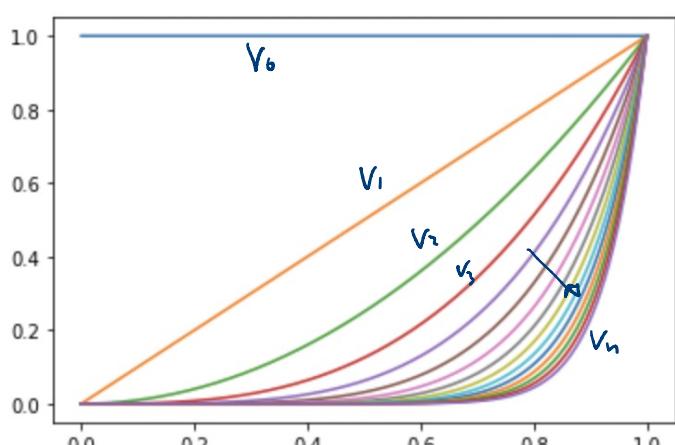
indeed  
this is a linear system  $\sum_{j=0}^n p^j v_j(x_i) = u(x_i)$

$\Downarrow_{ij} := v_j(x_i) \quad \Downarrow_{p} = u$

$$\{p^j = p^j\}$$

$$\{u(x_i) = u(x_i)\}$$

POWER BASIS IN  $[0, 1]$   
Example  $V : [a, b] = [0, 1]$



$$\text{Linearity implies } \text{II}(u+v) = \text{II}(u) + \text{II}(v)$$

$$\text{because } \text{II}(u) = (\mathbb{V}^{-1} u)^i v_i(x) \quad \text{II}(v) = (\mathbb{V}^{-1} v)^i v_i(x)$$

$$\text{II}(u+v) = (\mathbb{V}^{-1} (u+v))^i v_i(x)$$

the choice of this number of points is important (in principle you could have even higher  $n$ ) but then it would be more difficult to satisfy exactly the relationship which the points coincide at those higher  $n$  of points

The basis functions  $v_i \in C^0([a, b])$   
(basis functions can be chosen among polynomials, cos/sin, sigmoid)

$$v_i = \text{pow}(x, i) = x^i$$

if we fix  $n+1=5$   $\underline{p} := [1, 0, 0, 2, 3]$

$$\begin{aligned} p(x) &= 1 + 2x^3 + 3x^4 \\ &= p^i v_i(x) \end{aligned}$$

$$u+v = (u+v)$$

$$(u+v)_i = u(x_i) + v(x_i)$$

then here we use

Condition number for interpolation problem

We use as norm in  $C^0([a, b])$   $\|\cdot\|_\infty$ , the same in  $V$ :  
 explicit to be fact that the interval problem is the definition of the Chebyshev problem

$$\|\mathbb{I}(u) - \mathbb{I}(\hat{u})\|_\infty = \|\mathbb{I}(u - \hat{u})\|_\infty \leq \|\mathbb{I}\|_* \|u - \hat{u}\|_\infty$$

What is  $\|\mathbb{I}\|_*$ ?

Let's apply the definition

$$\|\mathbb{I}\|_* := \sup_{\text{domain of the function}} \text{norm of } V$$

"simplified" by the problem

$$\text{of } u \in C^0([a, b])$$

$$\|\mathbb{I}(u)\|_\infty = \sup_{\text{of } u \in X} \frac{\|\mathbb{I}(u)\|_\infty}{\|u\|_\infty}$$

by def. condition number

norm of matrix  $V$

$$\|\mathbb{I}\|_* = \sup_{\text{of } u \in X} \frac{\|\mathbb{I}(u)\|_\infty}{\|u\|_\infty} = \sup_{\text{of } u \in X} \frac{\|\mathbb{V}^{-1} u\|_\infty}{\|u\|_\infty} = \|\mathbb{V}^{-1}\|_\infty$$

$$\leq \sup_{\text{of } u \in X} \|\mathbb{V}^{-1}\|_\infty \|u\|_\infty \max_i \|v_i\|_\infty$$

$$\leq \|\mathbb{V}^{-1}\|_\infty \max_i \|v_i\|_\infty$$

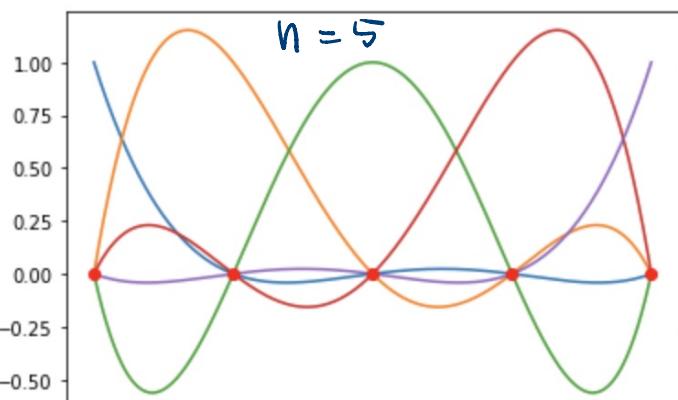
because for some the real minimum of the solution value of  $n$  is greater than all the values of  $n$  at  $x_i$  in the interval  $[a, b]$ . All roots in the interval  $[a, b]$  are simple. Therefore we can choose  $x_i$  with  $x_i = \frac{a+b}{n+1} i$ .  
 Two things:  $\{x_i\}_{i=1}^n$ ,  $\{v_i\}_{i=0}^n$  are

For power basis:

```

n = 1 : 1.0
n = 2 : 2.6180339887498953
n = 3 : 15.099657722502098
n = 4 : 98.86773850722759
n = 5 : 686.4349418185955
n = 6 : 4924.371056611224
n = 7 : 36061.16088021232
n = 8 : 267816.7009075794
n = 9 : 2009396.3800421846
n = 10 : 15193229.677753646
n = 11 : 115575244.54431371
n = 12 : 883478687.0721825
n = 13 : 6780588492.454725
n = 14 : 52214926084.1525
n = 15 : 403234616528.72504

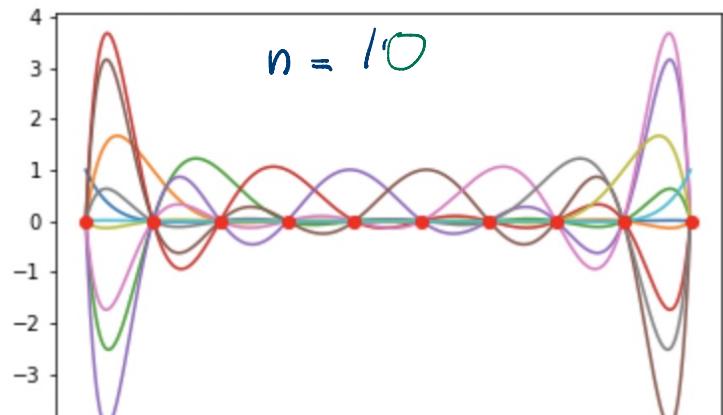
```

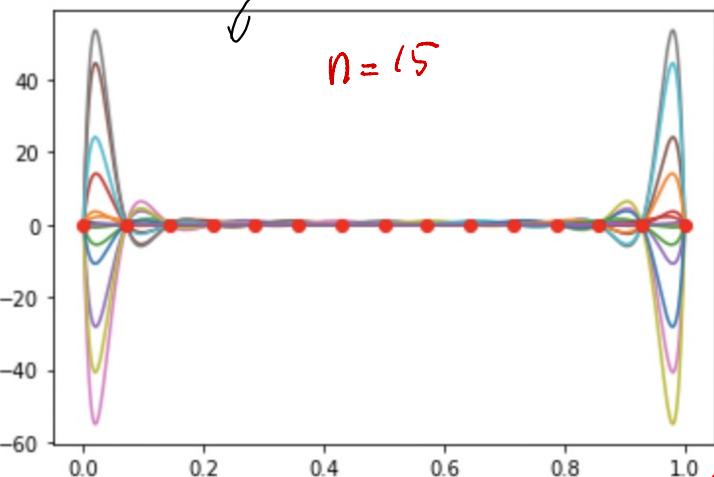


the Lagrange basis functions:  
 $\ell_i := \prod_{j=0}^n \frac{(x - x_j)}{(x_i - x_j)}$   
 for power basis  
 $\Rightarrow \mathbb{V}_{ij} := \ell_j(x_i) = \delta_{ij} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$

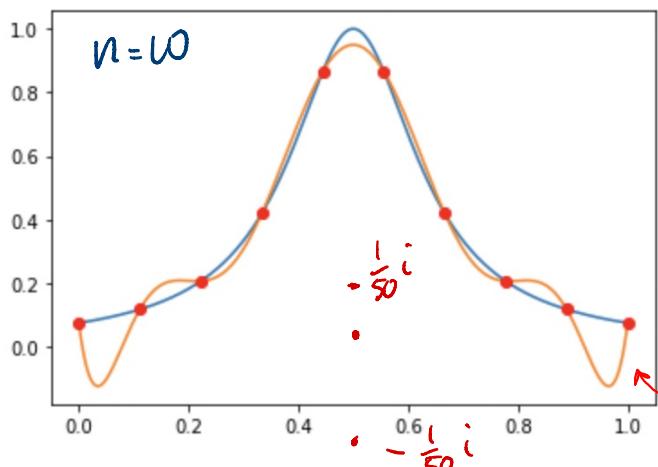
$\mathbb{V} = \mathbb{I}^\dagger$  by construction

We should choose  $x_i$  s.t.  $\max_i \|\ell_i\|_\infty$  is small





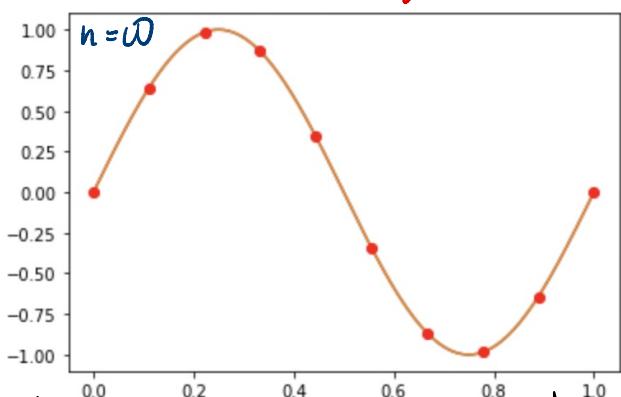
then you start getting bad oscillations  
 ↳ Is interpolation  
 a good strategy  
 for approximating  
 a function?



The orange line represents using  $n$  LAGRANGE interpolation  
 $\sum_{i=0}^n u(x_i) l_i(x) = p(x)$  with Lagrange basis

blue line is the interpolated function  
 $f(x) = \frac{1}{1 + 50(x - \frac{1}{2})^2}$

BAD BEHAVIOR HERE (away from the center while at the center things are good but this is FUNCTION DEPENDENT)  
 $f(x) = \sin(2\pi x)$



for sine results are perfect

Can we do better?

What can we say for smooth functions about interpolation polynomials of  $n$  intervals

Thus  $f \in C^{(n+1)}([a, b])$  if  $x_i \in (a, b)$ ,  $\forall x \in (a, b) \exists \xi \in (a, b)$   
 where  $n+1$  pts.  $\Rightarrow$   $n+1$  derivatives

$$(f - I^n f)(x) = \frac{1}{(n+1)!} f^{(n+1)}(\xi) \omega(x)$$

where  $\omega(x)$  is the characteristic polynomial of  $\{x_i\}_{i=1}^n$

$$\omega(x) = \prod_{i=0}^n (x - x_i)$$

can we make these two numbers small?

First consequence:

$$\|I^n f - f\|_\infty \leq \|w\|_\infty \|f^{(n+1)}\|_\infty$$

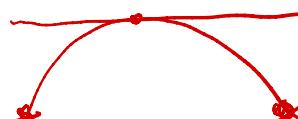
If you control the  $n+1$  derivative of the function ANYWHERE  $\rightarrow$  you have a bound on the norm of the  $n+1$  derivative and on the norm of the  $w$  function

Proof: Let  $x$ , define  $G(t)$  s.t. it is equal to 0 in all points, in which the polynomials and the function coincide plus the point  $x$  where the polynomials and the function coincide. Then  $G(t) = (f(t) - p(t)) w(x) - (f(x) - p(x)) w(t)$

$$\text{where } p(t) = \sum_{i=0}^n f(x_i) e_i(t) = (\text{If})(t) \quad \text{~\textcircled{1}~ interpolation function}$$

$G(t)$  has  $n+2$  zeros:  $\{x_i\}_{i=0}^n \cup \{x\}$

According to Rolle's theorem:  $\exists \xi \in (a, b)$  s.t.  $\frac{d^{n+1} G(\xi)}{dt^{n+1}} = 0$



With two zeros  $\Rightarrow \exists \xi \in \text{int}(a, b)$  s.t.  $f'(\xi) = 0$

writing explicitly the derivative

$$\frac{d^{n+1}}{dt^{n+1}} G(\xi) = f^{(n+1)}(\xi) \cdot w(x) - (f(x) - p(x))(n+1)! = 0$$

$\frac{d^{n+1}}{dt^{n+1}} (w(t))$

$\Rightarrow f(x) - (\text{If})(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} w(x) \quad \square$  QED

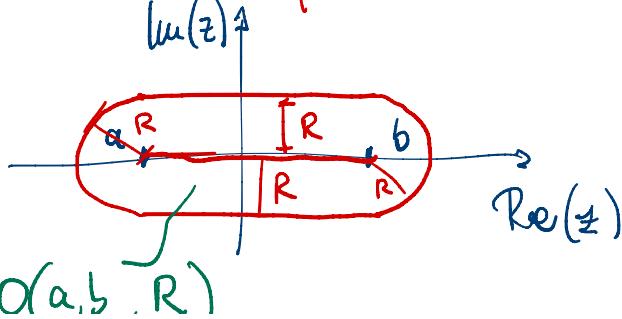
How to apply this?

Theorem 2: If  $f$  is analytically extendible in an oval  $O(a, b, R)$  with  $R > 0$

Then  $\|f^{(n+1)}\|_\infty \leq \frac{(n+1)!}{R^{n+1}} \|\tilde{f}\|_{L^\infty(O(a,b,R))}$

Analytically extendible means:  $\tilde{f}$

Take  $z$ , and replace  $t$  with  $z \in \mathbb{C}$



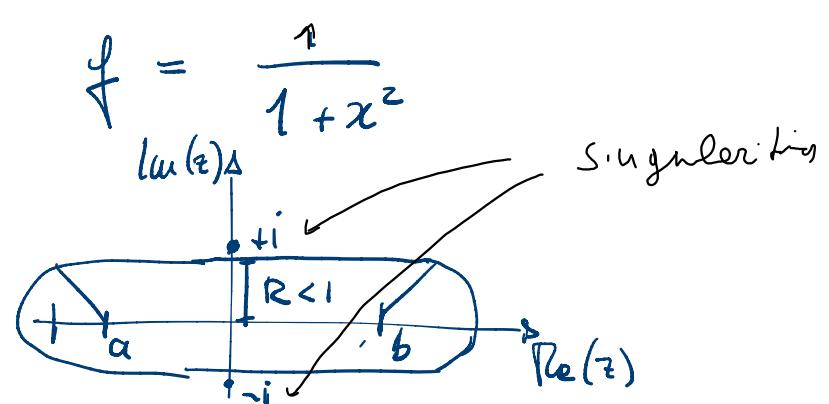
We are asking that  $f \in C^\infty$  in  $O(a, b, R)$

EXAMPLE

Runge function:

extension to  $\mathbb{C}$

$$f(z) = \frac{1}{1+z^2}$$



$f$  is analytical in  $O(-\infty, +\infty, 1-\varepsilon)$   $\forall \varepsilon > 0$

What are the consequences?

$$\|f^{(n)}\| \leq \frac{(n+1)!}{R^{n+1}} \|f\|_\infty$$

$R < 1$  stabilizing  
at finite singularities

So, by the previous theorem:

$$\|\rho - f\|_\infty \leq \|f^{(n+1)}\|_\infty \frac{\|\omega\|_\infty}{(n+1)!} \leq \frac{(n+1)!}{(n+1)!} \frac{\|\omega\|_\infty}{R^{n+1}} \|f\|_\infty$$

where

$$\|\omega\|_\infty = \left\| \prod_{i=0}^n (x-x_i) \right\| \leq \underbrace{|b-a|^{n+1}}_{\text{since } x \in (a,b)} \quad x \in (a,b)$$

$$\|\rho - f\|_\infty \leq \left( \frac{|b-a|}{R} \right)^{n+1} \|f\|_\infty$$

$\hookrightarrow$  good only when  $|b-a| < R$ , otherwise  $\|\rho - f\|_\infty$  is NOT BOUNDED

$\Rightarrow$  Runge:

$$\frac{1}{1+x^2} \quad \text{in } (-1, 1)$$

everything ok

Can we do better?

Best Approximation. Definition:

Restrict ourselves

$\mathcal{X}$  Banach, reflexive, strictly convex  
✓ subspace of  $\mathcal{X}$

We call  $p$  the best approximation in  $V$  of  $f$  in  $\mathcal{X}$   
when :

$$\|f-p\| \leq \|f-q\| \quad \forall q \in V$$

in other words

$$\|f-p\| = E(f) = \inf_{q \in V} \|f-q\|$$

Theo:  $\exists! p$  satisfying these  
certain (mandatory conditions)

How does polynomial interpolation compare to best approximation?

We can say that

$$\|f - I^n f\|_\infty = \|f - p + p - I^n f\|_\infty =$$

add on subtract best approx

$$= \|f - p + I^n(p - f)\|_\infty$$

$p$  is a poly of od n  
 $\Rightarrow$  it's poly interpolation  
is itself

$$\text{triangle inequality} \leq \|f - p\|_\infty + \|I^n\|_* \|p - f\|_\infty$$

$$\text{property of norms} \leq (1 + \|I^n\|_*) \|p - f\|_\infty$$

For Lagrange interpolation:

$$\|I^n_*\| := \sup_{\substack{u \in C^0([a, b]) \\ \|\sum_{i=0}^n |e_i(x)|\|_\infty}} \frac{\|\sum_{i=0}^n u(x_i) e_i(x)\|_\infty}{\|u\|_\infty} \leq$$

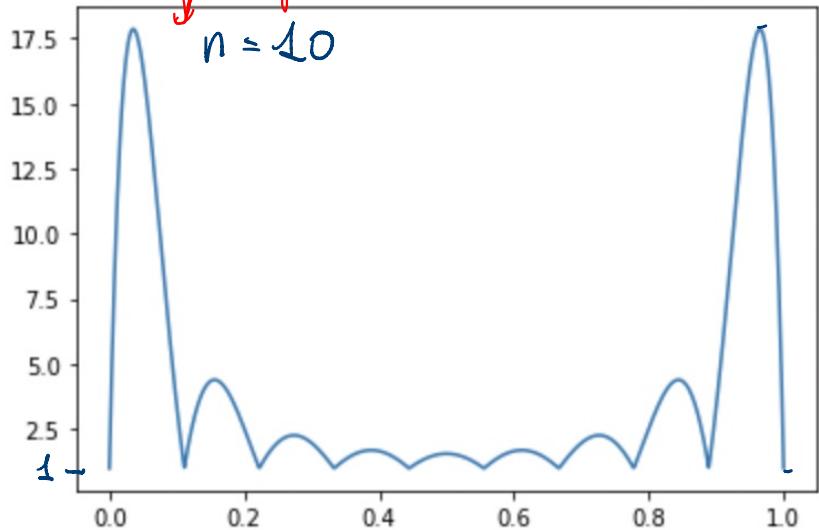
how away we  
are from best  
interpolation

$$= \|\sum_{i=0}^n |e_i(x)|\|_\infty \quad \text{def } L(x) := \sum_{i=1}^n |e_i(x)|$$

LEBESGUE Funktion

$$\Rightarrow \|\mathcal{I}_*^n\| = \|\mathcal{L}\|_\infty$$

Lebesgue functions are such that  $\Rightarrow$  bad results. Indeed  $\|\mathcal{L}\|_\infty$  for equispaced points



$\|\mathcal{L}\|_\infty$  for equispaced points

$$\|\mathcal{L}\|_\infty \leq \frac{2^{n+1}}{\pi} \text{ grows quickly}$$

A collection of points:  $\{x_i\}_{i=0}^n$  is a list of  $n+1$  points with increasing size -

Endos

$\nexists$  collection of points  $\{x_i\}_{i=0}^n$

$\exists c > 0$  s.t.

Faber

$$\|\mathcal{L}\|_\infty \geq \frac{2}{\pi} \log(n+1) - c$$

of independently of how you choose the points, at least you get logarithmic with the number of the points

$\nexists$  collection of points  $\{x_i\}_{i=0}^n$

$\exists \delta$  s.t.

$$\lim_{n \rightarrow \infty} \|\mathcal{I}^n f - g\|_\infty \rightarrow \infty$$

( $\hookrightarrow$  you can always find a pt that destroys all)

The best is  $\alpha := \{x_i\}_{i=0}^n$  s.t.

$$\|\mathcal{L}^\alpha\|_\infty \leq \|\mathcal{L}^\beta\|_\infty$$

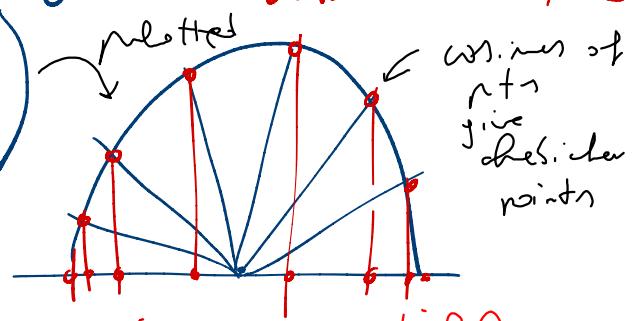
$\nexists \beta = \{x_i\}_{i=0}^n$   
such  $\alpha := \{x_i\}_{i=0}^n$

are called Chebyshev points:

between  $[-1, 1]$

defined as  $x_i = \cos\left(\frac{(2i+1)\pi}{2n+2}\right)$

$$\|\mathcal{L}^\alpha\|_\infty \leq \frac{2}{\pi} \log(n+1) + 1$$



Best case scenario

still grows exponentially

$$\frac{2}{\pi} \log(n+1) - c \leq \|\mathcal{L}^\alpha\|_\infty \leq \frac{2}{\pi} \log(n+1) + 1$$

$\Rightarrow$  INTERPOLATION IS GOOD BUT NOT  $\pi$  SO WEBS

# Nonlinear equations



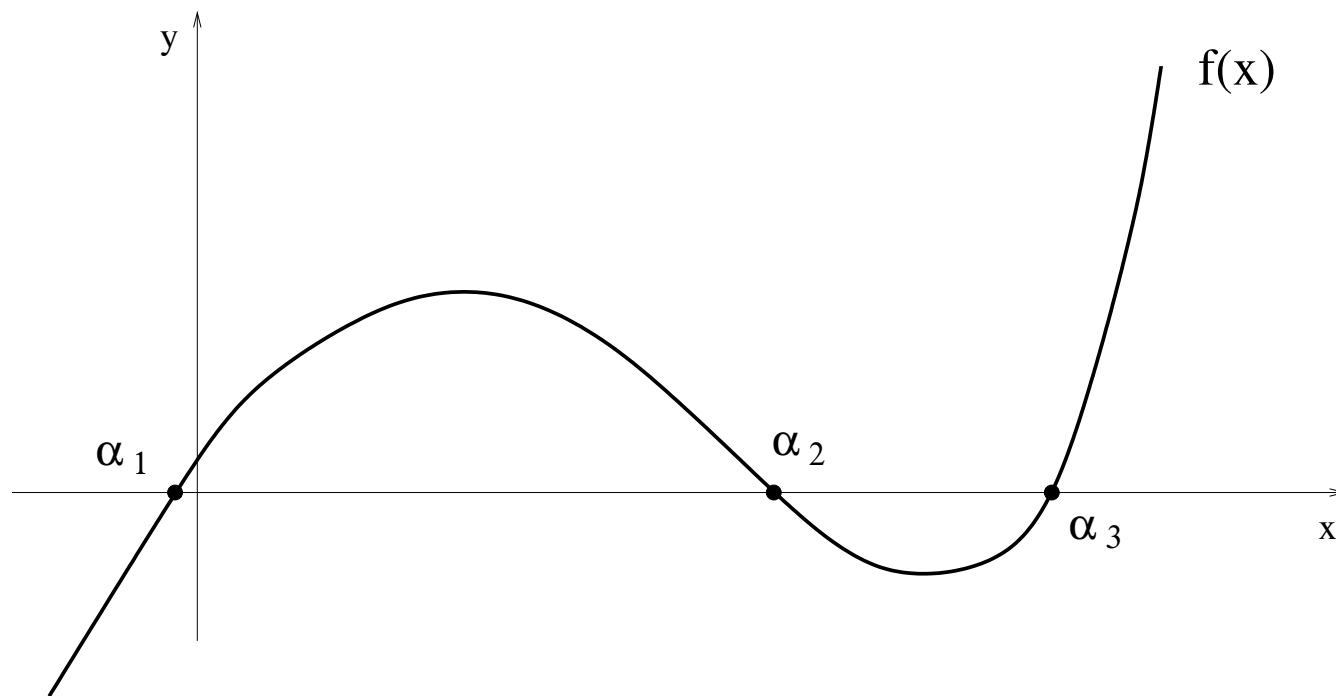
## Numerical Analysis

Profs. Gianluigi Rozza - Luca Heltai

2019-SISSA mathLab Trieste

# Nonlinear equations

**Objective:** Find the root of scalar (or vector) non-linear functions, i.e., find  $\alpha \in \mathbb{R}$  such that  $f(\alpha) = 0$ .



# Examples and motivation

**Example 1** (Interest rates). We want to compute the mean interest rate  $I$  of a portfolio over several years. We invest  $v = 1000$  Euro every year. After 5 years we end up with  $M = 6000$  Euro.

The relation between  $M$ ,  $v$ ,  $I$  and the number of years  $n$  is

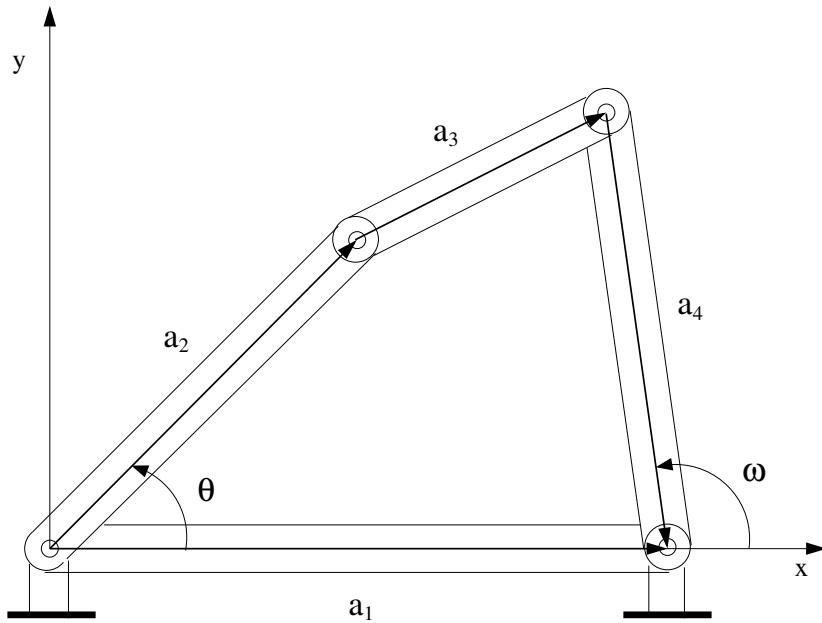
$$M = v \sum_{k=1}^n (1 + I)^k = v \frac{1 + I}{I} [(1 + I)^n - 1]$$

This can be rewritten as: *find  $I$  such that*

$$f(I) = M - v \frac{1 + I}{I} [(1 + I)^n - 1] = 0 \quad (1)$$

Therefore we have to solve a nonlinear equation in  $I$ , for which we can't find an analytical solution.

**Example 2** (Rods system). Let us consider the mechanical system represented by four rigid rods



For any admissible angle  $\omega$ , we want to compute the angle  $\theta$  between  $\mathbf{a}_1$  and  $\mathbf{a}_2$ .

Thanks to the vector identity

$$\mathbf{a}_1 - \mathbf{a}_2 - \mathbf{a}_3 - \mathbf{a}_4 = 0$$

and keeping  $\mathbf{a}_1$  on the  $x$ -axis, we can derive the following equation: involving  $\omega$  and  $\theta$ :

$$\frac{a_1}{a_2} \cos(\omega) - \frac{a_1}{a_4} \cos(\theta) - \cos(\omega - \theta) = -\frac{a_1^2 + a_2^2 - a_3^2 + a_4^2}{2a_2 a_4} \quad (2)$$

where  $a_i$  is the length of the  $i$ th rod. Equation (2) is nonlinear and can be solved only for particular values of  $\omega$ . For a general  $\omega$  it is not possible to find an analytic solution.

**Example 3** (State equation of a gas). We want to determine the volume  $V$  occupied by a gas at temperature  $T$  and pressure  $p$ . The state equation (i.e. the equation that relates  $p$ ,  $V$  et  $T$ ) is

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

where  $a$  and  $b$  are two coefficients that depend on the specific gas,  $N$  is the number of molecules which are contained in the volume  $V$  and  $k$  is the Boltzmann constant. We need therefore to solve a non-linear equation whose root is  $V$ .

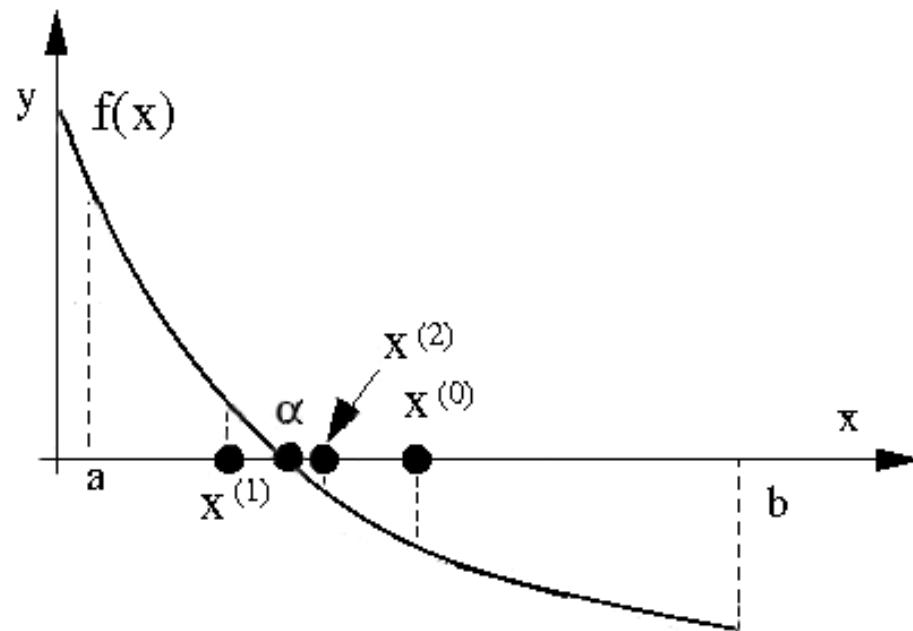


# Bisection method

(Book Sec. 2.2)

This method is used to compute the root of a **continuous** function  $f$ , i.e., the point  $\alpha$  such that  $f(\alpha) = 0$ . We can build a sequence  $x^{(0)}, x^{(1)}, \dots, x^{(k)}, (x^{(0)}$  such that  $\lim_{k \rightarrow \infty} x^{(k)} = \alpha$ .

We assume that  $f : (a, b) \rightarrow \mathbb{R}$  and  $a < b$ . If  $f(a)f(b) < 0$ , since  $f$  is continuous, we know that there exists (at least) one root  $\alpha$  of  $f$  in the interval  $(a, b)$ .



Then

1. We set  $a^{(0)} = a$ ,  $b^{(0)} = b$  and  $x^{(0)} = \frac{a^{(0)} + b^{(0)}}{2}$ ,
2. if  $f(x^{(0)}) = 0$ , then  $x^{(0)}$  is the zero.
3. if  $f(x^{(0)}) \neq 0$ , then:
  - (a) if  $f(x^{(0)})f(a^{(0)}) > 0 \Rightarrow$  the zero  $\alpha \in (x^{(0)}, b^0)$  and we define  $a^{(1)} = x^{(0)}$ ,  $b^{(1)} = b^{(0)}$  and  $x^{(1)} = (a^{(1)} + b^{(1)})/2$
  - (b) if  $f(x^{(0)})f(a^{(0)}) < 0 \Rightarrow$  the zero  $\alpha \in (a^{(0)}, x^{(0)})$  and we define  $b^{(1)} = x^{(0)}$ ,  $a^{(1)} = a^{(0)}$  et  $x^{(1)} = (a^{(1)} + b^{(1)})/2$

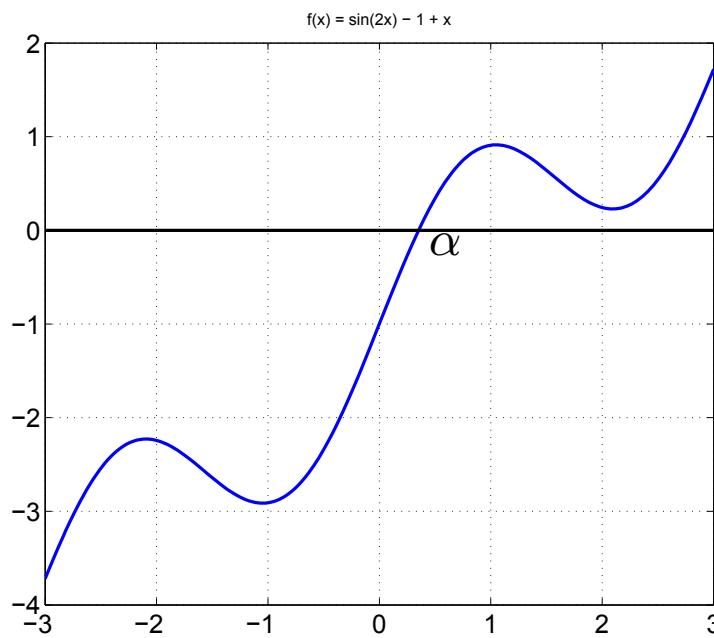
By the divisions of this type, we construct the sequence  $x^{(0)}, x^{(1)}, \dots, x^{(k)}$  that satisfies for all  $k$ ,

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

because we divide the interval by 2 at every step, starting from step 0

**Example 4.** We want to find the zero of the function  $f(x) = \sin(2x) - 1 + x$ . We draw the graph of the function  $f$  using the following commands in Matlab/Octave:

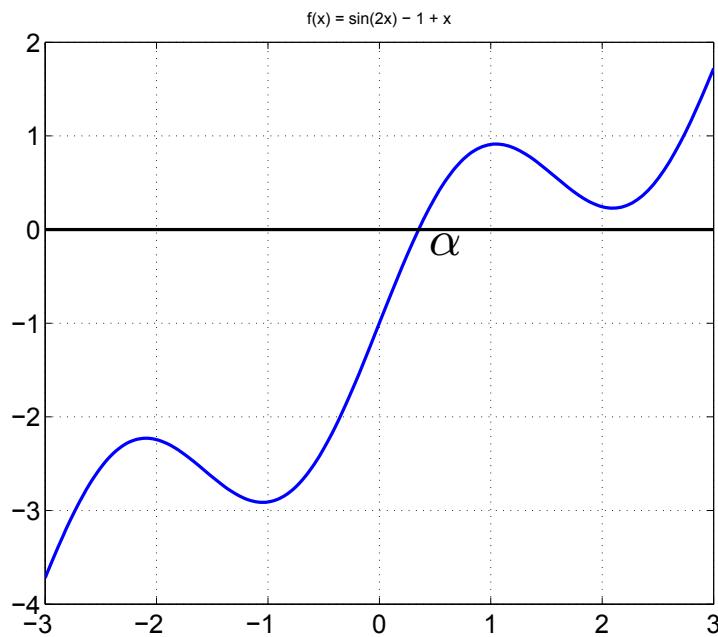
```
>> f = @(x) sin(2*x) - 1 + x;
>> x=-3:0.1:3;
>> plot(x,f(x)); grid on;
```



If we apply the bisection method in the interval  $[-1, 1]$  with a tolerance  $10^{-8}$  and maximum number of iterations  $k_{max} = 1000$

```
>> [zero,res,niter]=bisection(f,-1,1,1e-8,1000);
```

We find the value  $\alpha = 0.352288462$  after 27 iterations.



## APRIORI KNOWLEDGE OF # OF MAX ITERATION

Assume stopping criterion is based on size of interval.

At  $k$ :

$$\frac{b-a}{2^k} \leq \varepsilon \quad \rightarrow \text{solve for } k: \quad k \geq \log_2 \frac{b-a}{\varepsilon} = \frac{\ln(b-a) - \ln(\varepsilon)}{\ln 2}$$

$\uparrow$   
set tolerance

# Newton's method

(Chapt. 2.3 of the book)

Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a differentiable function.

Let  $x^{(0)}$  be an initial guess. Let us consider the equation  $y(x)$  which passes through the point  $(x^{(k)}, f(x^{(k)}))$  and which has the slope  $f'(x^{(k)})$ :

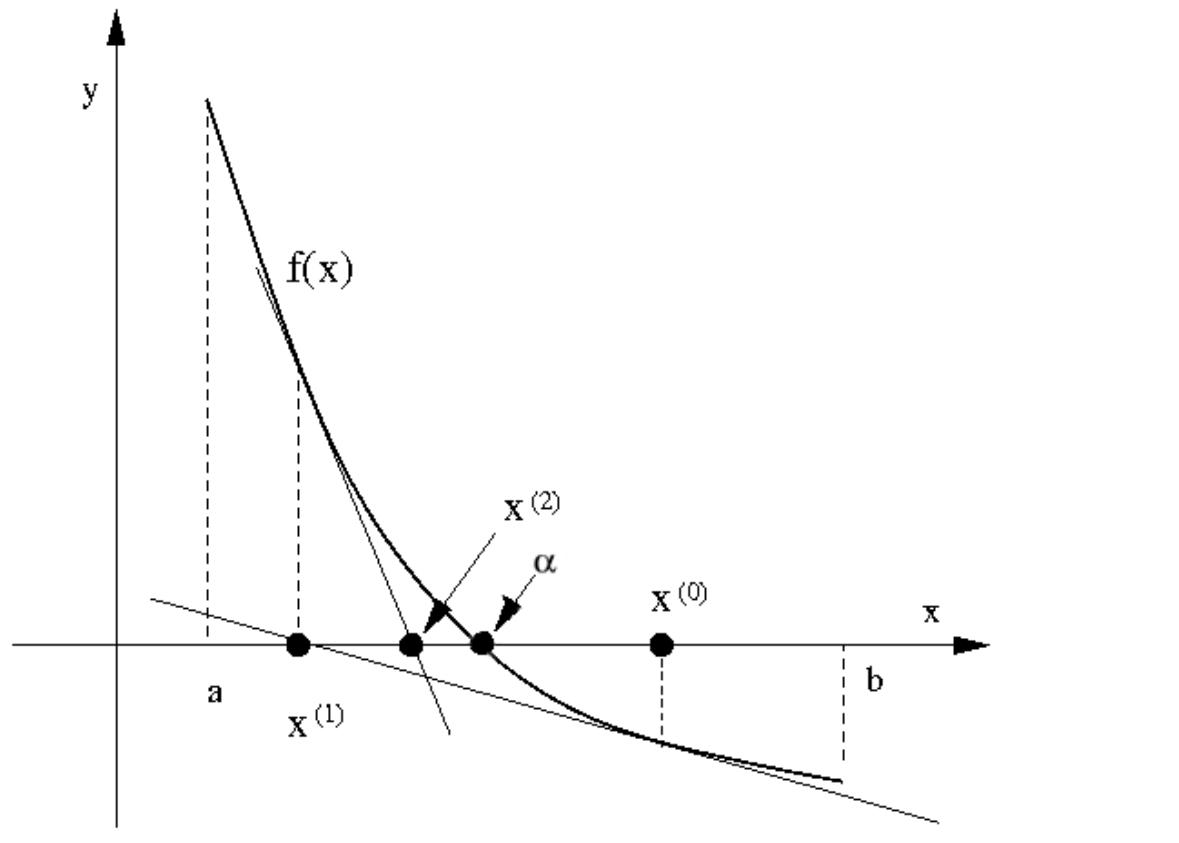
$$y(x) = f'(x^{(k)})(x - x^{(k)}) + f(x^{(k)}).$$

We define  $x^{(k+1)}$  by the point where this line intersects the axis  $x$ , i.e.  $y(x^{(k+1)}) = 0$ . We deduce that:

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}, \quad k = 0, 1, 2, \dots \quad (3)$$

# Newton's method

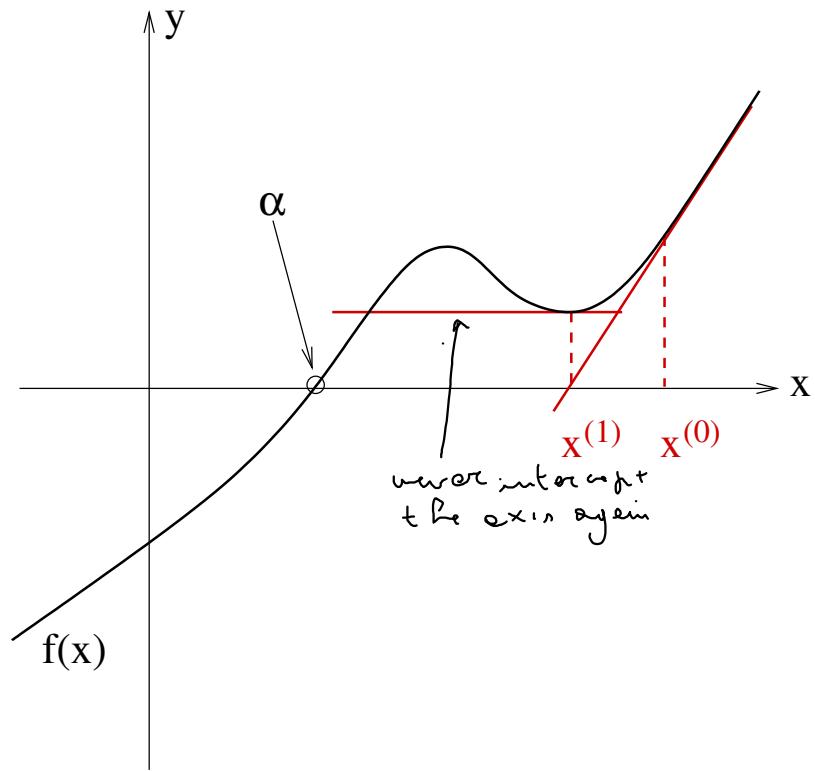
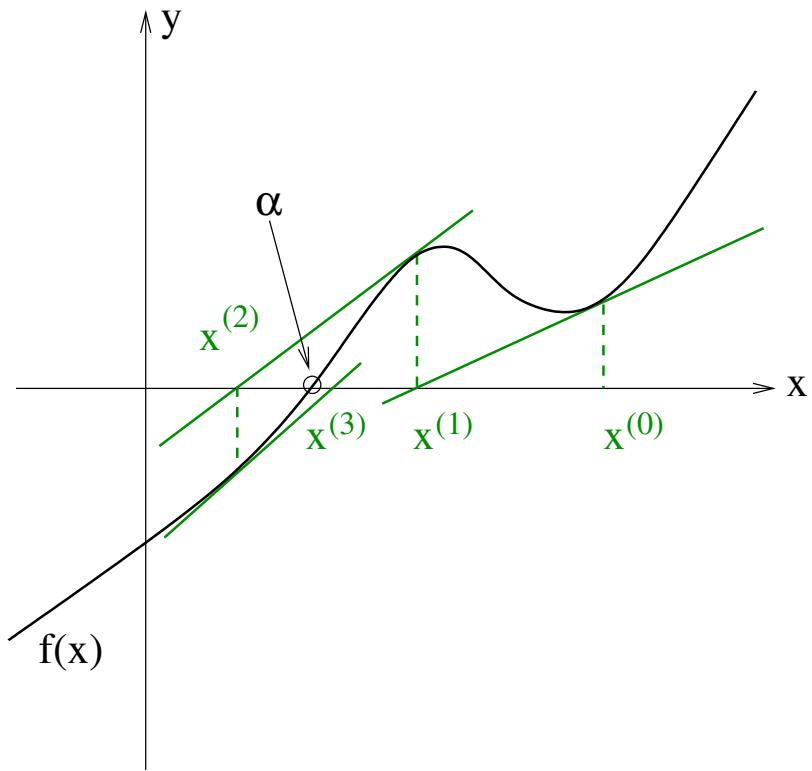
Starting from the point  $x^{(0)}$ , the sequence  $\{x^{(k)}\}$  converges to the root of  $f$



# Convergence?

Does this method always converge?

- it depends on the **property of the function**;
- and on the **initial guess**.



# Fixed point iterations.

(Chapt. 2.4 in the book)

A general method for finding the roots of a nonlinear equation  $f(x) = 0$  is the transformation<sup>of the problem in</sup> an equivalent problem  $x - \phi(x) = 0$ , where the auxiliary function  $\phi : [a, b] \rightarrow \mathbb{R}$  must have the following property :

$$\phi(\alpha) = \alpha \quad \text{if and only if} \quad f(\alpha) = 0.$$

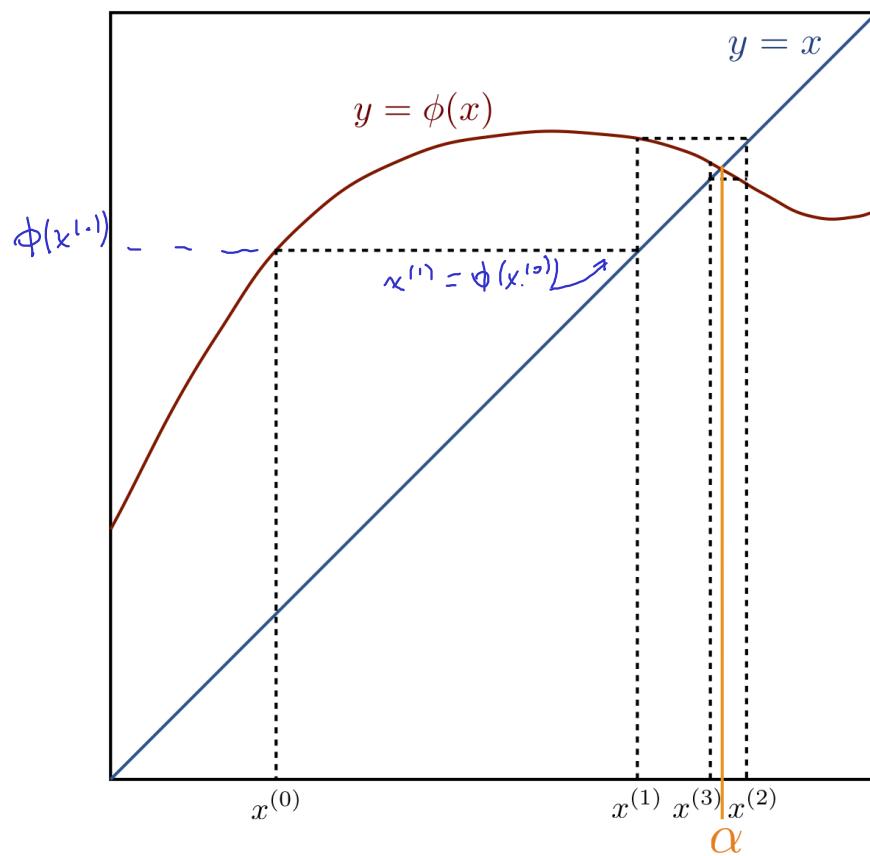
*↙ and  $\phi$  is called CONTRACTION*

The point  $\alpha$  is called *a fixed point* of  $\phi$ . Searching the zeros of  $f$  is reduced to the problem of determining the fixed points of  $\phi$ .

**Idea :** It could be computed by the following algorithm:  $x^{(k+1)} = \phi(x^{(k)})$ ,  $k \geq 0$ . Indeed, if  $x^{(k)} \rightarrow \alpha$  and if  $\phi$  is continuous on  $[a, b]$ , then the limit  $\alpha$  satisfies  $\phi(\alpha) = \alpha$ .

$$x^{(k+1)} = \phi(x^{(k)}) \xrightarrow{x^{(k)} \rightarrow \alpha} \phi(\alpha) \xrightarrow{\text{follows from the definition of ft } \phi} \alpha$$

Starting from the point  $x^{(0)}$ , the sequence  $\{x^{(k)}\}$  converges to the fixed point  $\alpha$

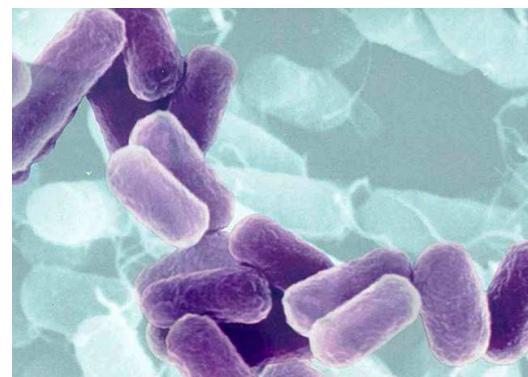


## Example 5 (Population dynamics).

In the survey population (e.g. bacteria), we want to establish a link between the number of individuals in the generation  $x$  and the number of individuals in the next generation  $x^+$ :

$$x^+ = \phi(x) = xR(x), \quad (4)$$

where  $R(x)$  represents the rate of growth (or decay) of the considered population.



Several models are available for  $R(x)$ :

- The Malthusian growth model (Thomas Malthus 1766-1834),

$$x^+ = \phi_1(x) = xR_1(x) \text{ with } R_1(x) = r, \quad r \text{ is a positive constant}$$

- the model of growth with limited resources (Pierre François Verhulst, 1804-1849),

$$x^+ = \phi_2(x) = xR_2(x) \text{ with } R_2(x) = r/(1 + x/K), \quad r > 0, K > 0$$

$K$  may be some kind of measurement for  
 available resources  $\rightarrow x/K$  is the number  
 of members of population per resource

that improve the Malthausian growth model by taking into account the growth of a population is limited by the resources.

- the predator/prey model

$$x^+ = \phi_3(x) = xR_3(x) \text{ with } R_3(x) = rx/(1 + (x/K)^2)$$

that represents the change of the Verhulst model by presence of an antagonist population.

The dynamics of a population is defined by a iterative process, starting from a given initial guess ( $x^{(0)}$ ),

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

where  $x^{(k)}$  represents the number of individuals in  $k$ -th generation.  
 In addition, the steady states (equilibriums)  $x^*$  of a considered population are identified by the following problem,

$$x^* = \phi(x^*), \tag{5}$$

or equivalently,

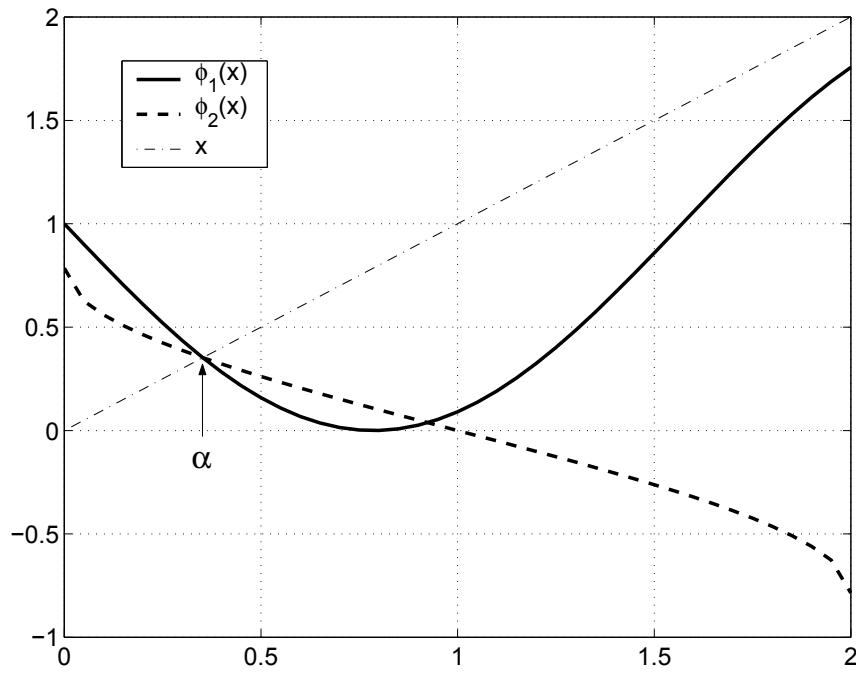
$$x^* = x^* R(x^*), \quad \text{i.e.} \quad R(x^*) = 1. \tag{6}$$

In both cases, we ask to solve a non-linear problem. In particular the problem is called the fixed point problem.

**Example 6.** We consider the equation  $f(x) = \sin(2x) - 1 + x = 0$ . We can rewrite it in two different fashions:

$$x = \phi_1(x) = 1 - \sin(2x)$$

$$x = \phi_2(x) = \frac{1}{2} \arcsin(1 - x), \quad 0 \leq x \leq 1$$



**Proposition 1.** (*Global convergence*)

that is: the pair  $(x, \phi(x))$   
 belongs to a square

1. Assume that  $\phi(x)$  is continuous on  $[a, b]$  and such that  $\phi(x) \in [a, b]$  for all  $x \in [a, b]$ ; then there exists at least one fixed point  $\alpha \in [a, b]$  of  $\phi$ .

2. If  $\exists L < 1$  such that  $|\phi(x_1) - \phi(x_2)| \leq L|x_1 - x_2| \quad \forall x_1, x_2 \in [a, b]$ ,

then there exists a unique fixed point  $\alpha \in [a, b]$  and the sequence

$x^{(k+1)} = \phi(x^{(k)})$ ,  $k \geq 0$  converges to  $\alpha$ , for any initial guess  $x^{(0)} \in [a, b]$ .

*Proof.*

$$\phi(x) \in [a, b] \quad \forall x \in [a, b]$$

because  $\phi(x)$  and  $x$  are continuous

1. The function  $g(x) = \phi(x) - x$  is continuous in  $[a, b]$  and, thanks to assumption made on the range of  $\phi$ , it holds  $g(a) = \phi(a) - a \geq 0$  and  $g(b) = \phi(b) - b \leq 0$ . By applying the theorem of zeros of continuous functions, we can conclude that  $g$  has at least one zero in  $[a, b]$ , i.e.  $\phi$  has at least one fixed point in  $[a, b]$ .
2. Indeed, should two different fixed points  $\alpha_1$  and  $\alpha_2$  exist, then

Applying I property of  $\phi$

$$|\alpha_1 - \alpha_2| = |\phi(\alpha_1) - \phi(\alpha_2)| \leq L|\alpha_1 - \alpha_2| < |\alpha_1 - \alpha_2|,$$

Roll?

which cannot be. There exists a unique fixed point  $\alpha \in [a, b]$  of  $\phi$ . □

Let  $x^{(0)} \in [a, b]$  and  $x^{(k+1)} = \phi(x^{(k)})$ . We have

$$0 \leq |x^{(k+1)} - \alpha| = \underbrace{|\phi(x^{(k)}) - \phi(\alpha)|}_{\text{1}} \leq L|x^{(k)} - \alpha| \leq \dots \leq L^{k+1}|x^{(0)} - \alpha|,$$

↗ by applying 1 and 2  
 iteratively

i.e.

$$\frac{|x^{(k)} - \alpha|}{|x^{(0)} - \alpha|} \leq L^k.$$

Because  $L < 1$ , for  $k \rightarrow \infty$ , we obtain

$$\lim_{k \rightarrow \infty} |x^{(k)} - \alpha| \stackrel{\text{dk, because } |x^{(0)} - \alpha| \text{ is constant}}{\leq} \lim_{k \rightarrow \infty} L^k = 0.$$

So,  $\forall x^{(0)} \in [a, b]$ , the sequence  $\{x^{(k)}\}$  defined by  $x^{(k+1)} = \phi(x^{(k)})$ ,  $k \geq 0$  converges to  $\alpha$  when  $k \rightarrow \infty$ .

### Remark 1.

If  $\phi(x)$  is differentiable in  $[a, b]$  and

$\exists K < 1$  such that  $|\phi'(x)| \leq K \ \forall x \in [a, b]$ ,

*just a matter  
of scale*

then the condition 2 of the proposal (1) is satisfied. This assumption is stronger, but is more often used in practice because it is easier to check.

**Definition 1.** For a sequence of real numbers  $\{x^{(k)}\}$  that converges,  $x^{(k)} \rightarrow \alpha$ , we say that the convergence to  $\alpha$  is **linear** if exists a constant  $C < 1$  such that, for  $k$  that is large enough

$$|x^{(k+1)} - \alpha| \leq C |x^{(k)} - \alpha|.$$

If exists a constant  $C > 0$  such that the inequality

$$|x^{(k+1)} - \alpha| \leq C |x^{(k)} - \alpha|^2$$

is satisfied, we say that convergence is **quadratic**.

In general, the convergence is **with order  $p$** ,  $p \geq 1$ , if exists a constant  $C > 0$  (with  $C < 1$  when  $p = 1$ ) such that the following inequality is satisfied

  
 =>  
 $\overset{C>1}{\text{dove}} \text{mt}$   
 converg

$$|x^{(k+1)} - \alpha| \leq C |x^{(k)} - \alpha|^p.$$

**Proposition 2.** (*Local convergence - Theorem 2.1 in the book*)

Let  $\phi$  be a continuous and **differentiable** function on  $[a, b]$  and  $\alpha$  be a fixed point of  $\phi$ . If  $|\phi'(\alpha)| < 1$ , then there exists  $\delta > 0$  such that, for all  $x^{(0)}$ ,  $|x^{(0)} - \alpha| \leq \delta$ , the sequence  $\{x^{(k)}\}$  defined by  $x^{(k+1)} = \phi(x^{(k)})$  converges to  $\alpha$  when  $k \rightarrow \infty$ .

Moreover, it holds

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

Note that, if  $0 < |\phi'(\alpha)| < 1$ , then for any constant  $C$  such that  $|\phi'(\alpha)| < C < 1$ , if  $k$  is large enough, we have:

$$|x^{(k+1)} - \alpha| \leq C |x^{(k)} - \alpha|.$$

Proof: expand  $\phi(x^{(k)})$  around  $\alpha$ , proceed as in proof of next part

$\frac{x^{(k+1)} - \alpha}{x^{(k)} - \alpha}$  will reach  $\phi'(\alpha)$   
 from above or below  
 $\Rightarrow$  a number, it will be  $\leq C$

**Proposition 3.** (Proposition 2.2 in the book) *use the result of the previous preparation*

Let  $\phi$  be a twice differentiable function on  $[a, b]$  and  $\alpha$  be a fixed point of  $\phi$ . Let us consider that  $x^{(0)}$  converges locally to  $\alpha$ . If  $\phi'(\alpha) = 0$  and  $\phi''(\alpha) \neq 0$ , then the fixed point iterations converges with order 2 and

QUADRATIC CONVERGENCE

$$\lim_{k \rightarrow \infty} \frac{x^{(k+1)} - \alpha}{(x^{(k)} - \alpha)^2} = \frac{\phi''(\alpha)}{2}.$$

*Proof.* Using the Taylor series for  $\phi$  with  $x = \alpha$ , we have

$$x^{(k+1)} - \alpha = \phi(x^{(k)}) - \phi(\alpha) = \underbrace{\phi'(\alpha)}_0 (x^{(k)} - \alpha) + \frac{\phi''(\eta)}{2} (x^{(k)} - \alpha)^2$$

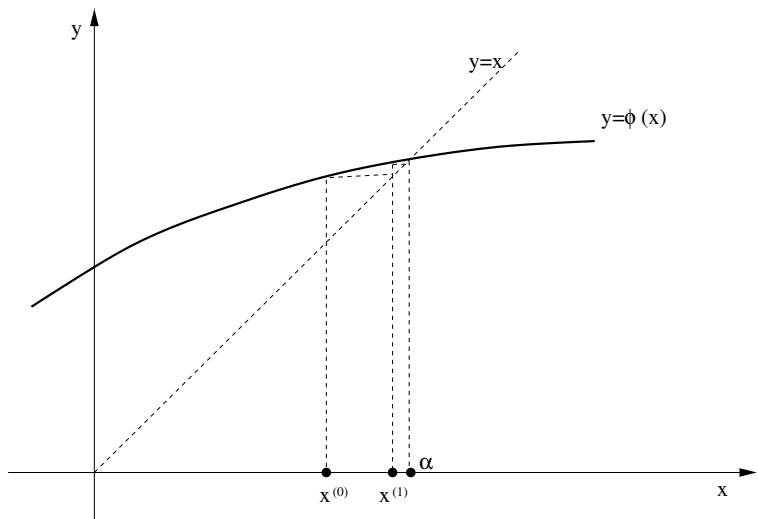
where  $\eta$  is between  $x^{(k)}$  and  $\alpha$ . So, we have

$$\lim_{k \rightarrow \infty} \frac{x^{(k+1)} - \alpha}{(x^{(k)} - \alpha)^2} = \lim_{k \rightarrow \infty} \frac{\phi''(\eta)}{2} = \frac{\phi''(\alpha)}{2}.$$

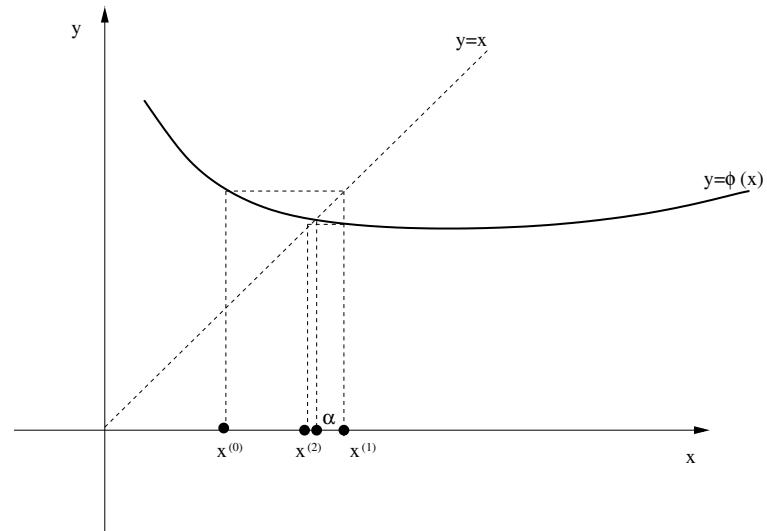
If we didn't have this, Taylor would be meaningless since  $x^{(k+1)}$  and  $x^{(k)}$  could be anywhere but near  $\alpha$ .

Some examples on how the value  $| \phi'(\alpha) |$  influences the convergence  
**Convergent cases:**

$$0 < \phi'(\alpha) < 1,$$



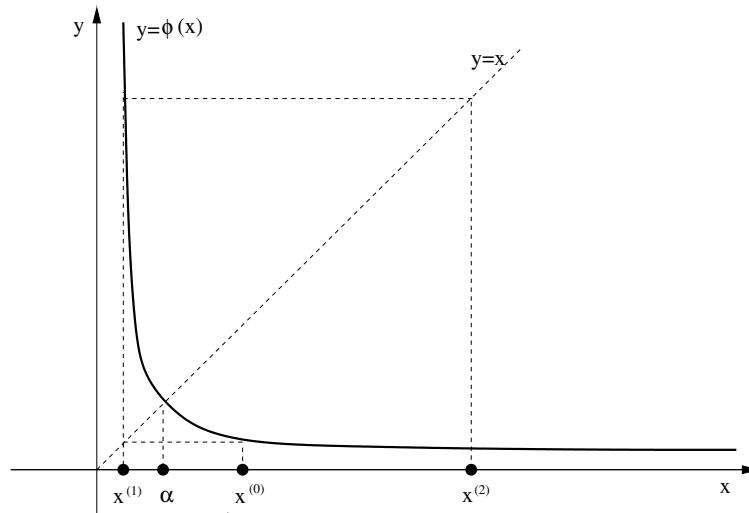
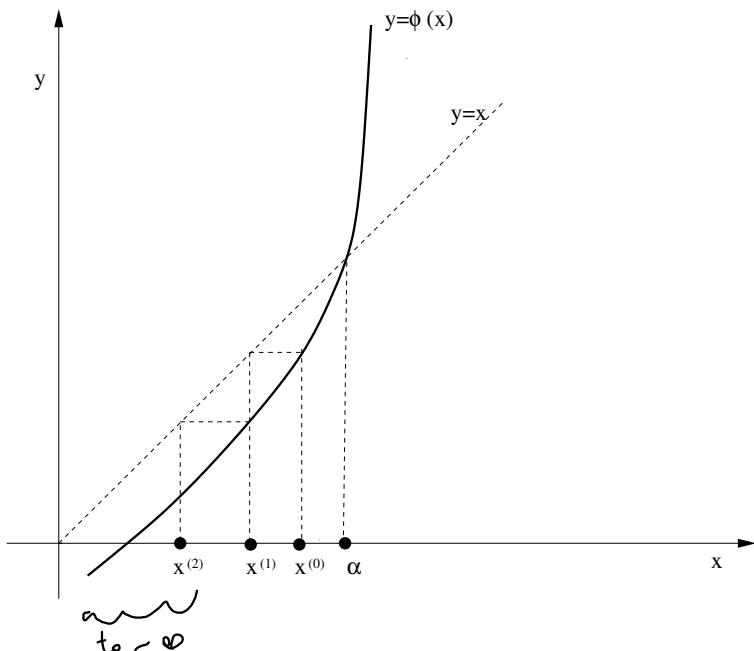
$$-1 < \phi'(\alpha) < 0.$$



## Divergent cases:

$$\phi'(\alpha) > 1,$$

$$\phi'(\alpha) < -1.$$



**Example. 5 (suite)** We apply the fixed point iterations on functions  $\phi_2(x) = rx/(1 + x/K)$  and  $\phi_3(x) = rx^2/(1 + (x/K)^2)$  that represent the Verhulst model and predator/prey model respectively, with  $K = 1.5$  and  $r = 2$ . We consider the starting point  $x^{(0)} = 1.0$ .

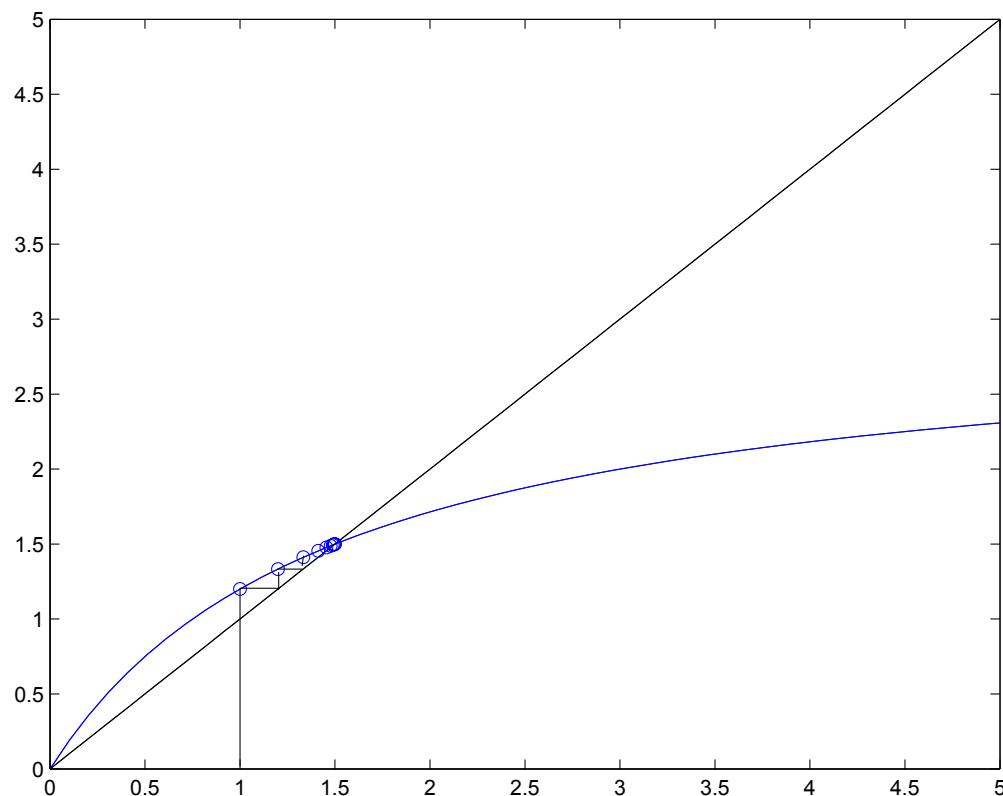
```

>> phi2=@(x) x.*(2./(1+(x./1.5)));
>> phi3=@(x) x.*(2*x./(1+(x./1.5).^2));
>> x=linspace(0,5,50);
>> figure I ght be plotted I pt t-be plotted
>> plot(x,phi2(x), 'b', x, x, 'k');
>> [p2,res2,niter2]=fixpoint(phi2,1,1e-6,1000);
>> figure
>> plot(x,phi3(x), 'b', x, x, 'k'); A for which to tollerance find fixed pt max # steps
>> [p3,res3,niter3]=fixpoint(phi3,1,1e-6,1000);

```

We find the stationary points  $\alpha_2 = 1.5$  and  $\alpha_3 = 3.9271$ .

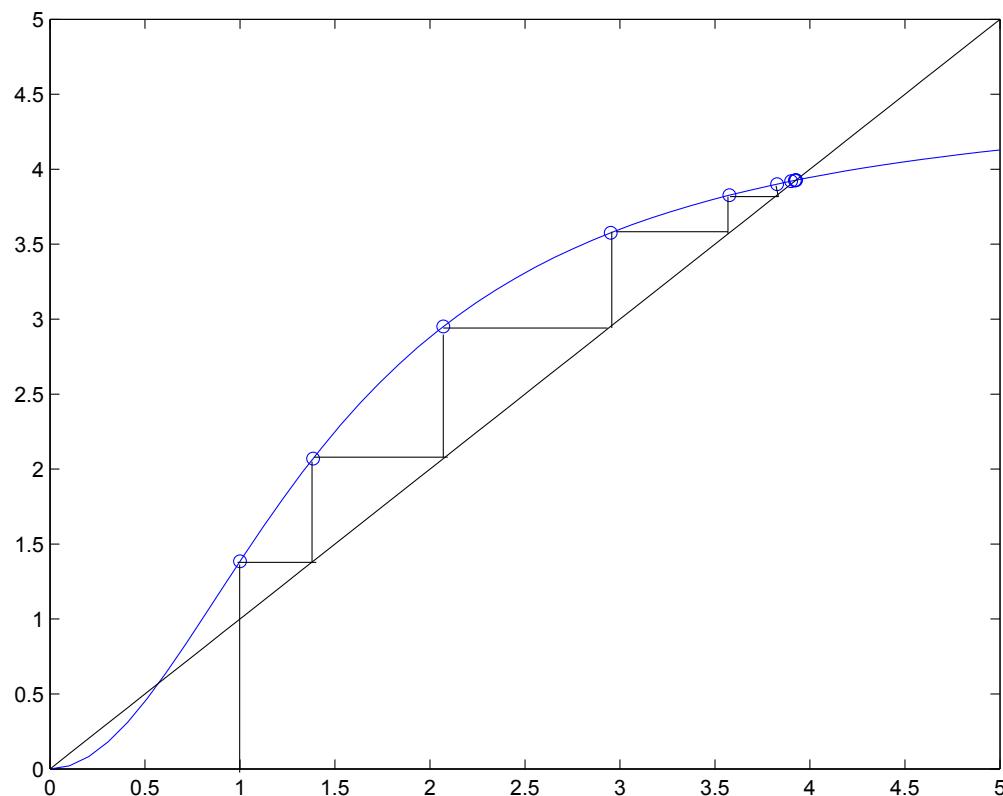
Function  $\phi_2(x)$ :



Careful to  
what first index  
is (0 or 1 !)

$x^{(0)} = 1.0000,$ $x^{(1)} = 1.2000,$ $x^{(2)} = 1.3333,$ $x^{(3)} = 1.4118,$	$ x^{(0)} - \alpha_2  = 0.5000$ $ x^{(1)} - \alpha_2  = 0.3000$ $ x^{(2)} - \alpha_2  = 0.1667$ $ x^{(3)} - \alpha_2  = 0.0882$
--	--

Function  $\phi_3(x)$ :



$$\begin{aligned} x^{(0)} &= 1.0000, \\ x^{(1)} &= 1.3846, \\ x^{(2)} &= 2.0703, \\ x^{(3)} &= 2.9509, \end{aligned}$$

$$\begin{aligned} |x^{(0)} - \alpha_3| &= 2.9271 \\ |x^{(1)} - \alpha_3| &= 2.5424 \\ |x^{(2)} - \alpha_3| &= 1.8568 \\ |x^{(3)} - \alpha_3| &= 0.9761 \end{aligned}$$

**Example. 6 (cont)** We have used the fixed point algorithms using the two functions  $\phi_1$  and  $\phi_2$  with initial value  $x^{(0)} = 0.7$ . Remember that both have the same fixed point  $\alpha$ .

easier, but  $\rightarrow x = \phi_1(x) = 1 - \sin(2x)$  because  $|\phi'_1(x)| > 1$   
 not working!

more difficult  $\rightarrow x = \phi_2(x) = \frac{1}{2} \arcsin(1 - x), \quad 0 \leq x \leq 1$   
 but works

! FIRST THING  
 TO CHECK  
 is I deriv!

```

>> [p1,res1,niter1]=fixpoint(phi1,0.7,1e-8,1000);
>> [p2,res2,niter2]=fixpoint(phi2,0.7,1e-8,1000);
  
```

The fixed point algorithm with the first function does not converge, while with the second one it converges to  $\alpha = 0.352288459558650$  in 44 iterations. Indeed,  $\phi'_1(\alpha) = -1.5237713$  and  $\phi'_2(\alpha) = -0.65626645$ .

# More about the Newton method.

↑ Take the best of Newton method and extend it in the whole fixed point framework

Very good convergence property (QUADRATIC!)

The Newton method is a fixed point method:  $x^{(k+1)} = \phi(x^{(k)})$  for the function

$$\boxed{\phi(x) = x - \frac{f(x)}{f'(x)}}.$$

Let  $\alpha$  be a zero of  $f$ , i.e. such that  $f(\alpha) = 0$ . Note that  $\underbrace{\phi'(\alpha) = 0}$ , when  $f'(\alpha) \neq 0$ . Indeed,

$\phi'(x) = 1 - \frac{[f'(x)]^2 - f(x)f''(x)}{[f'(x)]^2}.$

↓  
 if  $f'(\alpha) = 0$   
 UNLUCKY CASE

↓  
 $\phi'(\alpha) = 1 - \frac{[\overbrace{f'(\alpha)}^{\text{if } f'(\alpha)=0}]^2 - \overbrace{f(\alpha)}^{\text{if } f(\alpha)=0} \overbrace{f''(\alpha)}^{\text{if } f''(\alpha)\neq 0}}{[\overbrace{f'(\alpha)}^{\text{if } f'(\alpha)=0}]^2}$   
 $= 1 - \frac{[\overbrace{f'(\alpha)}^{\text{if } f'(\alpha)=0}]^2}{[\overbrace{f'(\alpha)}^{\text{if } f'(\alpha)=0}]^2} = 0$

so good because we obtain fixed point  
 with II procedure,  
 (it converges),  
 and it is consistent with  
 Newton method (II order)

use it one algorithm  
 ↓  
 merge fixed point  
 taking the property of  
 Newton method

↓  
 QUADRATIC CONVERGENCE  
 to the fixed point  
 which is "justified" by  
 obtaining  $\phi'(\alpha) = 0$

It is important that  $\phi''(\alpha) \neq 0$

If this fails to happen  
failure is of Newton  
method (it is going to  
be linear)



Solve by change a little  
bit the  $\phi$  ft  
(DEFLATION METHOD)

{  
modify Newton scheme  
by making sure that  
(the II derivative  $\neq 0$ )

**Theorem 1.** If  $f$  is twice differentiable,  $f(\alpha) = 0$  and  $f'(\alpha) \neq 0$ , then there exists  $\delta > 0$  such that, if  $|x^{(0)} - \alpha| \leq \delta$ , the sequence defined by the Newton method converges to  $\alpha$ .

Moreover, the convergence is quadratic; more precisely

$$\lim_{k \rightarrow \infty} \frac{x^{(k+1)} - \alpha}{(x^{(k)} - \alpha)^2} = \frac{\overbrace{f''(\alpha)}^{\text{this represents } \phi''(\alpha)}}{2\overbrace{f'(\alpha)}^{\text{the second derivative}}}$$

*Proof.* The property of convergence comes from the Proposition 2, while the quadratic convergence is a consequence of the Proposition 3, because

$$f'(\alpha) = 0 \text{ and } \frac{\phi''(\alpha)}{2} = \frac{f''(\alpha)}{2f'(\alpha)}.$$

$\Rightarrow$  We prove quadratic convergence thanks to inferring from the functions, without the second derivative of  $\phi(x) = x - \frac{f(x)}{f'(x)}$

**Definition 2.** Let  $\alpha$  be a zero of  $f$ .  $\alpha$  is said to have multiplicity  $m$ ,  $m \in \mathbb{N}$ , if  $f(\alpha) = \dots = f^{(m-1)}(\alpha) = 0$  and  $f^{(m)}(\alpha) \neq 0$ .

A zero that has multiplicity  $m = 1$  is called simple zero.

!  $\frac{x^{(k+1)} - \alpha}{x^{(k)} - \alpha}$  IS UNCOMPUTABLE, because

every ft or pty related to  $\alpha$  is  
UNCOMPUTABLE, since  $\alpha$  is NOT known

**Remark 2.** If  $f'(\alpha) = 0$ , the convergence of the Newton method is linear, not quadratic. We can use the modified Newton method:

$$x^{(k+1)} = x^{(k)} - m \frac{f(x^{(k)})}{f'(x^{(k)})}, \quad k = 0, 1, 2, \dots \quad (7)$$

where  $m$  is the multiplicity of  $\alpha$ .

at least = 2 since  $f'(\alpha) = 0$

If the multiplicity  $m$  of  $\alpha$  is unknown, there are other methods, the *adaptive methods*, which can recover the quadratic order of convergence.

Amplifying  $f$  by  $m$  or reduce  $f'$  by  $m$   
 $\Rightarrow$  IMP: you can try to avoid updating  $f'$   
 at each iteration and delay the update  
 of the most expensive computational part  
 not free!  
 you may  
 ↳ ADAPTIVITIES ↳ ACCURACY

# A stopping criterion for Newton

When to stop the Newton method? A good stopping criterion is the **control of the increment** : the iterations is completed when

$$|x^{(k+1)} - x^{(k)}| < \epsilon \quad (8)$$

where  $\epsilon$  is a fixed tolerance.

Indeed, if we denote  $e^{(k)} = \alpha - x^{(k)}$  is the error of the iteration  $k$ , we have

$$e^{(k+1)} = \alpha - x^{(k+1)} = \phi(\alpha) - \phi(x^{(k)}) \stackrel{\text{1 order Taylor expansion}}{=} \phi'(\xi^{(k)}) e^{(k)},$$

we bind error  
 at following  
 iteration

C between  $x^{(k)}$  and  $\alpha$

where  $\xi^{(k)}$  is between  $x^{(k)}$  and  $\alpha$ , and

$$x^{(k+1)} - x^{(k)} \stackrel{?}{=} \alpha - x^{(k)} - \alpha + x^{(k+1)} = e^{(k)} - e^{(k+1)} = \left(1 - \phi'(\xi^{(k)})\right) e^{(k)}. \quad (9)$$

Assuming that if  $k$  is large enough, we have  $\phi'(\xi^{(k)}) \approx \phi'(\alpha)$  and knowing that the Newton method for  $\phi'(\alpha) = 0$ , if  $\alpha$  is a simple zero, we find the estimation

$$|e^{(k)}| \approx |x^{(k+1)} - x^{(k)}|.$$

In this case, one can be goodly approx or bad approx by difference of subsequent iteration

The error that we commit when we adopt the criterion (8) is smaller than the fixed tolerance.

↓  
 NEWTON:  
 nice jfns  
 in perf. mean  
 ce (probabilic  
 convergence)  
 and nice re-  
 met for our  
 estimation

# Stopping criteria: the general case

In general, for all discussed methods, we can use two different stopping criteria: the iterations is completed when

$$|x^{(k+1)} - x^{(k)}| < \epsilon \quad (\text{control of the increment}),$$

or

$$|f(x^{(k)})| < \epsilon \quad (\text{control of the residual}),$$

where  $\epsilon$  is a fixed tolerance.

$f(x^{(n)})$  is the RESIDUAL, because  
 $f(x) = 0$ , so there we should  
 guarantee the consistency of the  
 methodology

Using fixed point iterations we obtain the following estimation:

$$e^{(k)} \approx \frac{1}{(1 - \phi'(\alpha))} (x^{(k+1)} - x^{(k)}).$$

*as a fit of  $\phi'$*

*for  $k$  large enough*

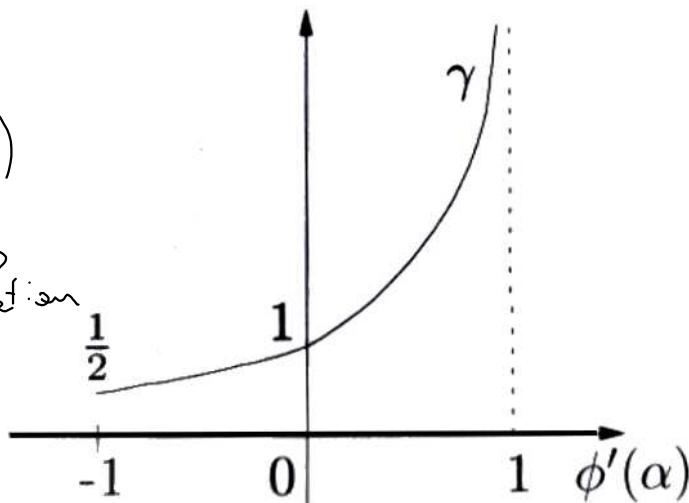
We can plot the graph of  $\frac{1}{(1 - \phi'(\alpha))}$  and comment on the relevance of the stopping criterium based on the increment:

- if  $\phi'(\alpha)$  is near to 1 the test is not satisfactory (esymptotic behavior)
- for methods of order 2 ( $\phi'(\alpha) = 0$ ), the criterium is optimal,  $\leftarrow$  Newton  $\Rightarrow$  best approximation
- if  $-1 < \phi'(\alpha) < 0$  the criterium is still all right.

*error is  $\frac{1}{2}$  of increment*

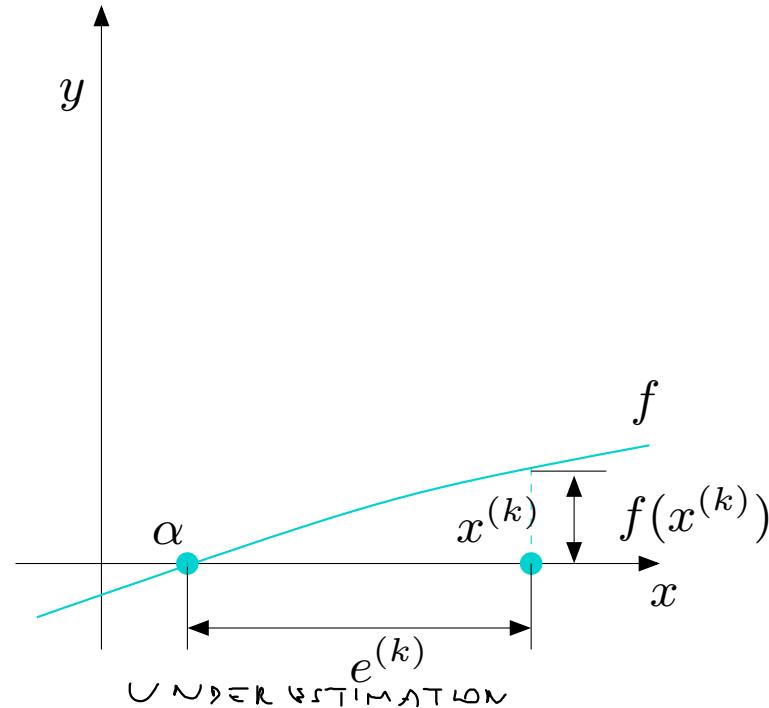
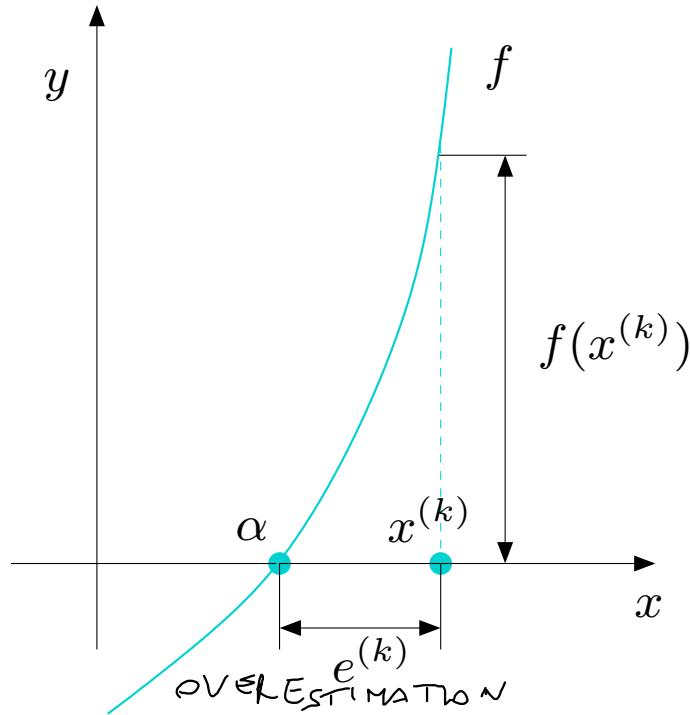
$\Rightarrow$  good, because it means we are in a

*SAFE ZONE* (good when things are worse than reality)



# Stopping Criteria

The stopping criterium based on the control on the residual  $|f(x^{(k)})| < \epsilon$  is satisfactory only if  $|f'| \simeq 1$  near the root  $\alpha$ . Otherwise it is too strong (if  $|f'| \gg 1$ ) or too weak (if  $|f'| \ll 1$ ):



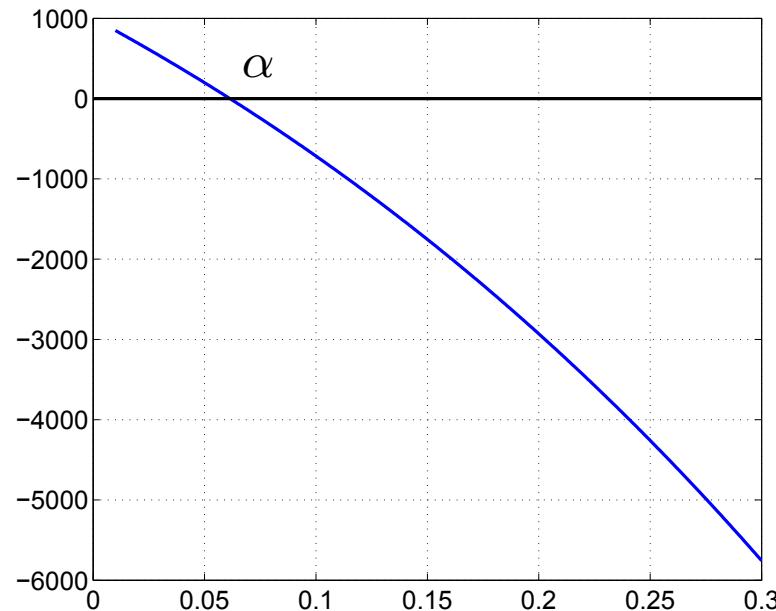
Two cases where the residual is a bad estimator of the error:  $|f'(x)| \gg 1$  (left),  $|f'(x)| \ll 1$  (right) with  $x$  near to  $\alpha$

my DEA... if divide  
 $f(x^{(n)})$  by  $f'(x^{(n)})$   
 non-iteration

# Applications

**Example. I (suite)** We draw the graph of  $f(I) = M - v \frac{1+I}{I} [(1+I)^n - 1]$  on the interval  $[0.01:0.3]$  with  $M = 6000$ ,  $v = 1000$  and  $n = 5$ :

```
>> f=@(x) 6000-1000*(1+x).*((1+x).^5 - 1)./x;
>> I = [0.01:0.001:0.3];
>> grid on;plot(I,feval(f,x));
```



The root of  $f$  is between 0.05 and 0.1.

We can apply the bisection method on the interval  $[0.05, 0.1]$  with a tolerance  $10^{-5}$

```
>> [zero,res,niter]=bisection(f,0.05,0.1,1e-5,1000);
```

The approximate solution after 12 iterations is  $\bar{x} = 0.061407470703125$ .

We can apply the Newton method with initial guess  $x^{(0)} = 0.05$

```
>> df=@(x) 1000*((1+x).^5.* (1-5*x) - 1)./(x.^2);
>> [zero,res,niter]=newton(f,df,.05,1e-5,1000);
```

The result is approximately the same, but we need only 3 iterations

The interest rate is 6.14%.

*! If you use Newton Method : You have to give also the I derivative of the nonlinear equation*

**Example. 2 (cont)** We would like to plot the angle  $\theta$  as function of  $\omega$  for  $0 \leq \omega \leq \pi$  with  $a_1 = 10\text{ cm}$ ,  $a_2 = 13\text{ cm}$ ,  $a_3 = 8\text{ cm}$ ,  $a_4 = 10\text{ cm}$ .  
 For each  $\omega$ , we have to solve the nonlinear problem

$$f(\theta) = \frac{a_1}{a_2} \cos(\omega) - \frac{a_1}{a_4} \cos(\theta) - \cos(\omega - \theta) + \frac{a_1^2 + a_2^2 - a_3^2 + a_4^2}{2a_2a_4} = 0. \quad (10)$$

To start with, we plot the graph of  $f(\theta)$  for  $\omega = \pi/3$ :

```
>> F = @(x, a1, a2, a3, a4, omega) ...  

    (a1/a2)*cos(omega) - (a1/a4)*cos(x) - cos(omega-x) ...  

    + ( (a1.^2 + a2.^2 - a3.^2 + a4.^2) / (2*a2*a4) );  

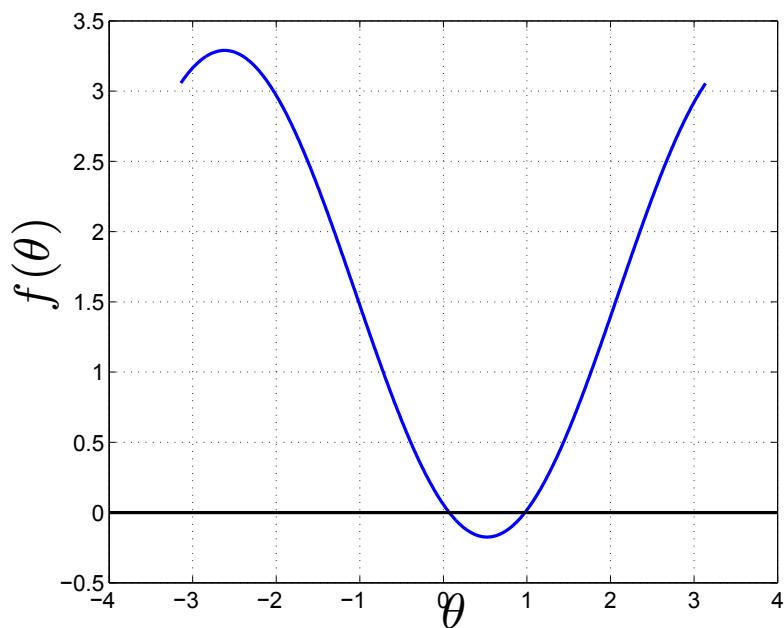
>> a1=10; a2=13; a3=8; a4=10; omega=pi/3;  

>> f = @(x) F(x,a1,a2,a3,a4,omega);  

>> x = [-pi:0.01:pi];  

>> plot( x, f(x) ); grid on;
```

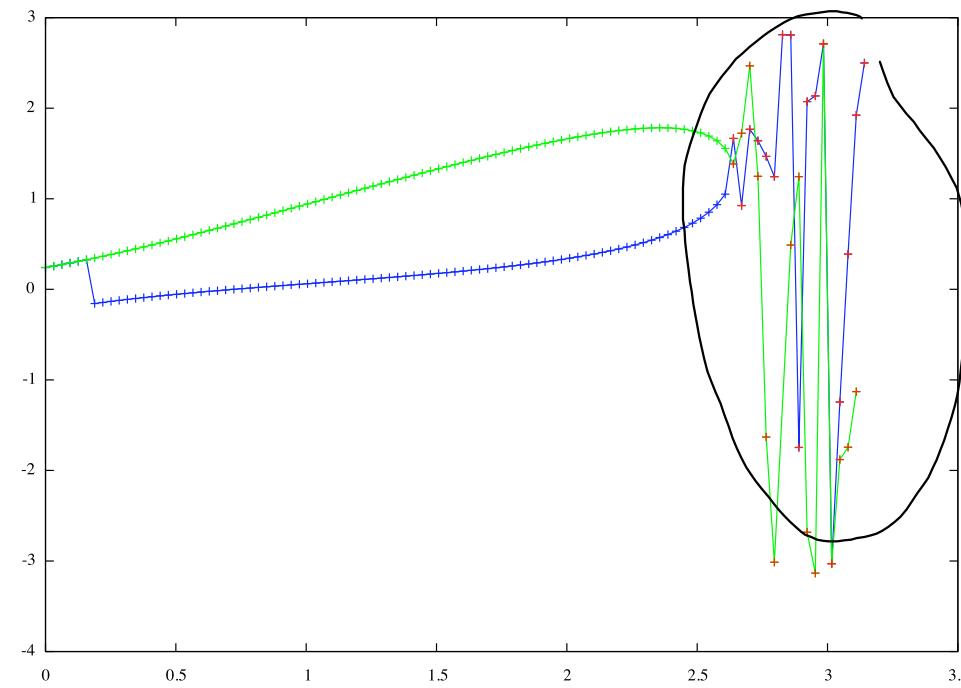
Function  $f$  with  $\omega = \pi/3$ .



We have two roots, which means that we have two possible configurations.

Now we chose 101 different values for  $\omega$ ,  $\omega_k = k \frac{\pi}{100}$ ,  $k = 0, \dots, 100$  and for each of them solve the nonlinear problem (10) with the Newton method. Since we know that we may have two distinct solutions, we give two different initial guesses to our Newton method. For example  $\theta_{01} = -0.1$  and  $\theta_{02} = 2/3\pi$

The following figure shows the solutions as function of  $\omega$ . For  $\omega > 2.6358$ , the Newton algorithm does not converge anymore. In fact, with these values there is no configuration possible



no solution,  
but INCONSISTENT  
BEHAVIOR OF  
MECHANICAL  
SYSTEM  
oscillation on  
the to numerical  
wise : NOT PHYSICAL

Here are the Matlab/Octave commands that we have used:

```
>> n=101; x01=-0.1; x02=2*pi/3; nmax=100;
>> dF = @(x,a1,a2,a3,a4,w) a1/a4*sin(x)-sin(w-x);
>> for k=1:1:n
    omega(k) = (k-1)*pi/100;
    f = @(x) F(x,a1,a2,a3,a4,omega(k));
    df = @(x) dF(x,a1,a2,a3,a4,omega(k));
    [theta1(k),res,niter] = newton(f,df,x01,1e-5,nmax);
    [theta2(k),res,niter] = newton(f,df,x02,1e-5,nmax);
end
>> plot(omega,theta1,'b:',omega,theta2,'g-')
```

**Example. 3(suite)** We consider the carbon dioxide ( $\text{CO}_2$ ), for which  $a = 0.401 \text{ Pa m}^6$  and  $b = 42.7 \cdot 10^{-6} \text{ m}^3$ .

We search the volume occupied by  $N = 1000$  molecules of  $\text{CO}_2$  in temperature  $T = 300 \text{ K}$  and pressure  $p = 3.5 \cdot 10^7 \text{ Pa}$ . We know that the Boltzmann constant is  $k = 1.3806503 \cdot 10^{-23} \text{ Joule K}^{-1}$ .



We draw the graph of the function

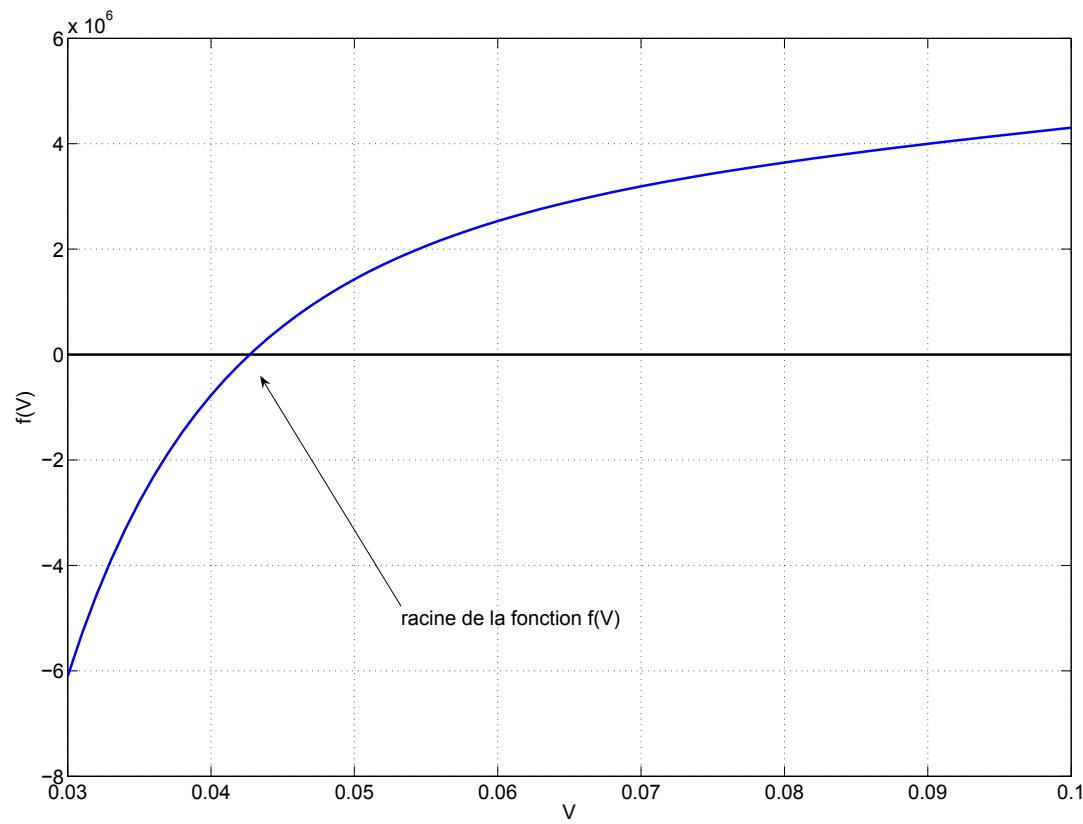
$$f(V) = \left[ p + a \left( \frac{N}{V} \right)^2 \right] (V - Nb) - kNT$$

for  $V > 0$ . We do not consider  $V < 0$  (it does not have physical meaning), because  $V$  is the volume of gas.

We use the commands in Matlab/Octave:

```
>> a=0.401; b=42.7e-6; p=3.5e7; T=300; N=1000; k=1.3806503e-23;
>> f = @(x,p,T,a,b,N,k)(p+a*((N./x).^2)).*(x-N*b)-k*N*T;
>> x=[0.03:0.001:0.1];
>> plot(x,f(x,p,T,a,b,N,k))
>> grid on
```

We obtain the graph of the function  $f(V)$ :



We see that there is a zero for  $0.03 < V < 0.1$ . If we apply the bisection method on the interval  $[0.03, 0.1]$  with a tolerance  $10^{-12}$ :

```
[zero,res,niter]=bisection(f,0.03,0.1,1e-12,1000,p,T,a,b,N,k);
```

then we find, after 36 iterations, the value  $V = 0.0427$ .

If we use the Newton method with the same tolerance, starting from the initial point  $x^{(0)} = 0.03$ ,

```
>> df = @(x,p,T,a,b,N,k) -2*a*N^2/(x^3)*(x-N*b)+(p+a*((N./x).^2));
>> [zero,res,niter]=newton(f,df,0.03,1e-12,1000,p,T,a,b,N,k);
```

then we find the same solution after 6 iterations.

The conclusion is that the volume  $V$  occupied by the gas is  $0.0427 \text{ m}^3$ .

# The rope method

→ way to avoid I derivative  
could be a Jacobian  
 IP in N-dimension  
 $N > 1$

This method is obtained by replacing  $f'(x^{(k)})$  by a fixed  $q$  in the Newton method:

$$x^{(k+1)} = x^{(k)} - \frac{1}{q} f(x^{(k)}), \quad k = 0, 1, 2, \dots \quad (11)$$

We can take, for example,  $q = f'(x^{(0)})$  or  $q = \frac{f(b) - f(a)}{b - a}$ , in the case when we search a zero in the interval  $[a, b]$ .

**Example. 6 (suite)** We apply the rope method and the Newton method to find the zero of  $f$ .

The rope method in the interval  $[-1, 1]$ , with  $x^{(0)} = 0.7$  :

```
>> [zero,res,niter]=chord(f,-1,1,0.7,1e-8,1000)
```

We find the result after 15 iterations.

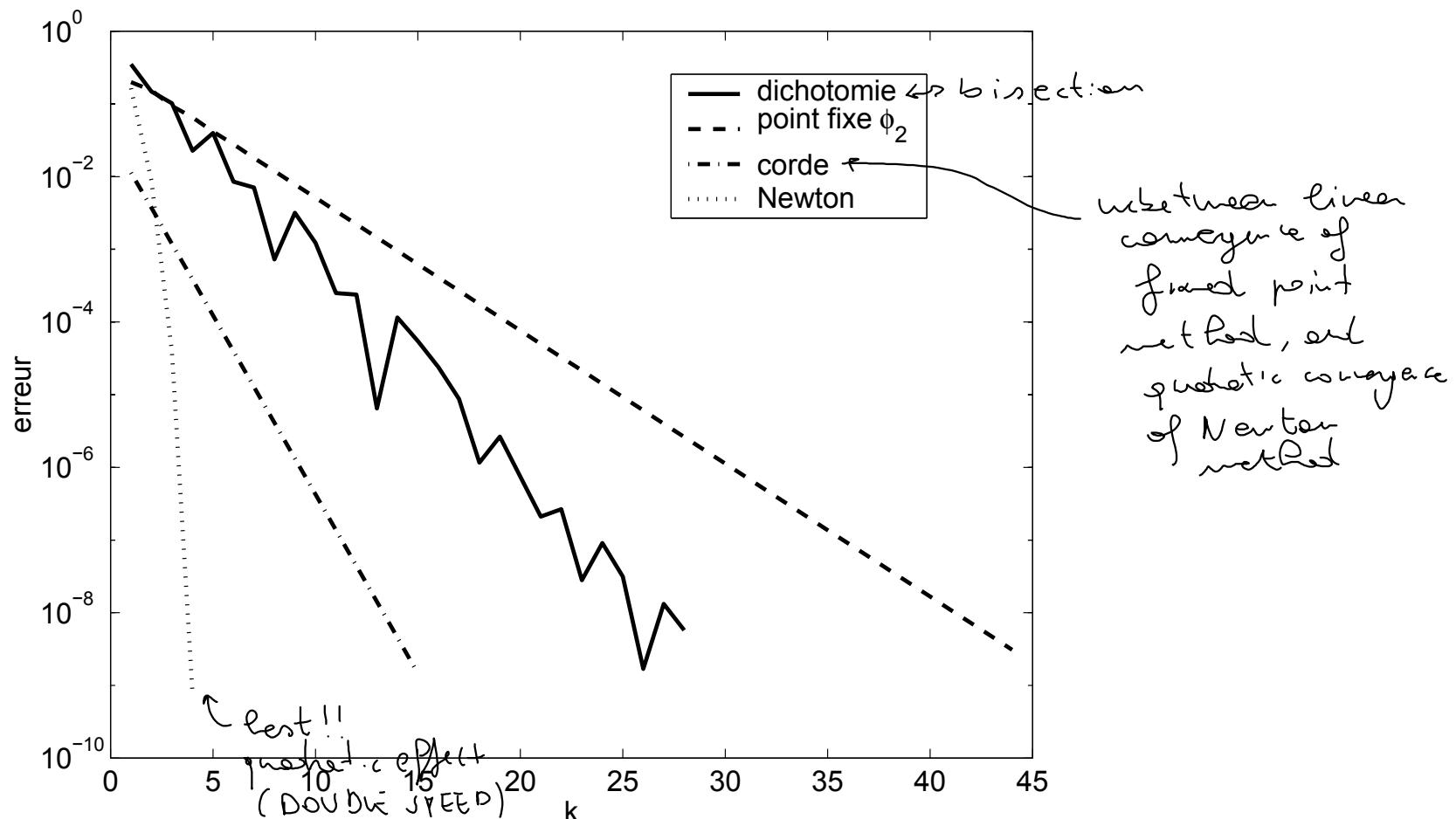
The Newton method with the same  $x^{(0)}$  :

```
>> df = @(x) 2*cos(2*x) + 1;
>> [zero,res,niter]=newton(f,df,0.7,1e-8,1000);
```

We find the result after 5 iterations. Much faster.

Exercice  
question

Values of the errors plotted versus the number of iterations for 4 methods: bisection, fixed point  $\phi_2$ , rope and Newton. There is logarithmic scale on the axis  $y$ .



**Remark 3.** *The rope method is also a fixed point method for*

$$\phi(x) = x - \frac{1}{q}f(x).$$

*So, we have  $\phi'(x) = 1 - \frac{1}{q}f'(x)$  and thanks to the Proposition 2, we obtain that the method converges if the following condition is satisfied:*

$$| 1 - \frac{1}{q}f'(\alpha) | < 1 .$$

we want  $n$  to be ~~at~~  
~~of interpolation~~  $\Rightarrow$  ~~at~~  $n+1$   
~~at less in P^n~~

Recap on interpolation : 1) a set of interpolation points  $\{x_i\}_{i=0}^n$   
 2) a set of basis functions  $\{v_i\}_{i=0}^n$

Interpolation of  $f \in C^0([a, b])$   $x_i \in [a, b] \forall i$

on the space  $V = \text{span}\{v_i\}_{i=0}^n$  is defined as the solution to :

find  $p$  in  $V$  s.t.  $p(x_i) = f(x_i) \forall i$

Given  $\{x_i\}_{i=0}^n$  best choice for  $\{v_i\}_{i=0}^n$  is Lagrange basis

$$\Rightarrow v_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^{n} \frac{(x - x_j)}{(x_i - x_j)}$$

So the interpolation poly becomes

$$p = \sum_{i=0}^n v_i(x) f(x_i) := L^n f$$

Next step was:

Try to find the "best" set of  $\{x_i\}_{i=0}^n$

Remember that  $L^n$  is a smooth poly.  
 Every smooth poly  $\underbrace{\text{last approx}}$

$$\|L^n f - f\|_\infty = \|L^n(f - p) + p - f\|_\infty \leq (\|L^n\|_\infty + 1) \|p - f\|_\infty$$

$$\|L^n\|_\infty \leq \|L^n(\{x_i\}_{i=0}^n)\|_\infty$$

We try to min the ~~\*~~ of the absolute values of the basis ft ( $\Leftrightarrow$  Lebesgue

$$\lambda := \sum_{i=0}^n |v_i|$$

$\hookrightarrow$  Lebesgue function

For Chebyshev points:

$$\frac{2}{\pi} \log(n+1) - c \leq \|\lambda\|_\infty \leq \frac{2}{\pi} \log(n+1) + 1$$

! True only if the function is continuous

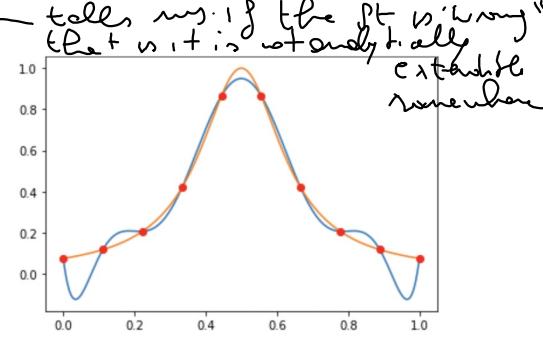
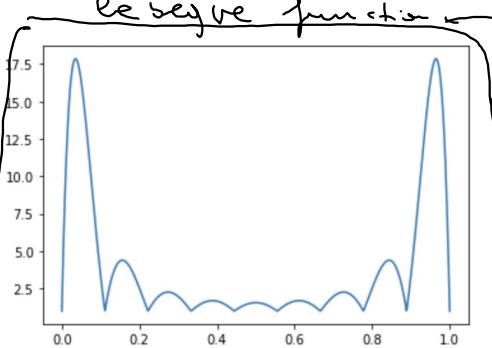
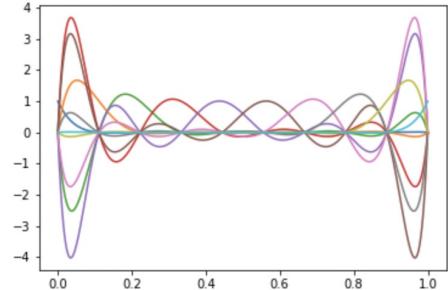
If you add more regularity, you get

Better things (Result for analytically  
extendible f) with continuous derivatives  
up to (the n-th derivative)

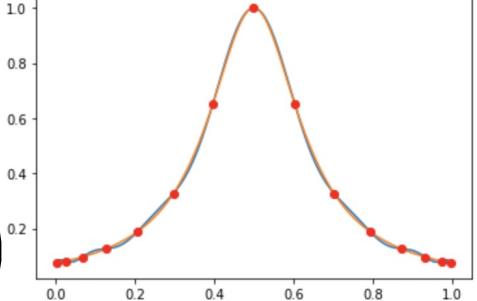
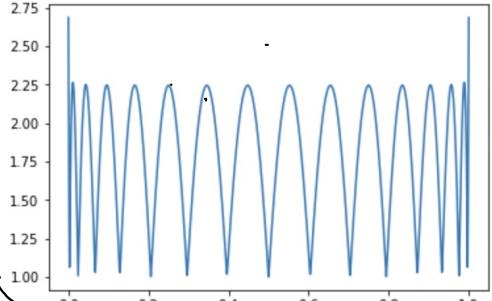
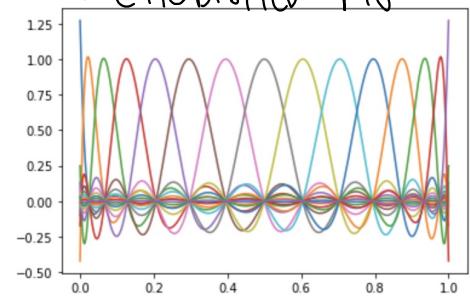
! INTERPOL IS EXACT AT POINT OF  
INTERPOL

=> try to get exact approx every where  
in some way

USING EQUISPACED PTS



USING CHEBYSHEV PTS



If interpolation is not good, what can we do in  $C^0([a,b])$ ?

Ans: Weierstrass Approximation Theorem

$\forall f \in C^0([a,b])$ ,  $\forall \epsilon > 0$ ,  $\exists p \in P^n$ , s.t.  $\|f-p\|_{\infty} < \epsilon$

~~No matter what f, you can always find a p such that p approx f~~

Proof: Let  $B_n$  be a sequence of linear operators s.t.

- 1)  $B_n$  positive linear operators  $C([a,b]) \rightarrow P^n([a,b])$
- 2)  $B_n q \rightarrow q$  pointwise  $\forall q \in P^2([a,b])$  that is only for quadrat. polys

Then:

$B_n f$  converges uniformly to  $f$   $\forall f \in C^0([a,b])$

N.B.:  $B_n$  is "positive" if  $q(x) \geq 0 \Rightarrow (B_n q)(x) \geq 0$  in  $[a,b]$   
(if  $q(x) \leq 0$  it remains negative if  $B$  is positive)

Part 1: .  $\forall f \in C^0([a,b])$ ,  $\forall x_0 \in (a,b)$ :

- . construct  $q^\pm \in P^2$  s.t.  $q^-(x_0) < f(x_0) < q^+(x_0)$
- . use  $B_n q^\pm \rightarrow q^\pm$
- . use  $B_n(f - q^-)(x_0) > 0$        $B_n(q^+ - f)(x_0) > 0$

$$f \in C^0([a, b]) \Rightarrow \forall \varepsilon > 0, \exists \delta > 0$$

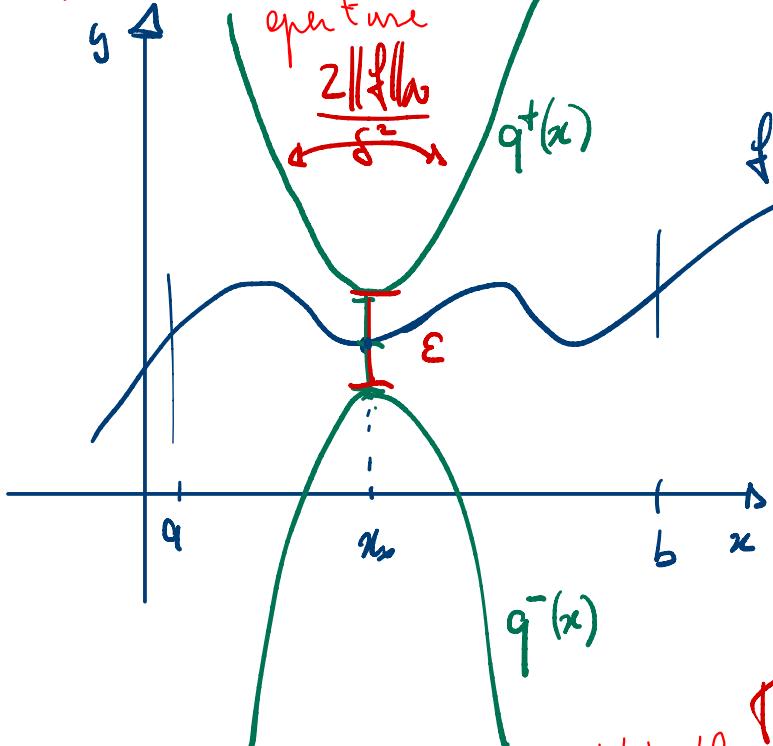
} def of continuity!

$$|x_1 - x_2| \leq \delta \rightarrow |f(x_1) - f(x_2)| \leq \varepsilon$$

Define

$$q^\pm := f(x_0) \pm \left( \frac{\varepsilon}{2} + \frac{2\|f\|_\infty}{\delta^2} (x - x_0)^2 \right)$$

Busting up in this way guarantees that they are always exact below &



for varying  $x_0$

$$q^\pm(x) = a^\pm x + b^\pm x + c^\pm$$

$a^\pm, b^\pm, c^\pm$  depend on

$x_0, \|f\|_\infty, \delta, \varepsilon$

Take  $M = \max_{x_0 \in [a, b]} (|a^\pm|, |b^\pm|, |c^\pm|)$

I take the max over the entire interval

$M$  depends on  $\|f\|_\infty, \delta, \varepsilon$ , but Not on  $x_0$

Choose  $N$  large enough s.t.

( $\exists N$  s.t. this is possible)  
Since we know that  $x_n$  converges uniformly for all of the two (HYPOTHESIS)

$$\|B_n x_n^i - x_i^i\|_\infty \leq \frac{\varepsilon}{6M}$$

$\forall n \geq N$   
 $\forall i = 0, 1, 2$

$$\Rightarrow \forall x_0 \quad * \|B_n q^\pm - q^\pm\| \leq \frac{\varepsilon}{2}$$

$\forall n \geq N$   
 $M$  controls every coeff  
 $\Rightarrow$  this allows to say that  $\frac{\varepsilon}{6M} \leq \frac{\varepsilon}{2}$  (?)

in  $x_0$   $f(x_0) - \varepsilon \leq q^-(x_0) - \frac{\varepsilon}{2} \leq B_n q^-(x) \leq B_n f(x_0)$

And repeat for the reversed inequality:

$$B_n f$$

positivity

$$B_n f \leq B_n q^+ \leq q^+ + \frac{\varepsilon}{2} \leq f(x) + \varepsilon$$

$$\|B_n f - f\| \leq \varepsilon$$

How to define  $B_n$ ? Let  $(a, b) = [0, 1]$

Trick : Consider  $\underbrace{\left( (1-x) + x \right)^n}_{\stackrel{=1}{\text{defines bin. functions}}} = \sum_{i=0}^n \binom{n}{i} x^i (1-x)^{n-i}$  poly's of  $n$

$$= \sum_{i=0}^n v_i(x)$$

$v_i(x) = \binom{n}{i} x^i (1-x)^{n-i}$

Bernsteine polynomials.

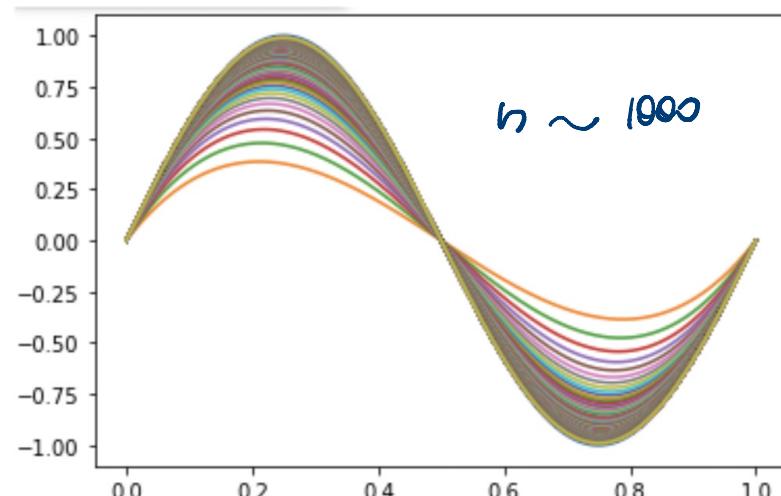
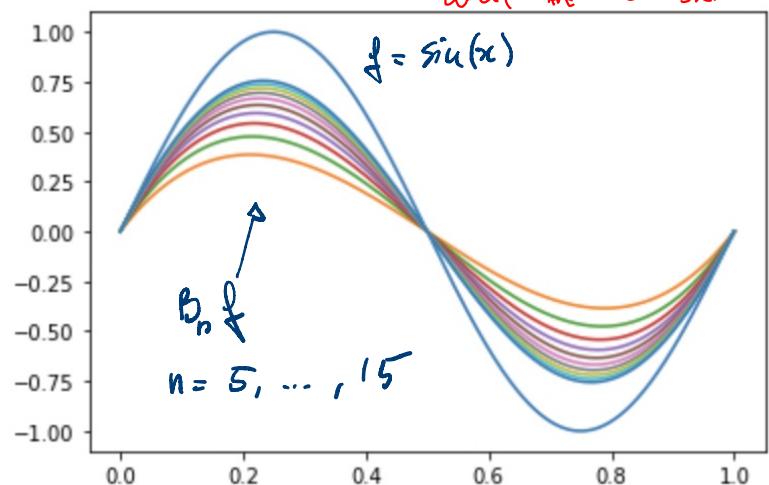
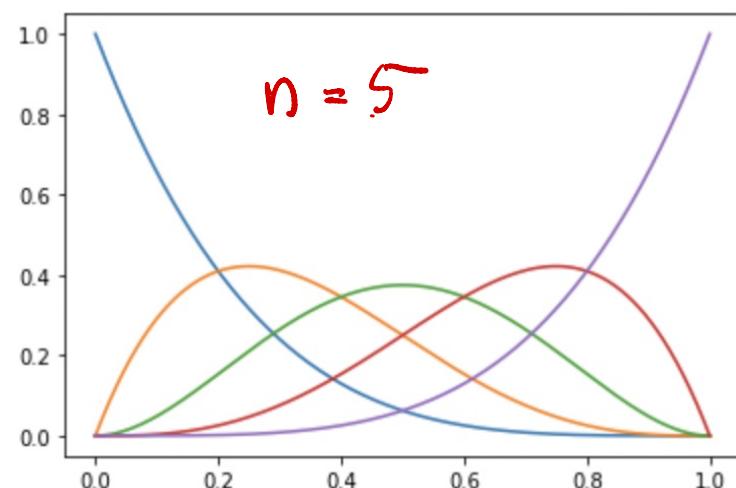
$B_n f := \sum_{i=0}^n v_i(x) f\left(\frac{i}{n}\right)$

$\rightarrow$  but always positive in  $[0, 1]$   
 or also from 0 in middle  
 because  $v_i$  positive  
 difference between  
 interpolation and thin  
 $f \geq 0 \Rightarrow B_n f \geq 0$   
 because  $v_i$  positive  
 you cannot  
 say that  
 the  $v_i$  are  
 interpolating  
 at any pt but  
 0 and 1

$B_n x = x \quad B_n 1 = 1$

$B_n x^2 = \left(\frac{n-1}{n}\right)x^2 + \frac{1}{n}x \neq x^2$

for large  $n$   
 but  $n \rightarrow \infty \quad B_n x^2 \rightarrow x^2$



We don't  
 have a quick  
 convergence

## Lecture 2 PROJECTION BASE APPROX METHOD 10. Nov. 2020

1) Approximate and measure using  $L^2$  norm (vs  $L^\infty$ ) derived from a scalar product

→ Scalar product :  $u, v \in L^2([0,1])$ ,

$$(u, v) := \int_0^1 u \cdot v \, dx \quad \|a+b\|^2 = (a+b, a+b) = \|a\|^2 + \|b\|^2 + 2(a, b)$$

$$\leadsto \|u\|_{L^2([0,1])}^2 \equiv \|u\|^2 := (u, u) = \int_0^1 |u|^2 \, dx$$

let's assume we have  $V := \text{span}\{v_i\}_{i=0}^n$ , and  $\|\cdot\|$ .

The best approximation of  $u \in L^2([0,1])$  in  $V$  satisfies:

$$p \in V \text{ s.t. } (u-p, v) = 0 \quad \forall v \in V$$

$\curvearrowleft (u-p, v) = 0 \quad \forall v \in V \iff p \text{ is B.A. of } u \text{ in } V$

Proof: Assume  $p$  is best approx  $\Rightarrow (u-p, v) = 0 \quad \forall v \in V$   
by def of B.A.

$$\|u-p\|^2 \leq \|u-p + tq\|^2 \quad \forall t > 0, \forall q \in V$$

$$\|u-p + tq\|^2 - \|u-p\|^2 \geq 0 \quad \begin{matrix} \text{if } v \text{ in any direction from} \\ \text{get larger distance from } p \text{ to } u \end{matrix}$$

$$(a+b)^2 - (a-b)^2 = 4(a, b)$$

$$\left\| \underbrace{(u-p + \frac{t}{2}q)}_{a+b} + \left( \frac{t}{2}q \right) \right\|^2 - \|u-p\|^2 \geq 0$$

$$a = u-p + \frac{t}{2}q$$

$$b = \frac{t}{2}q$$

$$2 \not\perp (u-p + \frac{t}{2}q, \frac{t}{2}q) \geq 0$$

Linearity of scalar prod in 1 argument

$$0 \leq 2(u-p, q)t + t^2(q, q)$$

We have a symmetry for:

$$q \rightarrow -q$$

$$-\frac{t}{2}\|q\|^2 \leq (u-p, q) \leq \frac{t}{2}\|q\|^2$$

$$\forall t > 0$$

$$\Rightarrow (u-p, q) = 0$$

$$\forall q \in V$$

Assume  $\star$   $(u - p, v) = 0 \quad \forall v \in V \Rightarrow p$  is B.A.

$$\|u - v\|^2 = \|u - p + p - v\|^2 = \|u - p\|^2 + \|p - v\|^2 + 2(u - p, p - v)$$

$$\Rightarrow \|u - p\|^2 \leq \|u - v\|^2 \quad \forall v \in V$$

that is  $p \in$  B.A. of  $\frac{u-p}{\|u-p\|}$

by \*  
and the fact  
that  $u-p \perp p$

B.A. <sup>(weak)</sup> Find  $p \in V$  s.t.

$$\text{1) } (p, v) = (u, v) \quad \forall v \in V$$

given  
Write  $p(x) = \sum p^j v_j(x)$

$$\Rightarrow 2) \quad (p^j v_j, v_i) = (u, v_i)$$

MASS  
MATRIX

$$\underline{\underline{M}} \quad \underline{p} = \underline{U}$$

$$M_{ij} \quad p^j = U_i$$

$$M_{ij} = (v_i, v_j)$$

$$U_i = (u, v_i)$$

computing  $\Rightarrow$  the prob  
 $p^j = \underline{\underline{M}}^{-1} \underline{U}_i$

Example :

$$V = \text{span} \left\{ \text{pow}(x, i) \right\}_{i=0}^n \quad p = \sum_j p^j x^{(j)} \quad \text{power basis}$$

then  $M_{ij} = \int_0^1 x^i \cdot x^j \, dx = \int_0^1 x^{(i+j)} \, dx = \frac{1}{i+j+1}$

Here <sup>"called"</sup> Hilbert Matrix : cond. ( $H$ )  $\approx \frac{(1+\sqrt{2})^{4n}}{\sqrt{n}}$

it is Tenu-  
ely condi-  
tioned  
matrix  
If  $i=j$  very  
large  $\Rightarrow$  column  
should be same  
 $(i+j+1) \times i = 1 + j + 1 \dots$

What's the best possible "M" matrix?

Identity :  $M_{ij} = \delta_{ij} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$

Def B.A. in a space with norm defined  
for scalar product, select basis functions

and try to compute mass matrix



B.A. in a finite dim space (spanned by  
some basis  $\{f_i\}$ ) what you obtain is

$$\underline{M} \underline{p} = \underline{U}$$

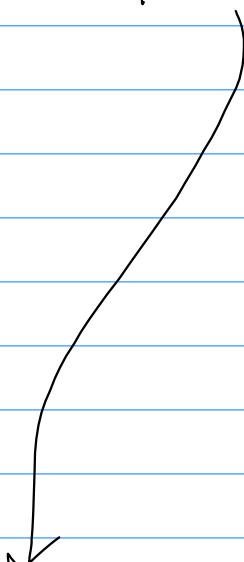
$\Rightarrow$  Best pt is the one in which you  
don't have to write any matrix  $\Rightarrow$  possible if  $M=11$

$$\Rightarrow p_i = (u, v_i)$$



Simplest procedure

GRAM SCHMIDT PROCEDURE



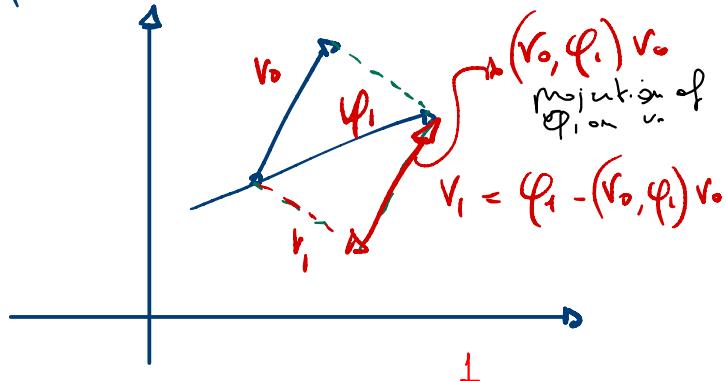
Let's try to find an orthonormal basis

Set  $v_0 = 1$   $\rightarrow \left( \int_0^1 (1)^2 dx \right)^{\frac{1}{2}} = 1$

$$q_{i+1} = x v_i - \sum_{j=0}^{i-1} (x v_i, v_j) v_j$$

and normalize

$$v_{i+1} = \frac{q_{i+1}}{\|q_{i+1}\|}$$



Possible problems:

$$\begin{aligned} \|q_i\|^2 &= \|x v_0\|^2 + (x v_0, v_0)^2 \|v_0\|^2 - 2(x v_0, v_0)^2 \\ &= \|x v_0\|^2 - (x v_0, v_0)^2 \end{aligned}$$

This can be very small!!

We have another solution which does not ask same norm:  $\varphi_{ij} \neq \delta_{ij}$   
but it is diagonal:

Bonnet recursion:

$$P_0 = 1 \quad P_1 = x$$

$$(n+1) P_{n+1}(x) = (2n+1) x P_n(x) - n P_{n-1}(x)$$

Are orthogonal in  $[-1, 1]$ , and

$$P_n(1) = 1 \quad \forall n$$

dim (?)  
(start from fact that  $P_0, P_1$  are OK)

"The only" thing which is left to do is to compute

$$\int_0^1 u P_i dx$$

Approximate integrals.

Between functions you def the error as  $\| P(x) - u(x) \|$  where  $\| \cdot \|$  is  
the norm used.

## Lecture 13

12.11.2020

Recap of Lec 12: Given  $u \in L^2([0,1])$ , find  $\rho \in V = \text{span}\{v_i\}_{i=0}^n$   
 s.t.  $(\rho, v) = (u, v) \quad \forall v \in V$

$$\Rightarrow \sum_j (\rho^T v_j, v_i) = (u, v_i) \quad i = 0, \dots, n$$

Rewritable in a system

$$\underline{\underline{\rho}} = \underline{\underline{v}}$$

$$M_{ij} := (v_i, v_j)$$

$$v_i = (u, v_i)$$

$$(a, b) := \int_0^1 a \cdot b \, dx$$

Induced a norm

$$\Rightarrow \|a\|^2 = (a, a) \quad \Rightarrow \|a\| = \left( \int_a^b a^2 \, dx \right)^{\frac{1}{2}}$$

Today: Approximate

$$\int_a^b u(x) \, dx$$

For now:  $u \in C^0([a,b])$

Rough approx schemes for  $C^0([a,b])$  functions

For example:

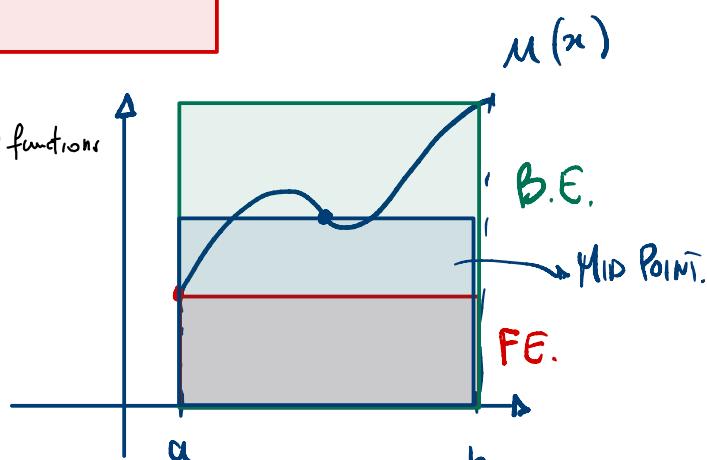
$\begin{cases} FE \\ BE \\ \text{MIDPOINT} \end{cases}$	$\begin{cases} \text{(Forward Euler)} \\ \text{(Backward Euler)} \end{cases}$
---	---

We have, for each scheme

$$FE: \int_a^b u(x) \, dx \approx u(a)(b-a)$$

$$BE: \int_a^b u(x) \, dx \approx u(b)(b-a)$$

$$\text{MIDPOINT}: \int_a^b u(x) \, dx \approx u\left(\frac{a+b}{2}\right)(b-a)$$



you don't use just one of them, but you apply all three, one for each special case (such as particular sub intervals of  $[a,b]$ )

Interpolatory Quadrature Rules:

$$I(u) := \int_a^b u(x) \, dx$$

$$I_n(u) := \sum_{i=0}^n u(q_i) w_i$$

quadrature points      weights

$$\approx \int_a^b u(x) \, dx = I(u)$$

if you choose correctly both  $q_i$  and  $w_i$

If we want to integrate  
we can choose  $\{q_i\}_{i=0}^n$  and  $\{w_i\}_{i=0}^n$  2(n+1) unknowns

Most popular choice: Given  $\{q_i\}_{i=0}^n$ , construct  
interpolatory quadrature rule  $L_u^n$  and integrate (exactly)  $L_u^n$

Lagrange interpolation of  $u$  function

$$I(u) = \int_0^1 (L_u^n)(x) dx$$

where  $L_u^n = \sum_{i=0}^n u(q_i) e_i(x)$  choose the weights  $e_i = \prod_{j=0, j \neq i}^n \frac{(x - q_j)}{(q_i - q_j)}$

$$\Rightarrow I(u) = \int_0^1 \sum_{i=0}^n u(q_i) e_i(x) dx = \sum_{i=0}^n u(q_i) \int_0^1 e_i(x) dx$$

$$= \sum_{i=0}^n u(q_i) w_i$$

number independent  
in the function  $I^n$   
trying to integrate

FE, B.E., Mid Point are all interpolatory quadrature  
rules with  $P^0$  polynomial approximation  
or that is with only 1 quadrature pt

### Trees 1

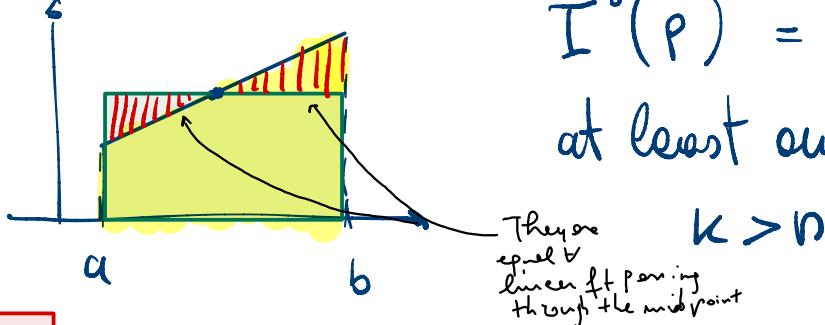
Interpolatory quadrature rules with  $n+1$  points  
are exact at least for polynomials of order  $\leq n$

Proof:  $\forall p \in P_{odd}^{n+1}$ ,  $L_p^n = p$  by construction

$$\Rightarrow I(p) = I^n(p) \quad \forall p \in P^n$$

Degree of accuracy of  $I^n$  is the maximum integer  
 $K$  s.t.  $I(p) = I^n(p) \quad \forall p \in P^K$

For the midpoint quadrature rule



$$I^0(p) = I(p) \quad \forall p \in P^1$$

at least one case where

Theo

Given  $\{q_i\}_{i=0}^n$ , degree of accuracy is  $< 2(n+1)$

Proof : construct  $w(x) = \prod_{i=0}^n (x - q_i)$ ,  $w \in P^{n+1}$

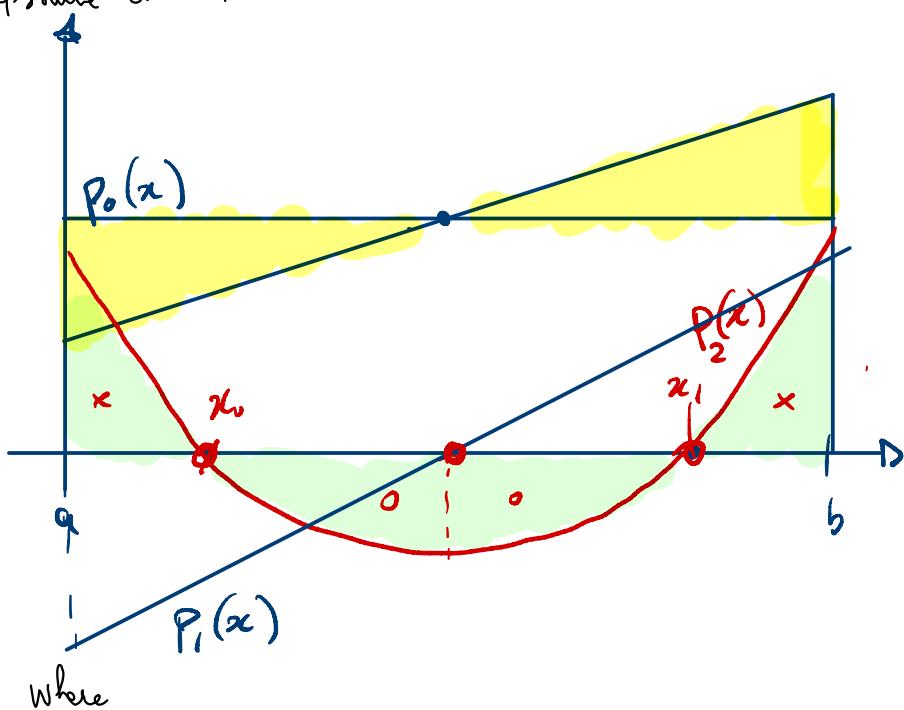
$$w^2(x) > 0 \quad \forall x \neq q_i, w(q_i) = 0 \quad i = 0, \dots, n$$

$$\Rightarrow I(w^2) > 0 \quad I^b(w^2) = 0 \quad \text{because by construction.}$$

$$\Rightarrow \exists p = w^2 \in P^{2n+2} \text{ s.t. } I(p) + I^b(p)$$

Can the degree of accuracy be  $k = 2n+1$ ? Yes

Assume current situation



Areas marked with  $x$  and with  $o$  must be equal

equivalent to saying  $\equiv$

average of  $P_2(x)$  is zero

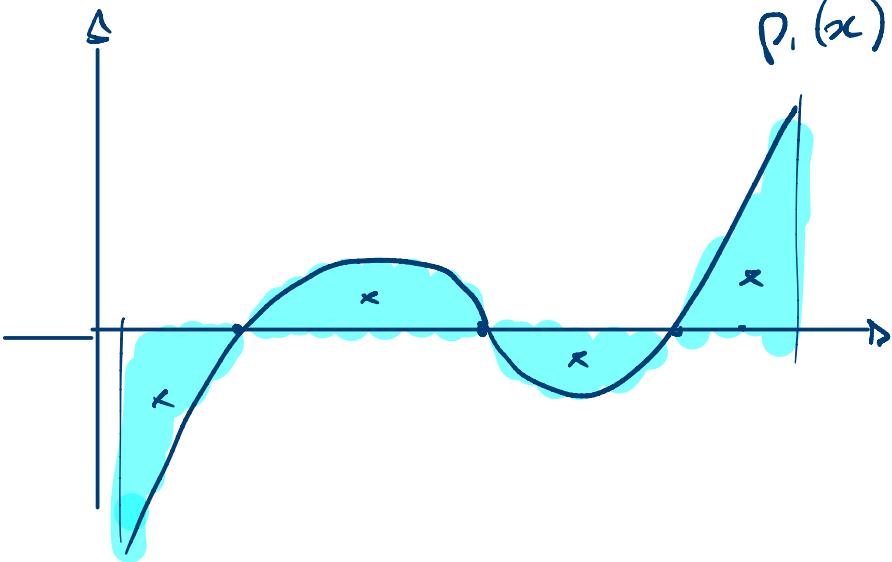
where

$$P_0(x) = 1 \quad \xrightarrow{\frac{a+b}{2}}$$

$$P_1(x) = (x - 0.5)$$

$$P_2(x) = (x - x_0)(x - x_1) \quad \text{s.t. } \int P_2(x) = 0$$

Now take the pt  
 $p_1(x) \cdot p_2(x) \in P^3(x)$



the areas marked  
with  $x$  are  
equal

Theo

Let  $u \in P^{n+m}$

$m \leq n+1$

then  $I^n(u) = I(u)$

$\Leftrightarrow$

$I^n$  has degree of accuracy  $k = n+m$

$\Leftrightarrow$

$$\int_a^b \omega(x) p = 0 \quad \forall p \in P^{m-1}$$

can be chosen as  $\perp$  to every  $p$

$$\omega = \prod_{i=0}^n (x - q_i) \in P^{n+1}$$

we are asking  $\omega$  to be Legende polynomial

Proof

Any polynomial  $p \in P^{n+m}$  can be

written

as :

$$p(x) = \underbrace{\omega(x)}_{\text{by}} \underbrace{\prod(x)}_{P^{n+1}} + q(x) \quad \underbrace{P^{m-1}}_{P_m}$$

basically  
division + rest

Ruffini's theorem

Given  $g(x) \in \mathbb{P}^{m-1} \subseteq \mathbb{P}^n$  ( $m \leq n+1$ )

Then  $I(g) = I^n(g)$

$$I(p) = I(\omega\pi) + I(g) \quad \text{by linearity of interpolation}$$

$$I^n(p) = I^n(\omega\pi) + I^n(g) \\ = I(g)$$

$$\text{If } I(\omega\pi) = 0 = \int_0^1 \omega\pi dx \quad \text{then } I(p) = I^n(p)$$

Now we choose  $\{g_i\}_{i=0}^n$  as the roots of Legendre basis

of order  $n+1$ :  $\int_a^b \omega_{n+1} l_j dx = 0 \quad j < n+1$

Legendre basis of  
order  $n+1$

The roots of  $\omega_{n+1}$  are called Gauss quadrature points. And  $I^n$  in this case is a Gauss quadrature formula, and it has degree of accuracy  $K = 2n+1$

General Strategy : find  $\{q_i\}, \{w_i\}$  s.t.

$$\sum_{j=0}^n v_i(q_j) w_j = \int_a^b v_i(x) dx \quad \text{for } \{v_i\}_{i=0}^{2n+1} \text{ known}$$

you impose that the interpolatory quadrature formula is EXACT for a fixed set of  
orthogonal functions

$2(n+1)$  unknowns, in  $2n+2$  equations.

NON LINEAR SYSTEM OF  $2n+2$  EQUATIONS.

if you can solve this, you win.

for  $v_i = x^i$  you obtain  
the GAUSS QUADRATURE  
FORMULA

# Least Squares



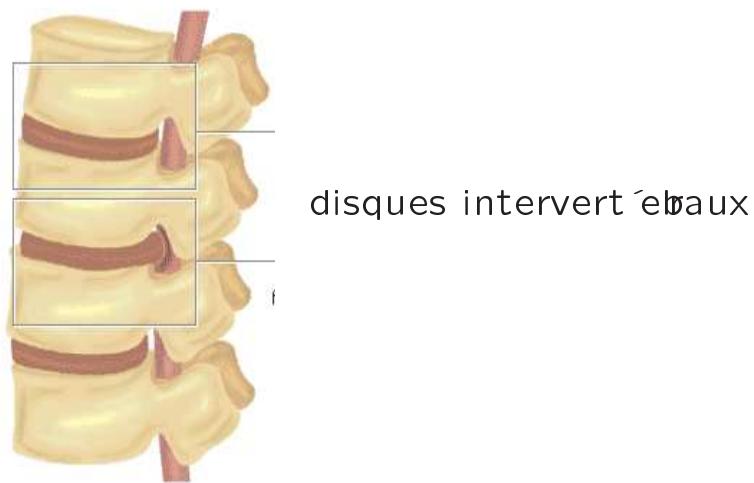
## Numerical Analysis

Profs. Gianluigi Rozza- Luca Heltai

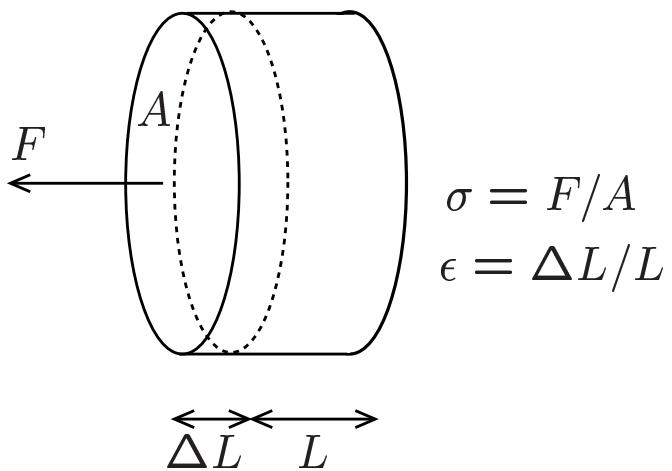
2019-SISSA mathLab Trieste

# Examples and motivations

**Example 1.** We consider a mechanical test to establish the link between stresses ( $MPa = 100N/cm^2$ ) and relative deformations ( $cm/cm$ ) of a sample of biological tissue (an intervertebral disc, taken from P. Komarek, Chapt. 2 of *Biomechanics of Clinical Aspects of Biomedicine*, 1993, J. Valenta ed., Elsevier).

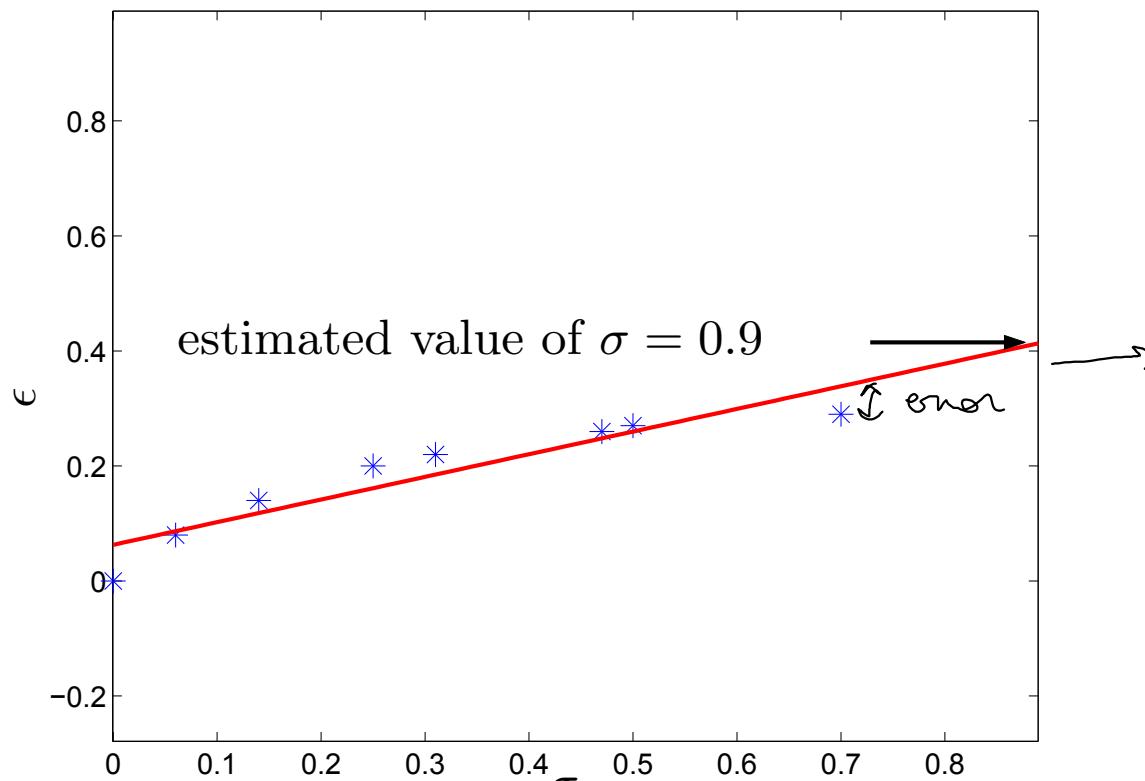


test $i$	pressure $\sigma$	deformation $\epsilon$
1	0.00	0.00
2	0.06	0.08
3	0.14	0.14
4	0.25	0.20
5	0.31	0.23
6	0.47	0.25
7	0.60	0.28
8	0.70	0.29



From these data, we want to estimate the deformation corresponding to a stress  $\sigma = 0.9$  MPa.

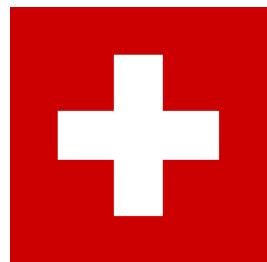
By the method of least squares, we get that best approximation for the data is  $p(x) = 0.3938x - 0.0629$ . The approximation can be used (called *a linear regression*) to estimate  $\epsilon$  when  $\sigma = 0.9$  MPa: We get  $p(0.9) \simeq 0.4$ .



we build a  
TREND by mini-  
mizing the distan-  
ce between  
data and line

! If you want to approx this date (could be  $10^2$ ,  $10^3$  data collected) with std interpolation you'll get a poly with high degree  $\rightarrow$  NOT GOOD!

**Example 2.** The result of census of the population of Switzerland between 1900 and 2010 (in thousands):

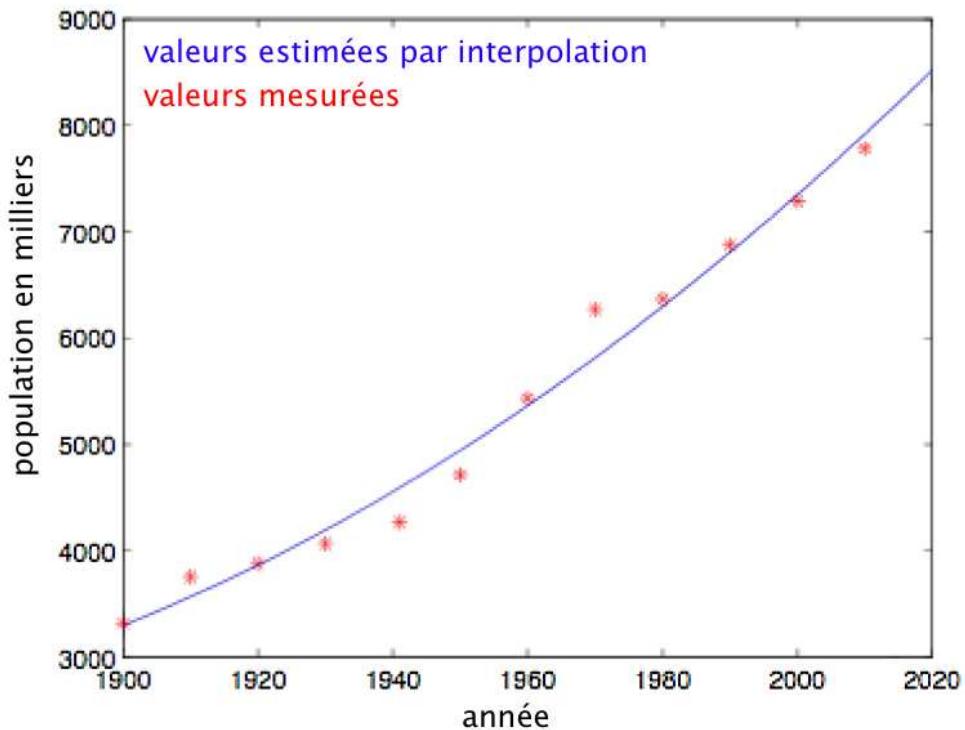


year	1900	1910	1920	1930	1941	1950
population	3315	3753	3880	4066	4266	4715
year	1960	1970	1980	1990	2000	2010
population	5429	6270	6366	6874	7288	7783

- Is it possible to estimate the number of inhabitants of Switzerland during the year when there has not been census, for example in 1945 and 1975?
- Is it possible to predict the number of inhabitants of Switzerland in 2020?

The polynomial of degree two (parabola) which approximates the data by the method of least squares is  $p(x) = 0.15x^2 - 549.9x + 501600$ .

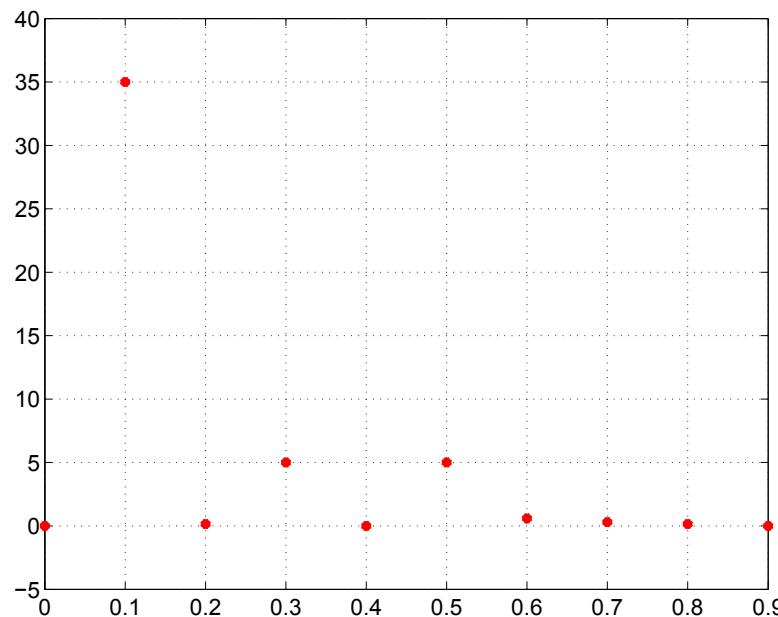
↑ II od poly. LEAST SQUARES



When fitting we should guarantee the minimization of a q'ty, like errors

**Example 3.** Points on the following figure represent the measurements of blood flow in a part of the common carotid artery during a heartbeat. The frequency of data collection is constant and equal to  $10 / T$  where  $T = 1$  sec. is the period of the beat.

We want to deduce from this discrete signal a continuous signal represented by a linear combination of known functions (eg. trigonometric functions - then we can adequately approximate a periodic signal).



reconstruction of the behavior of the flow rate

Let  $f(t)$  be a signal that we know.  $N = 10$  is a size of a sample  $[f(t_0), \dots, f(t_{N-1})]$ , where  $t_j = jT/N$ . We are looking for  $\{c_k\} \in \mathbb{C}$ ,  $k \in [0, N-1] \subset \mathbb{N}$  such that:

$$f(t_j) = \frac{1}{N} \sum_{k=0}^{N-1} c_k \omega_N^{-kj}, \quad j = 0, \dots, N-1,$$

← Basis function

↑ vector of weff is  
the only thing we  
have to get in order  
to get the expr.  
(1)  
of blood flow signal

where

$$\omega_N = e^{(-2\pi i)/N} = \cos(2\pi/N) - i \sin(2\pi/N),$$

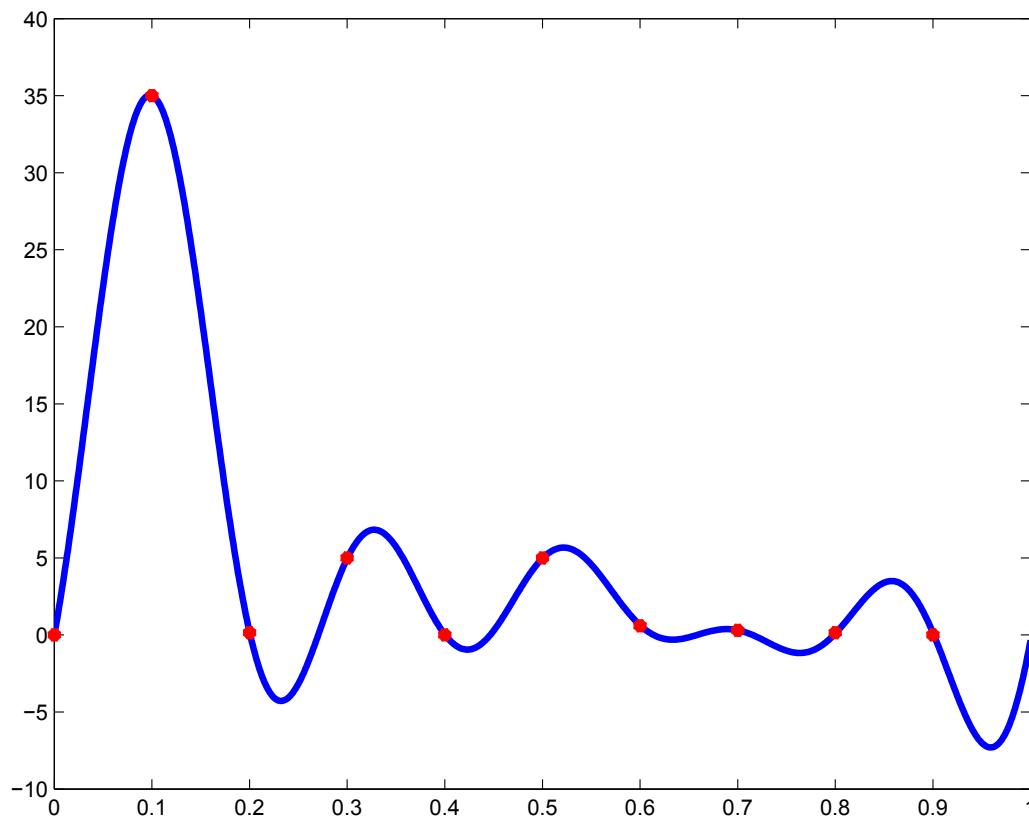
$i$  is the imaginary unit. We can calculate a vector of the coefficients  $\mathbf{c} = [c_1, \dots, c_{10}]^T$  by the FFT algorithm (Fast Fourier Transform, Cooley & Tuckey, 1965), implemented in Matlab/Octave by the **fft** command.

# of data

Look for this  
in python library

Using the formula (I), there are several techniques for obtaining the value of  $f$  in each point  $t \in [0, T]$ .

Using one of techniques we obtain a result that is plotted on a following figure:



If you use  
higher degree  
expect bigger  
costs; but if  
you reduce it,  
expect less  
accurate results.

# Approximation by the least square method

(Chapt. 3.4 of the book)

Suppose we have  $n + 1$  points  $x_0, x_1, \dots, x_n$  and  $n + 1$  values  $y_0, y_1, \dots, y_n$ . We have seen that if  $n$  is large, the interpolating polynomial may show large oscillations (*high degree!*)

Instead of interpolating the values, it is possible to define a polynomial of degree  $m < n$  that approximates the data “at best”

**Definition 1.** We call *least squares polynomial approximation of degree m* the polynomial  $\tilde{f}_m(x)$  of degree m such that

$$\sum_{i=0}^n |y_i - \tilde{f}_m(x_i)|^2 \leq \sum_{i=0}^n |y_i - p_m(x_i)|^2 \quad \forall p_m(x) \in \mathbb{P}_m$$

**Remark 4.** Where  $y_i = f(x_i)$  ( $f$  is a continuous function) then  $\tilde{f}_m$  is called the approximation of  $f$  in the least squares sense.

↑ minimizing the error by selecting  
 the best polynomial of some degree

In other words, the least squares polynomial approximation is the polynomial of degree  $m$  that minimizes the distance to the data.  $\rightsquigarrow$  solution of a minimize problem

Let note  $\tilde{f}_m(x) = \underbrace{a_0 + a_1x + a_2x^2 + \dots + a_mx^m}_{\text{Unknowns}}$  and define the function

$$\Phi(a_0, a_1, \dots, a_m) = \sum_{i=0}^n |y_i - (a_0 + a_1x_i + a_2x_i^2 + \dots + a_mx_i^m)|^2$$

$m+1$  unknowns  $\Rightarrow$  we need a system of equations

Then the coefficients of  $\tilde{f}^m$  can be determined by the relation

CN of being an extreme point (min or max)

$$\frac{\partial \Phi}{\partial a_k} = 0, \quad k = 0, \dots, m, \tag{5}$$

i.e.,  $m+1$  linear equations in with  $m+1$  unknowns  $a_k$ ,  $k = 0, \dots, m$ .

Usually 2 pictures

Look for best fit approx ( $\Rightarrow$  minimization of error)

Take the next one step (  $\rightarrow$  GREEDY algorithm )

depends on application  $\hookrightarrow$  you take the maximum

**Example 8.** Let  $n = 2$  and  $m = 1$ , nodes  $x_0 = 1, x_1 = 3, x_2 = 4$  with values  $y_0 = 0, y_1 = 2, y_2 = 7$ . We want to compute the (*regression line*), i.e., the least squares polynomial approximation of degree 1,  $\tilde{f}_1(x) = a_0 + a_1 x$ .

We set  $\Phi(a_0, a_1) = \sum_{i=0}^2 [y_i - (a_0 + a_1 x_i)]^2$  and impose  $\frac{\partial \Phi}{\partial a_0} = 0$  and  $\frac{\partial \Phi}{\partial a_1} = 0$ :

$$\begin{aligned}\frac{\partial \Phi}{\partial a_0} &= -2 \sum_{i=0}^2 [y_i - (a_0 + a_1 x_i)] = -2 \left( \sum_{i=0}^2 y_i - 3a_0 - a_1 \sum_{i=0}^2 x_i \right) \\ &= -2(9 - 3a_0 - 8a_1) \\ \frac{\partial \Phi}{\partial a_1} &= -2 \sum_{i=0}^2 x_i [y_i - (a_0 + a_1 x_i)] = -2 \left( \sum_{i=0}^2 x_i y_i - a_0 \sum_{i=0}^2 x_i - a_1 \sum_{i=0}^2 x_i^2 \right) \\ &= -2(34 - 8a_0 - 26a_1)\end{aligned}$$

Hence the coefficients  $a_0$  and  $a_1$  are the solution of the system

$$\begin{cases} 3a_0 + 8a_1 = 9 \\ 8a_0 + 26a_1 = 34 \end{cases}$$

Generalize to sd m

Ideally, for  $\tilde{f}_m(x) = a_0 + a_1x + a_2x^2 + \dots$  we would like to impose  $\tilde{f}_m(x_i) = y_i$  pour  $i = 0, \dots, n$ .

This can be written as a linear system with unknowns  $a_k$ ,  $k = 0, \dots, m$ :  $B\mathbf{a} = \tilde{\mathbf{y}}$ , where  $B$  is a matrix of dimension  $(n+1) \times (m+1)$

$$B = \left( \begin{array}{cccc} 1 & x_0 & \dots & x_0^m \\ 1 & x_1 & \dots & x_1^m \\ \vdots & & & \vdots \\ 1 & x_n & \dots & x_n^m \end{array} \right) \quad \begin{array}{l} \text{dim of col of approx} \\ \text{dim of detent} \end{array}$$

A rectangular  $B$  matrix is NOT solvable as a linear system

Since  $m < n$ , the system is oversized. The solution to (5) is equivalent to the square system (*system of normal equations*)

$$B^T B \mathbf{a} = B^T \tilde{\mathbf{y}}$$

↑  
pre-multiply by  $B^T$  → you get a square matrix,  
now you can solve shit

**Example 9.** We have 8 measures ( $\sigma$ - $\epsilon$ ):

```
>> sigma = [0.00 0.06 0.14 0.25 0.31 0.47 0.50 0.70];
>> epsilon = [0.00 0.08 0.14 0.20 0.22 0.26 0.27 0.29];
```

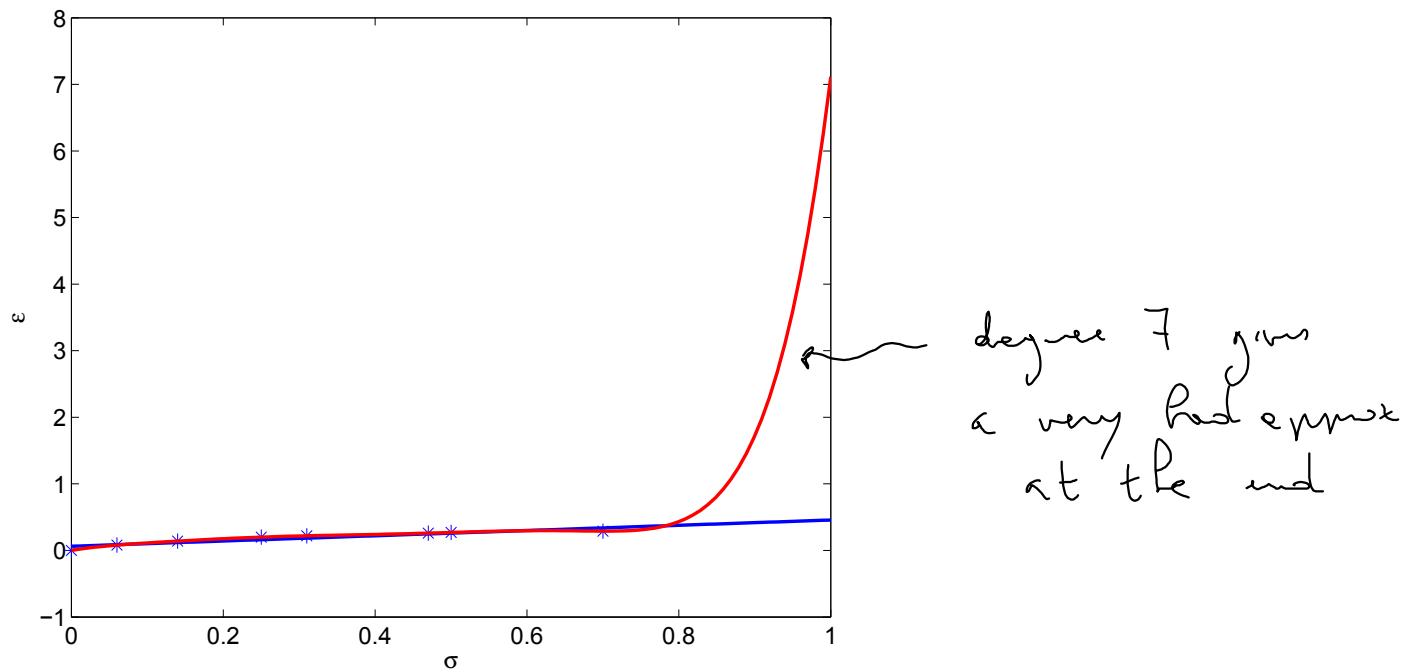
We want to extrapolate the value of  $\epsilon$  for  $\sigma = 0.4$ . We consider two ways:

- compute the interpolating polynomial  $\Pi_7$  of degree 7
- compute the least squares polynomial approximation of degree 1

On peut utiliser les commandes suivantes:

```
>> plot(sigma, epsilon, '*'); hold on; % plotting the known values
>> sigma_sample = linspace(0,1.0,100); de la
>> p7 = polyfit(sigma, epsilon, 7); degre
>> pol = polyval(p7, sigma_sample); % interpolating polynomial
>> plot(sigma_sample, pol, 'r'); somme tell me que tu es une ligne
>> p1 = polyfit(sigma, epsilon, 1); somme tell me que tu es une ligne
>> pol_mc = polyval(p1, sigma_sample); % least square
>> plot(sigma_sample, pol_mc, 'g'); hold off;
```

- the interpolating polynomial  $\Pi_7$  of degree 7 is in red
- the least squares polynomial approximation of degree 1 is in blue



For  $\sigma > 0.7$ , the behavior of the two polynomials are very different.

In particular, for  $\sigma = 0.9$  the values of  $\epsilon(\sigma)$  extrapolated with the two methods are

```
>> polyval(p7, 0.9)
```

```
ans =
```

1.7221

```
>> polyval(p1, 0.9)
```

```
ans =
```

0.4173

$\curvearrowleft$  evaluation of f.t at 0.9

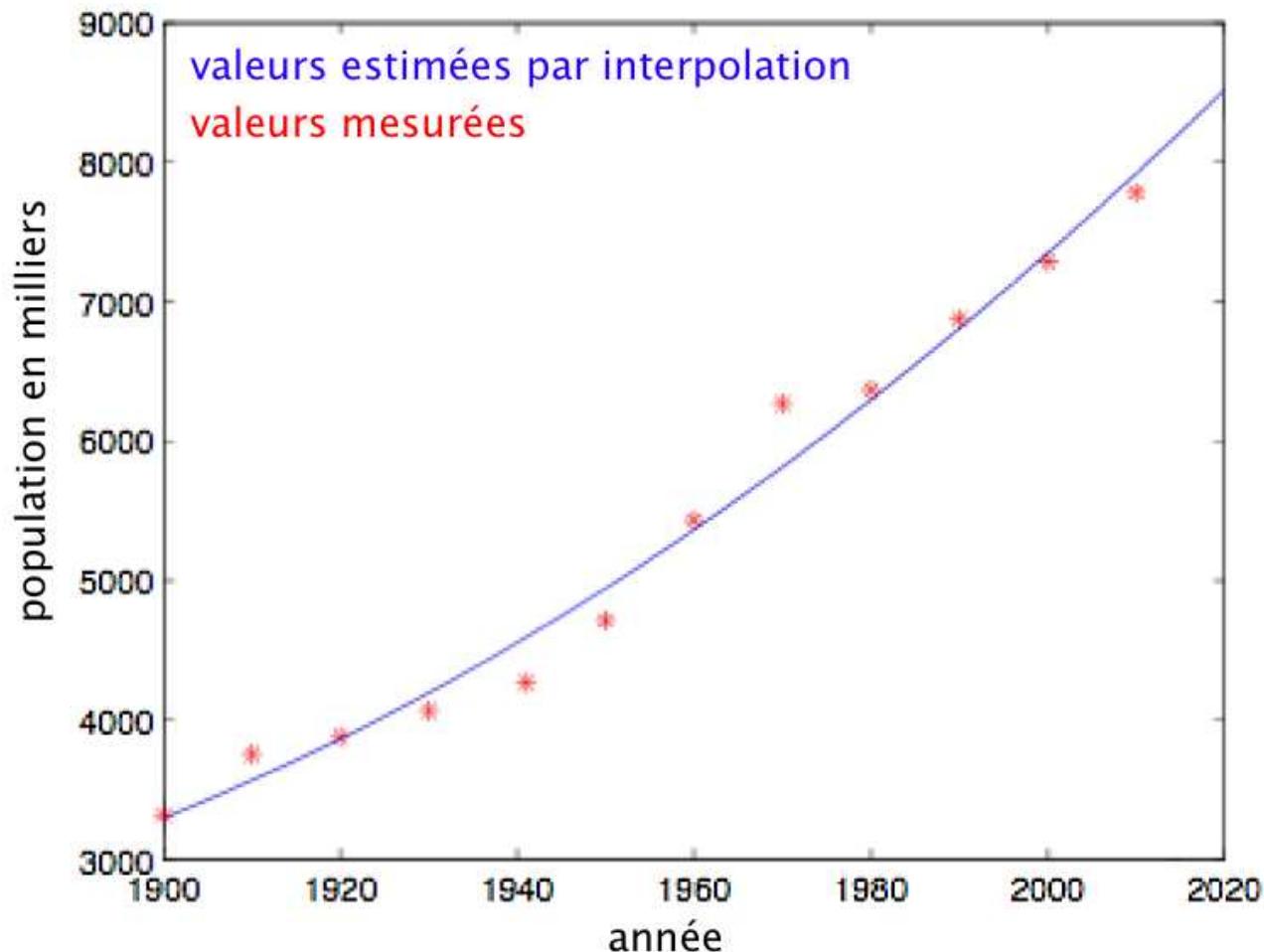
- The value obtained by  $\Pi_7$  (172.21%) is unrealistic
- On the contrary, the value obtained with the regression line is more appropriate to compute the value at  $\sigma = 0.9$ .

1D

Since the results, in EQUATIONS,  
multiplied by  $100$ , you have the  
percentage of elongation  
having  $\sim 200\%$ . of  
elongation would obtain  
an absurd non-physical result

**Example 10.** Swiss population Starting from the population at the 20th century decades, we extrapolate the Swiss population this year. We use a least square polynomial approximation of degree 2

```
>> year = [1900, 1910, 1920, 1930, 1941, 1950, ...
    1960, 1970, 1980, 1990, 2000];
>> population = [3315, 3753, 3880, 4066, 4266, 4715, ...
    5429, 6270, 6366, 6874, 7288];
>> p2 = polyfit(year, population, 2);
>> year_sample = [1900:2:2010];
>> vp2 = polyval(p2, year_sample);
>> plot(year, population, '*r', ...
    year_sample, vp2, 'b');
```

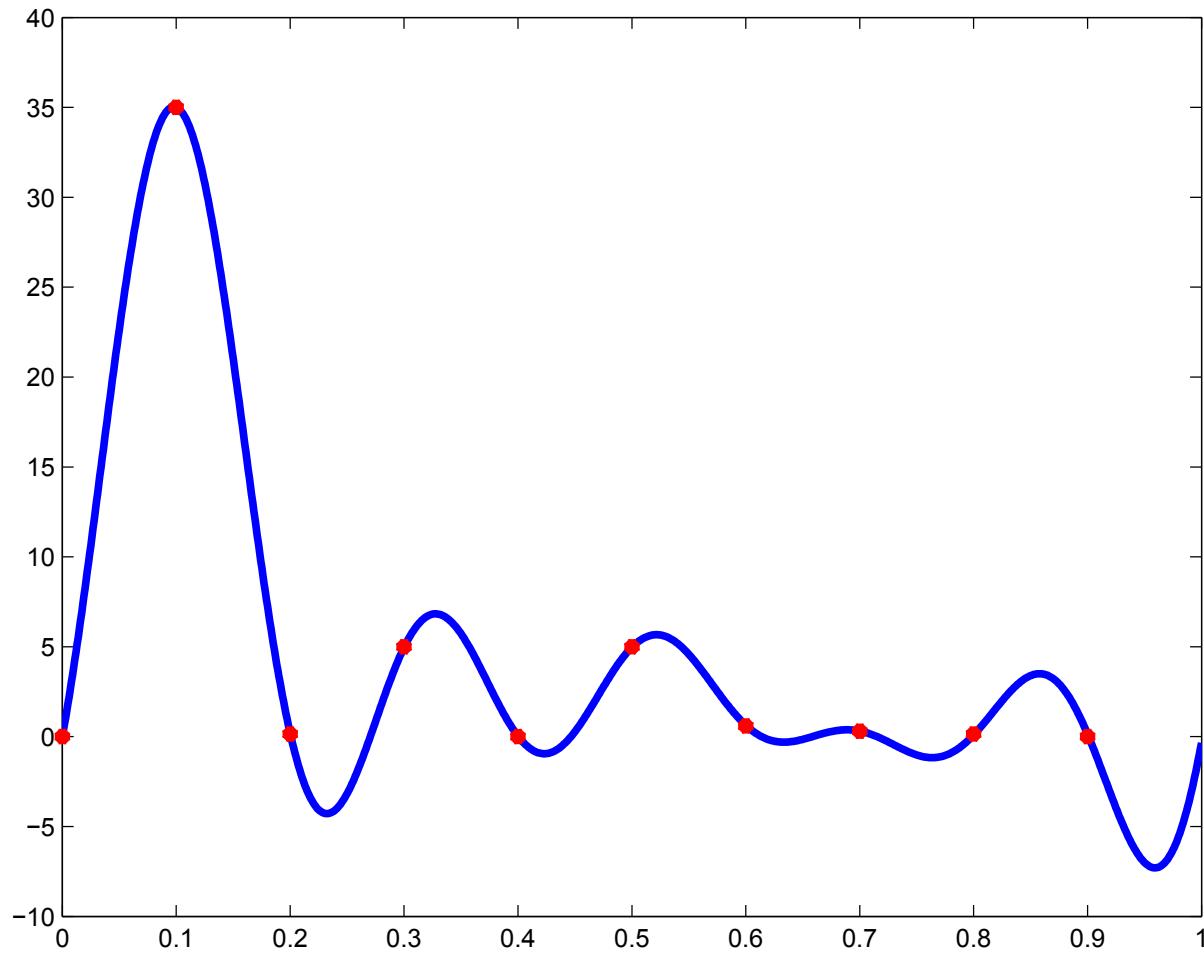


**Example. Carotide flow rate** We define 10 couples of data (time vs flowrate) with the vectors **t** and **deb**. With the help of the command **interpft**, we compute the value of the interpolating function at 1000 equidistributed points in the periodicity interval [0,1]:

```

t = [0 .1 .2 .3 .4 .5 .6 .7 .8 .9];
deb = [0 35 .15 5 0 5 .6 .3 .15 0];
f = interpft(deb, 1000); interpolation Fourier t-f
plot(linspace(0, 1, 1000), f); hold on;
plot(t, deb, 'r*')

```



# Processing of polynomials

In Matlab/Octave, there are specific commands for doing calculations with polynomials. Let  $x$  be a vector of abscissas,  $y$  be a vector of ordinates and  $p$  (respectively  $p_i$ ) be the vector of coefficients of a polynomial  $P(x)$  (respectively  $P_i$ ); then, we have following commands:

command	action
<code>y=polyval(p,x)</code>	$y = \text{values of } P(x)$
<code>p=polyfit(x,y,n)</code>	$p = \text{coefficients of the interpolating polynomial } \Pi_n$
<code>z=roots(p)</code>	$z = \text{zeros of } P \text{ such that } P(z) = 0$
<code>p=conv(p1,p2)</code>	$p = \text{coefficients of the polynomial } P_1 P_2$
<code>[q,r]=deconv(p1,p2)</code>	$q = \text{coefficients of } Q, r = \text{coefficients of } R$ $\text{such that } P_1 = QP_2 + R$
<code>y=polyderiv(p)</code>	$y = \text{coefficients of } P'(x)$
<code>y=polyinteg(p)</code>	$y = \text{coefficients of } \int P(x) dx$

# Linear systems - Direct methods



## Numerical Analysis

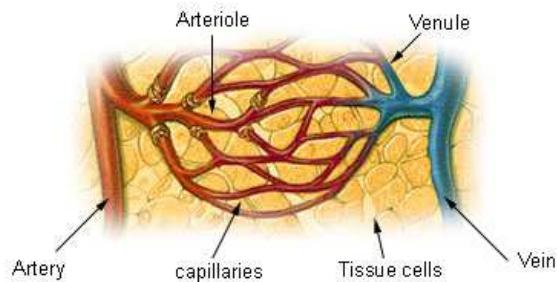
Profs. Gianluigi Rozza - Luca Heltai

2019-SISSA mathLab Trieste

# Examples and motivation

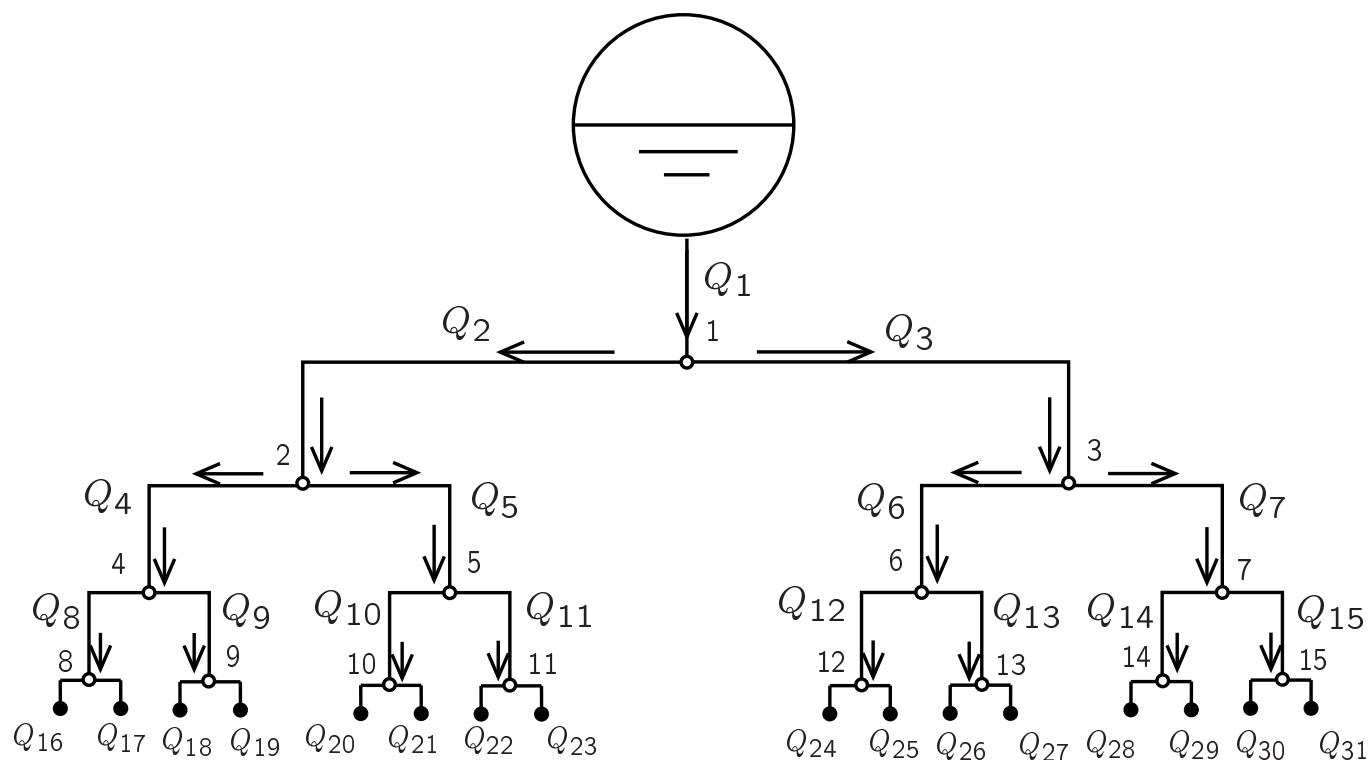
Capillaries are microscopic blood vessels, the smallest part of the circulatory system. The capillaries are grouped in networks called capillary beds, made of 10 to 100 capillaries, depending on the tissue.

- Example 1** (Capillary bed). Blood arrives to the capillary bed through arterioles. In the capillary bed an exchange of oxygen and waste molecules takes place. Then, the capillaries become venules and collect the blood in veins to be transferred back to the heart.



We consider a model of capillary bed:

- We can model a certain portion of capillary system as a hydraulic network where every capillary is represented as a pipe.
- We call *nodes* (small empty circles in the figure on the next page) the points where several capillaries meet.
- The artery that feeds the system is represented as source at a constant pressure of 50mmHg.
- We suppose the blood leaves the system through the venules (small black circles in the figure) at a constant pressure (venous pressure), fixed to the reference value zero (all other pressures will be referenced to this one).
- The blood flows from the source to the sinks (we suppose in a constant way) due to the pressure gradient.



We want to find the pressure distribution  $p_j$ ,  $j = 1, \dots, 15$ , and the flows  $Q_m$ ,  $m = 1, \dots, 31$ , in the capillary network. To do this, we consider that:

- in every branch  $m$  of the network,  $m = 1, \dots, 31$ , we have the following relation, called **constitutive relation**, between the blood flow  $Q_m$  ( $\text{mm}^3/\text{s}$ ) and the pressure (in mmHg) in the two nodes  $i$  and  $j$  at the end-points of each capillary

$$Q_m = \frac{1}{R_m L_m} (p_i - p_j), \quad (1)$$

where  $R_m$  is the hydraulic resistance per unit length ( $\text{mmHg s/mm}^4$ ) and  $L_m$  is the length of the capillary  $m$ ;

- at every node  $j$  of the network,  $j = 1, \dots, 15$ , we impose the **balance equation** between inflow and outflow:

$$\left( \sum_{m \text{ entering}} Q_m \right) - \left( \sum_{m \text{ exiting}} Q_m \right) = 0,$$

where outflows have a negative value.

For example, at the node 2 in the figure, the flow  $Q_2$  is an inflow and the flows  $Q_4$  and  $Q_5$  are outflows, being the balance:

$$Q_2 - Q_4 - Q_5 = 0.$$

By using for every flow the constitution relation (I) in the previous balance, we have

$$\frac{1}{R_2 L_2} (p_1 - p_2) - \frac{1}{R_4 L_4} (p_2 - p_4) - \frac{1}{R_5 L_5} (p_2 - p_5) = 0.$$

One such equation is obtained for every node in the network.

Remark: if we consider the balance for the node 1, we get

$$\frac{1}{R_1 L_1} (p_r - p_1) - \frac{1}{R_2 L_2} (p_1 - p_2) - \frac{1}{R_3 L_3} (p_1 - p_3) = 0.$$

Being the pressure  $p_r$  constant, we move it to the right-hand side.

$$-\frac{1}{R_1 L_1} p_1 - \frac{1}{R_2 L_2} (p_1 - p_2) - \frac{1}{R_3 L_3} (p_1 - p_3) = -\frac{1}{R_1 L_1} p_r.$$



In a similar way for the node 8

$$\frac{1}{R_8 L_8} (p_4 - p_8) - \frac{1}{R_{16} L_{16}} p_8 - \frac{1}{R_{17} L_{17}} p_8 = 0,$$

where we see that  $p_{16} = p_{17} = 0$ . Same process for the nodes  $9, \dots, 15$ .

Writing the balance for every node after having substituted the constitution relation, we get a system of linear equations for the pressures:

$$A\mathbf{p} = \mathbf{b}, \tag{2}$$

$\nwarrow$  A contains all the properties of our system

where  $\mathbf{p} = [p_1, p_2, \dots, p_{15}]^T$  is the unknown vector pressures at the nodes of the network,  $A \in \mathbb{R}^{15 \times 15}$  is the coefficient matrix of the system and  $\mathbf{b} \in \mathbb{R}^{15}$  is the array of known data.

*! Do evaluations of g(b) as you have for multiple use, in order to be more efficient and avoid "reputation" → save time and computational costs*

Lets suppose that all capillaries have the same hydraulic resistance  $R_m = R = 1$  and that capillary length halves at each bifurcation (if  $L_1 = 20$ , we'll have  $L_2 = L_3 = 10$ ,  $L_4 = L_5 = L_6 = L_7 = 5$  etc..), the following matrix A is generated:

$$\left( \begin{array}{cccccccccccccccc} -\frac{1}{4} & \frac{1}{10} & \frac{1}{10} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \frac{1}{10} & -\frac{1}{2} & 0 & \frac{1}{5} & \frac{1}{5} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \frac{1}{10} & 0 & -\frac{1}{2} & 0 & 0 & \frac{1}{5} & \frac{1}{5} & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \frac{1}{5} & 0 & -1 & 0 & 0 & 0 & 0.4 & 0.4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \frac{1}{5} & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0.4 & 0.4 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{5} & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0.4 & 0.4 & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{5} & 0 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0.4 & 0.4 & 0 \\ 0 & 0 & 0 & 0.4 & 0 & 0 & 0 & -2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.4 & 0 & 0 & 0 & 0 & -2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.4 & 0 & 0 & 0 & 0 & -2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0.4 & 0 & 0 & 0 & 0 & 0 & -2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.4 & 0 & 0 & 0 & 0 & 0 & -2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.4 & 0 & 0 & 0 & 0 & 0 & 0 & -2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.4 & 0 & 0 & 0 & 0 & 0 & 0 & -2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -2 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.4 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -2 \end{array} \right)$$

and

$$\mathbf{b} = [-5/2, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]^T.$$

The problem of determining the pressures and flows in the network requires to solve a linear system  $A\mathbf{p} = \mathbf{b}$ .

In this case, the matrix  $A$  is symmetric and definite positive. This last property ensures that the matrix  $A$  is not singular and thus that the system has a unique solution.

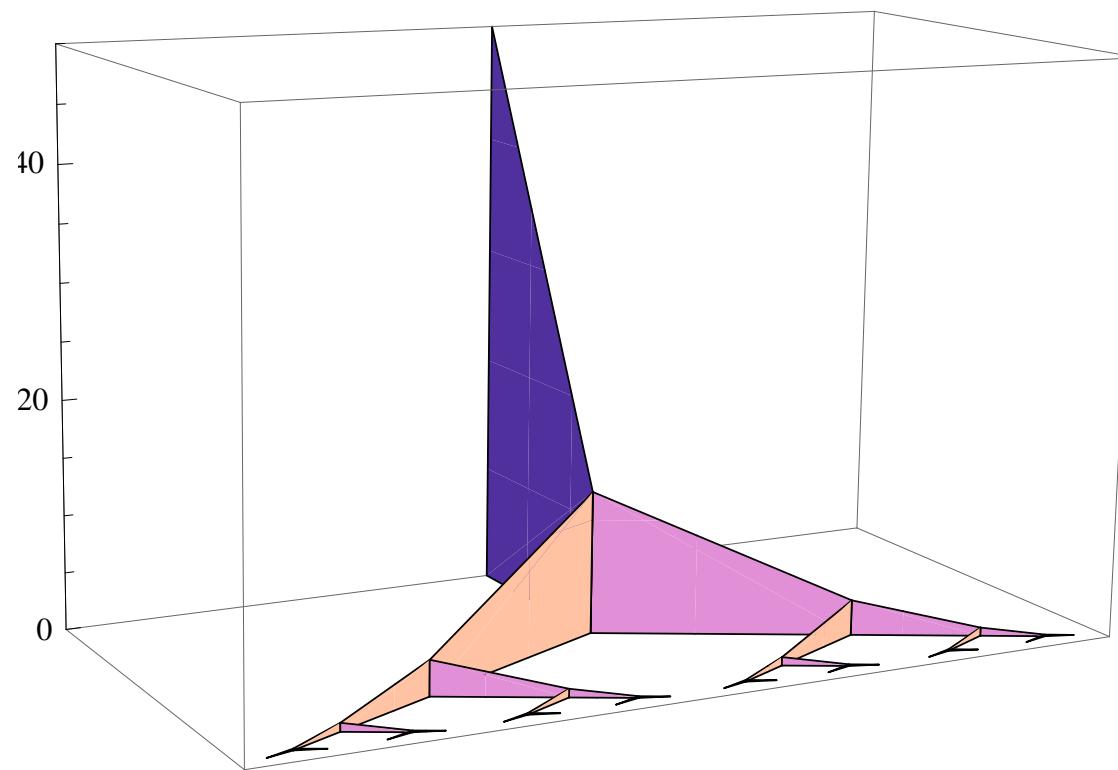
The solution is given by the vector

$$\mathbf{p} = [12.46, 3.07, 3.07, 0.73, 0.73, 0.73, 0.73, 0.15, 0.15, 0.15, 0.15, 0.15, 0.15, 0.15, 0.15]^\top.$$

Once we have the pressures, we can compute, using the relation (I), the flows:

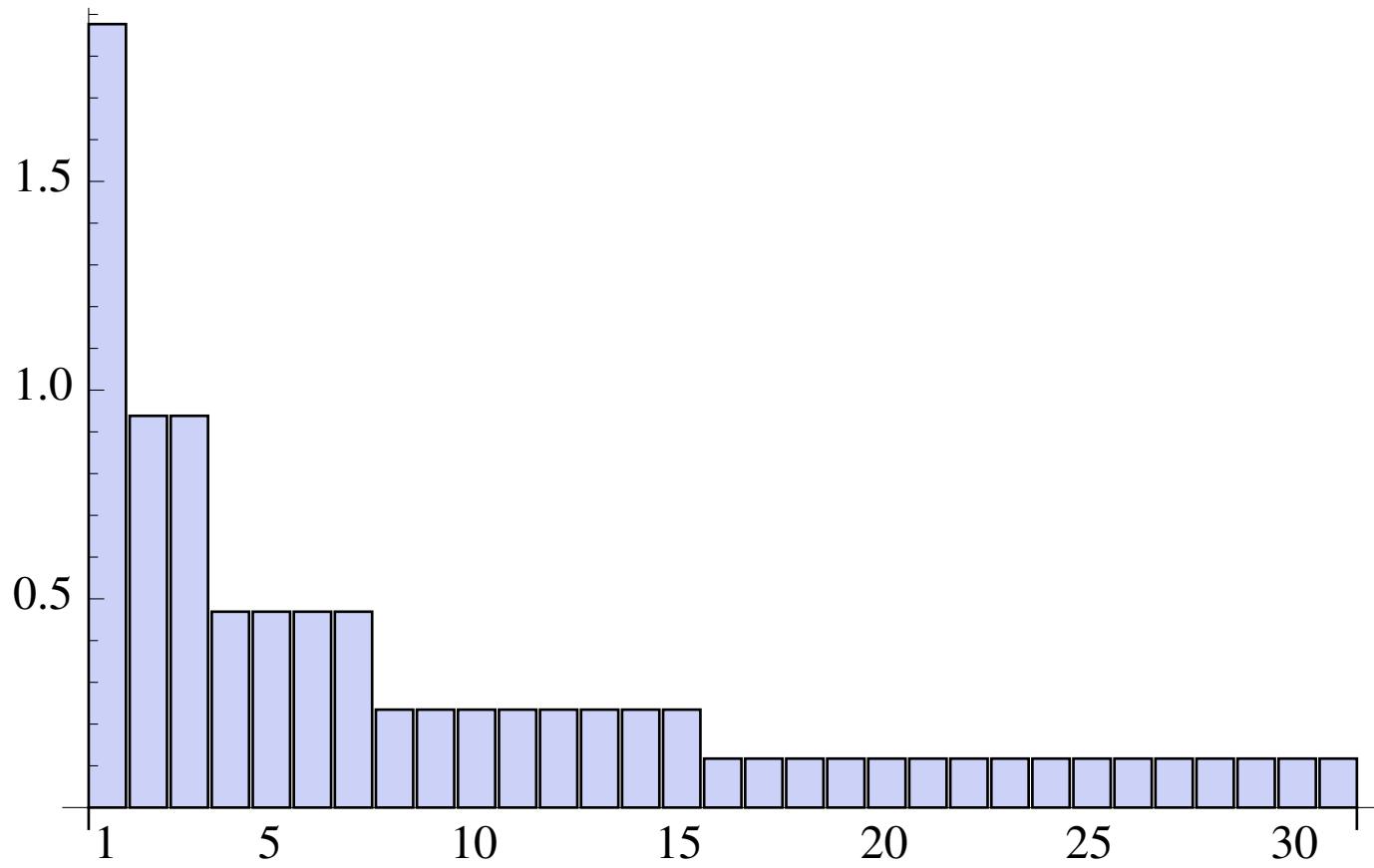
$$\begin{aligned}
 Q_1 &= 1.88 \\
 Q_{2,3} &= 0.94 \\
 Q_{4,\dots,7} &= 0.47 \\
 Q_{8,\dots,15} &= 0.23 \\
 Q_{16,\dots,31} &= 0.12
 \end{aligned}$$

Linear distribution of pressures in the capillary bed computed from the pressures at every node (solution of the system)

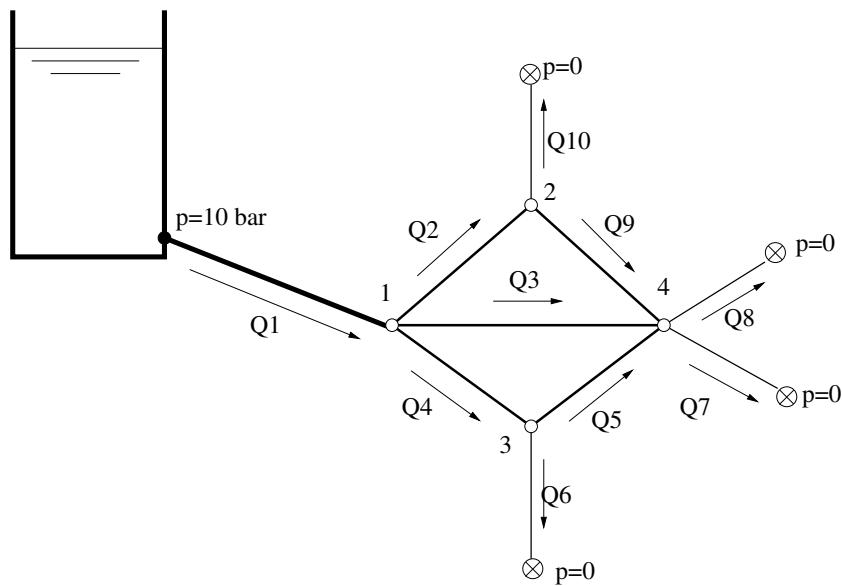


$\leftarrow$  system is dumping energy  
 we have a drop in pressure

## Distribution of the flows in the capillary bed



**Example 2** (Hydraulic network). Let us consider a hydraulic network made of 10 pipes.



The network is fed by a reservoir of water at constant pressure  $p_r = 10$  bar (pressure values in this exercise are the difference between the “real” pressure and the atmospheric pressure, in bar). For each pipe in the network, the following relation between the water flow  $Q$  ( $\text{m}^3/\text{s}$ ) and the pressure (in bar) holds at both pipe-ends

$$Q = \frac{1}{RL}(p_{\text{in}} - p_{\text{out}}), \quad (3)$$

where  $R$  is the hydraulic resistance per unit length ((bar s)/ $\text{m}^4$ ) and  $L$  is the length (in m) of the pipeline. The resistance  $R$  depends on the pipe radius  $r$  and the dynamic viscosity  $\mu$  of the liquid, following the relation

$$R = 8\mu/(\pi r^4).$$

but again  $R$  depends on geom qty ( $r$ ) too

At every node of the network, the balance between inflows and outflows can be set; for example for the node 2 in the figure, we have  $Q_2 - Q_9 - Q_{10} = 0$  (outflows have a negative sign).

We assume atmospheric pressure at the outlets ( $p = 0$  bar). A typical problem consists the pressure values and flows in the network. Using both given relations, a linear system can be built for the pressure of a form:

$$A\mathbf{p} = \mathbf{b}, \quad (4)$$

where  $\mathbf{p} = [p_1, p_2, p_3, p_4]^T$  is the vector of pressures (unknown) at the 4 nodes of the network.

↖ unknowns

The following table contains the relevant characteristics of the pipelines:

Pipe	R	L	Pipe	R	L	Pipe	R	L
1	0.2500	20	2	2.0000	10	3	1.0204	14
4	2.0000	10	5	2.0000	10	6	7.8125	8
7	7.8125	8	8	7.8125	8	9	2.0000	10
10	7.8125	8						

Correspondingly,  $A$  and  $\mathbf{b}$  are:

$$A = \begin{bmatrix} -0.370 & 0.050 & 0.050 & 0.070 \\ 0.050 & -0.116 & 0 & 0.050 \\ 0.050 & 0 & -0.116 & 0.050 \\ 0.070 & 0.050 & 0.050 & -0.202 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}.$$

↘ *symmetric!*

The array of pressure values at the nodes is:

$\mathbf{p} = A^{-1}\mathbf{b} = [8.1172, 5.9893, 5.9893, 5.7779]^T$ . We can finally compute, using the relation (II), the flows (in  $\text{m}^3/\text{s}$ ):

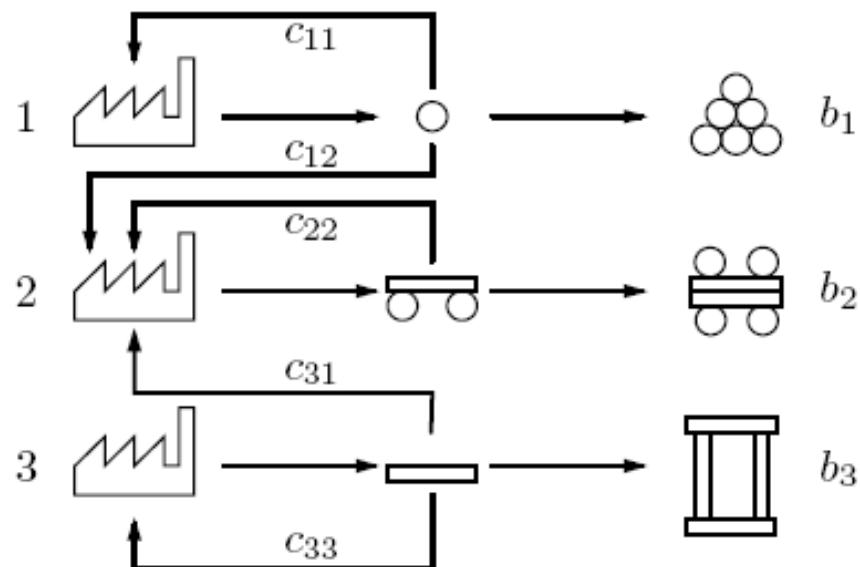
Q1	0.3766	Q2	1.0640	Q3	0.1638
Q4	0.1064	Q5	0.0106	Q6	0.0958
Q7	0.0924	Q8	0.0924	Q9	0.0106
Q10	0.0958				

**Example 3** (Economy/Logistic). We want to determine the situation of equilibrium between demand and offer of certain goods. Let us consider that  $m \geq n$  factories produce  $n$  different products. They have to adapt their productions to the internal demand (i.e. the goods needed as input by the other factories) as well as to the external demand, from the consumers.

$x_i$ ,  $i = 1, \dots, n$  is the total number of goods made by the factory  $i$ ,

$b_i$ ,  $i = 1, \dots, n$  is the corresponding demand from the market and

$c_{ij}$  the amount produced by the factory  $i$  needed for the factory  $j$  to make one unit of product.



*Interaction scheme between 3 factories and the market*

If we suppose that the relation between the different products is linear, the equilibrium is reached when the vector  $\mathbf{x} = [x_1, \dots, x_n]^T$  satisfies

$$\mathbf{x} = C\mathbf{x} + \mathbf{b},$$

interaction  
 identifying  
 internal and  
 of intermediate products

market demand

where  $C = (c_{ij})$  and  $\mathbf{b} = [b_1, \dots, b_n]^T$ . Consequently, the total production  $\mathbf{x}$  is solution of the linear system :

$$A\mathbf{x} = \mathbf{b}, \quad \text{where } A = I - C.$$

# Formulation of the problem

We call **linear system of order  $n$**  ( $n$  positive integer), an expression of the form

$$A\mathbf{x} = \mathbf{b},$$

↗ unknowns  
 ↑ vector of data (given)

where  $A = (a_{ij})$  is a given matrix of size  $n \times n$ ,  $\mathbf{b} = (b_j)$  is a **given vector** and  $\mathbf{x} = (x_j)$  is the **unknown vector** of the system. The previous relation is equivalent to the  $n$  equations

↗ characterizing solution of our problem

$$\sum_{j=1}^n a_{ij}x_j = b_i, \quad i = 1, \dots, n.$$

The matrix  $A$  is called **non-singular** if  $\det(A) \neq 0$ ; the solution  $\mathbf{x}$  will be **unique** (for any given vector  $\mathbf{b}$ ) if and only if the matrix associated to the linear system is non-singular.

⇒ must avoid singular matrices

↗ if determinant = 0 ⇒ it is not unique  
 ⇒ not deterministic outcome

In theory, if  $A$  is non-singular, the solution is given by the *Cramer's rule*:

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

where  $B_i$  is the matrix by substituting the  $i$ -th column of  $A$  by the vector  $\mathbf{b}$ :

$$B_i = \begin{bmatrix} a_{11} & \dots & b_1 & \dots & a_{1n} \\ a_{21} & \dots & b_2 & \dots & a_{2n} \\ \vdots & & \vdots & & \vdots \\ a_{n1} & \dots & b_n & \dots & a_{nn} \end{bmatrix}$$

↑  
 $i$

NOT EASY TO COMPUTE  $\det(B_i)$   
 Because it requires  
 $\sim N!$  flops (floating point operation)

⇒ COMPUTATIONAL COST OF CRAMER'S RULE IS  $\sim (N+1)!$

↑  $N=20 \Rightarrow$  you say is  
FUCKED!

Unfortunately, the application of this rule is unacceptable for the practical solution of systems because the **computational cost is of the order of  $(n + 1)!$  floating point operations (flops)**. In fact, every determinant requires  $n!$  flops.

For example, the following table gives the time required by different computers to solve a linear system using the Cramer rule (o.r. stands for “out of reach”):

$n$	Number of flops of the computer				
	$10^9$ (Giga)	$10^{10}$	$10^{11}$	$10^{12}$ (Tera)	$10^{15}$ (Peta)
10	$10^{-1}$ sec	$10^{-2}$ sec	$10^{-3}$ sec	$10^{-4}$ sec	negligible
15	17 hours	1.74 hours	10.46 min	1 min	$6 \cdot 10^{-2}$ sec
20	4860 years	486 years	48.6 years	4.86 years	1.7 days
25	o.r.	o.r.	o.r.	o.r.	38365 years

Alternative algorithms have to be developed. In the following sections, several methods will be analysed.

# Triangular systems

A matrix  $U = (u_{ij})$  is *upper triangular* if

$$u_{ij} = 0 \quad \forall i, j : 1 \leq j < i \leq n$$

and a matrix  $L = (l_{ij})$  is *lower triangular* if

$$l_{ij} = 0 \quad \forall i, j : 1 \leq i < j \leq n.$$

Respectively, the system to be solved is called *upper or lower triangular system*.

Remark: If a matrix  $A$  is non-singular and triangular, knowing that

$$\det(A) = \prod_{i=1}^n \lambda_i(A) = \prod_{i=1}^n a_{ii}$$

( $\lambda_i(A)$  being the  $i$ -th eigenvalue of  $A$ ), we can deduce that  $a_{ii} \neq 0$ , for all  $i = 1, \dots, n$ .

If  $L$  is lower triangular and non-singular, the linear system  $Ly = b$  corresponds to

$$\begin{cases} l_{11}y_1 &= b_1 \\ l_{21}y_1 + l_{22}y_2 &= b_2 \\ \vdots & \\ l_{n1}y_1 + l_{n2}y_2 + \dots + l_{nn}y_n &= b_n \end{cases}$$

Thus:

$$y_1 = b_1 / l_{11}, \quad [1 \text{ operation}]$$

and for  $i = 2, 3, \dots, n$

$$y_i = \frac{1}{l_{ii}} \left( b_i - \sum_{j=1}^{i-1} l_{ij} y_j \right).$$

already computed  
 at previous steps

[1 + 2(i - 1) operations]

This algorithm is called *forward substitutions algorithm*.

The forward substitutions algorithm requires  $n^2$  operations, where  $n$  is the size of the system:

$$\begin{aligned} 1 + \sum_{i=2}^n (1 + 2(i-1)) &= 1 + (n-1) + 2 \sum_{i=2}^n (i-1) \\ &= 1 + (n-1) + 2 \frac{n(n-1)}{2} \\ &= n^2. \end{aligned}$$

If  $U$  is upper triangular and non-singular, the system  $U\mathbf{x} = \mathbf{y}$  is:

$$\left\{ \begin{array}{lcl} u_{11}x_1 + \dots + u_{1,n-1}x_{n-1} + u_{1n}x_n & = & y_1 \\ & \vdots & \\ u_{n-1,n-1}x_{n-1} + u_{n-1,n}x_n & = & y_{n-1} \\ u_{nn}x_n & = & y_n \end{array} \right.$$

Thus:

$$x_n = y_n / u_{nn},$$

*we proceed with inverted order*

and for  $i = n - 1, n - 2, \dots, 2, 1$

$$x_i = \frac{1}{u_{ii}} \left( y_i - \sum_{j=i+1}^n u_{ij}x_j \right).$$

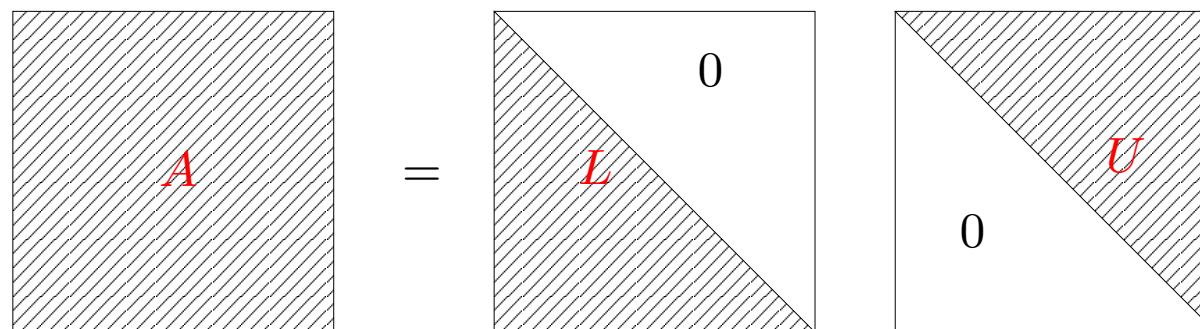
This algorithm is called *backward substitutions algorithm*. The cost is, once again,  $n^2$  operations.

# The $LU$ factorization method

(Sec. 5.3 of the book)

Let  $A = (a_{ij})$  be a non-singular  $n \times n$  matrix. Assume that there exist a matrix  $U = (u_{ij})$ , **upper triangular** and a matrix  $L = (l_{ij})$ , **lower triangular** such that

$$A = LU. \quad (5)$$



We call (5) a *factorization LU* of  $A$ .

If we know the factorization  $LU$  of  $A$ , solving the system  $A\mathbf{x} = \mathbf{b}$  is equivalent to solving two systems defined by *triangular* matrices. Indeed,

$$A\mathbf{x} = \mathbf{b} \Leftrightarrow LU\mathbf{x} = \mathbf{b} \Leftrightarrow \begin{cases} L\mathbf{y} = \mathbf{b}, \\ U\mathbf{x} = \mathbf{y}. \end{cases}$$

We can easily calculate the solutions of both systems:

- first, we use the forward substitutions algorithm to solve  $L\mathbf{y} = \mathbf{b}$  (order  $n^2$  flops);
- then, we use the backward substitutions algorithm to solve  $U\mathbf{x} = \mathbf{y}$  (order  $n^2$  flops).

It is required to find first (if possible) the matrices  $L$  and  $U$  (what requires a number of operations of the order  $\frac{2n^3}{3}$  flops).

**Example 4.** Lets try to find a factorization  $LU$  in the case where the size of the matrix  $A$  is  $n = 2$ . We can write the equation (5) as

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

Or equivalently:

- (a)  $l_{11}u_{11} = a_{11}$ ,    (b)  $l_{11}u_{12} = a_{12}$ ,
- (c)  $l_{21}u_{11} = a_{21}$ ,    (d)  $l_{21}u_{12} + l_{22}u_{22} = a_{22}$ .

We have then a system (non-linear) with 4 equations and 6 unknowns; in order to have the same number of equations and unknowns, we fix the diagonal of  $L$  by taking  $l_{11} = l_{22} = 1$ . Consequently, from (a) and (b) we have  $u_{11} = a_{11}$  and  $u_{12} = a_{12}$ ; finally, if we assume  $a_{11} \neq 0$ , we obtain  $l_{21} = a_{21}/a_{11}$  and  $u_{22} = a_{22} - l_{21}u_{12} = a_{22} - a_{21}a_{12}/a_{11}$  using the equations (c) and (d).

To determine a factorization  $LU$  of the matrix  $A$  of any size  $n$ , we apply the following method.

1. The elements of  $L$  and  $U$  satisfy the non-linear system

$$\sum_{r=1}^{\min(i,j)} l_{ir} u_{rj} = a_{ij}, \quad i, j = 1, \dots, n; \quad (6)$$

2. The system (6) has  $n^2$  equations and  $n^2 + \underbrace{n}_\text{we have the diagonals of L and U}$  unknowns. We can wipe out  $n$  unknowns if we set the  $n$  diagonal elements of  $L$  equal to 1:

$$l_{ii} = 1, \quad i = 1, \dots, n.$$

We will see that in this case there exists an algorithm (*Gauss factorization*) allowing us to efficiently compute the factors  $L$  and  $U$ .

# The Gauss elimination method

The Gauss elimination method transforms the system

$$Ax = \mathbf{b}$$

in an equivalent system (i.e. with the same solution) of the form:

$$Ux = \hat{\mathbf{b}}$$

If systems have also  
 Infinitely many solutions!

where  $U$  is an upper triangular matrix and  $\hat{\mathbf{b}}$  is a properly modified second member.

This system can be solved by a backward substitutions method.

In the transformation, we essentially use the property that says that we do not change the solution of the system if we add to a given equation a linear combination of other equations.

Let us consider an invertible matrix  $A \in \mathbb{R}^{n \times n}$  in which the diagonal element  $a_{11}$  is assumed to be non-zero. we set  $A^{(1)} = A$  and  $\mathbf{b}^{(1)} = \mathbf{b}$ . We introduce the *multiplier*

$$l_{i1} = \frac{a_{i1}^{(1)}}{a_{11}^{(1)}}, \quad i = 2, 3, \dots, n, \quad A^{(1)} = \begin{bmatrix} a_{11}^{(1)} & \dots & a_{1j}^{(1)} & \dots & a_{1n}^{(1)} \\ \vdots & & \vdots & & \vdots \\ a_{i1}^{(1)} & \dots & a_{ij}^{(1)} & \dots & a_{in}^{(1)} \\ \vdots & & \vdots & & \vdots \\ a_{n1}^{(1)} & \dots & a_{nj}^{(1)} & \dots & a_{nn}^{(1)} \end{bmatrix}$$

where the  $a_{ij}^{(1)}$  represent the elements of  $A^{(1)}$ . Example:

$$A = \begin{bmatrix} 2 & 4 & 6 \\ 4 & 8 & 10 \\ 7 & 8 & 9 \end{bmatrix} \implies l_{21} = \frac{4}{2}, l_{31} = \frac{7}{2}.$$

The unknown  $x_1$  can be removed from the rows  $i = 2, \dots, n$  by subtracting  $l_{i1}$  times the first row and doing the same at the right-hand side.

Let us define

$$a_{ij}^{(2)} = a_{ij}^{(1)} - l_{i1}a_{1j}^{(1)}, \quad i, j = 2, \dots, n,$$

$$b_i^{(2)} = b_i^{(1)} - l_{i1}b_1^{(1)}, \quad i = 2, \dots, n,$$

where the  $b_i^{(1)}$  are the components of  $\mathbf{b}^{(1)}$  and we get a new system of the form

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ \vdots & \vdots & & \vdots \\ 0 & a_{n2}^{(2)} & \dots & a_{nn}^{(2)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(2)} \\ \vdots \\ b_n^{(2)} \end{bmatrix},$$

which will be written as  $A^{(2)}\mathbf{x} = \mathbf{b}^{(2)}$  and that is equivalent to the system we had at the beginning.

Once again we can transform this system by removing the unknown  $x_2$  from the rows  $3, \dots, n$ . By repeating this step we obtain a finite series of systems

$$A^{(k)} \mathbf{x} = \mathbf{b}^{(k)}, \quad 1 \leq k \leq n, \quad (7)$$

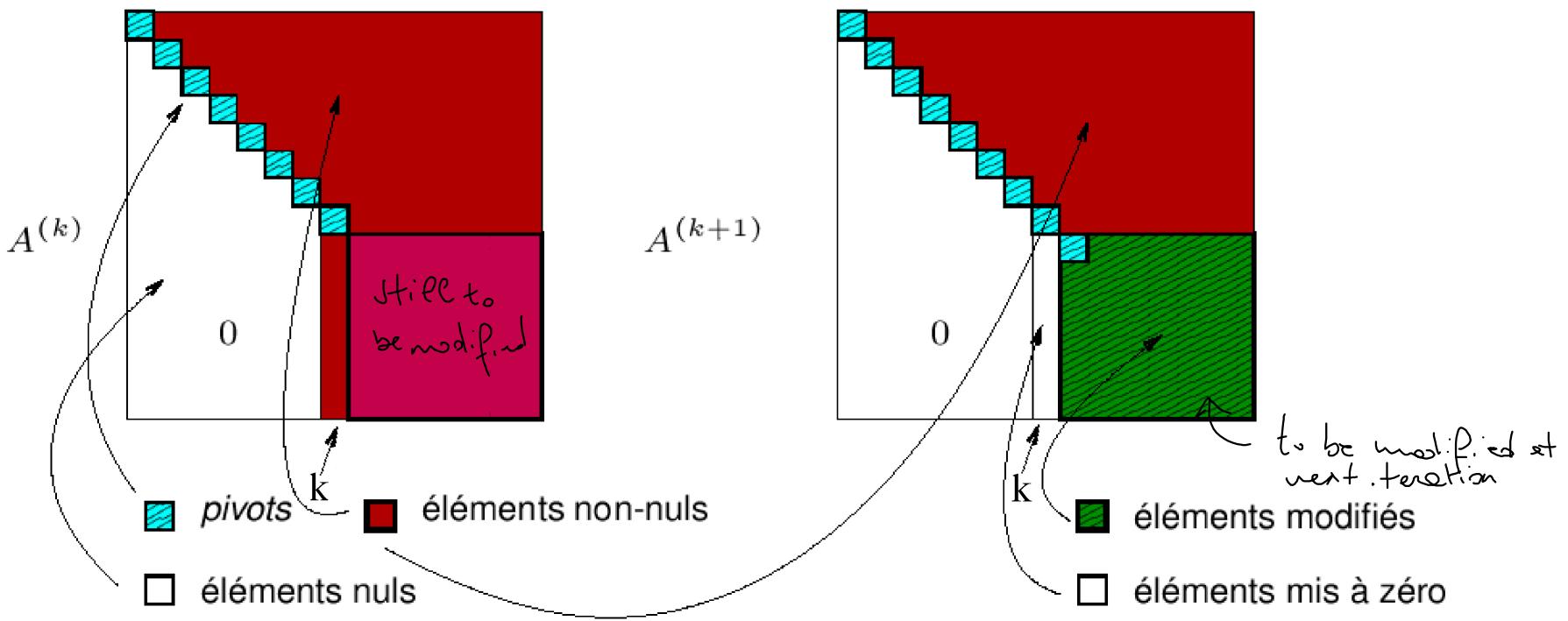
where, for  $k \geq 2$ , the matrix  $A^{(k)}$  is of the form

$$A^{(k)} = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & \dots & \dots & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & & & & a_{2n}^{(2)} \\ \vdots & & \ddots & & & \vdots \\ 0 & \dots & 0 & a_{kk}^{(k)} & \dots & a_{kn}^{(k)} \\ \vdots & & \vdots & \vdots & & \vdots \\ 0 & \dots & 0 & a_{nk}^{(k)} & \dots & a_{nn}^{(k)} \end{bmatrix},$$

where we assume  $a_{ii}^{(i)} \neq 0$  for  $i = 1, \dots, k - 1$ .

↳ otherwise the multiplier gives division by zero

Gauss elimination method: diagram showing how the matrix  $A^{(k+1)}$  is obtained from the matrix  $A^{(k)}$ .



It is clear that for  $k = n$  we obtain the following upper triangular system

$$A^{(n)} \mathbf{x} = \mathbf{b}^{(n)} :$$

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \cdots & \cdots & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & & & a_{2n}^{(2)} \\ \vdots & \ddots & \ddots & & \vdots \\ 0 & & \ddots & \ddots & \vdots \\ 0 & & & & a_{nn}^{(n)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(2)} \\ \vdots \\ \vdots \\ b_n^{(n)} \end{bmatrix}.$$

To be consistent with the previous notation, we write as  $U$  upper triangular matrix  $A^{(n)}$ . The elements  $a_{kk}^{(k)}$  are called *pivots* and have to be non-zero for  $k = 1, \dots, n - 1$ .

In order to make explicit the formulae to get from the  $k$ -th system to the  $k + 1$ -th, for  $k = 1, \dots, n - 1$ , we assume that  $a_{kk}^{(k)} \neq 0$  and we define the multiplier

$$l_{ik} = \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}}, \quad i = k + 1, \dots, n, \quad [(\mathbf{n} - k) \text{ operations}] \quad (8)$$

↓  
 number of multipliers needed

we set then

$$a_{ij}^{(k+1)} = a_{ij}^{(k)} - l_{ik} a_{kj}^{(k)}, \quad i, j = k + 1, \dots, n, \quad [2(n - k)^2 \text{ operations}] \quad (9)$$

$$b_i^{(k+1)} = b_i^{(k)} - l_{ik} b_k^{(k)}, \quad i = k + 1, \dots, n. \quad [2(n - k) \text{ operations}]$$

**Remark 1.** To perform the Gauss elimination,

$$\begin{aligned}
 & \text{final operation doesn't count because we have , problem is just of constant size} \\
 2 \sum_{k=1}^{n-1} \underbrace{(n-k)^2}_p + 3 \sum_{k=1}^{n-1} \underbrace{(n-k)}_p = \\
 2 \sum_{p=1}^{n-1} p^2 + 3 \sum_{p=1}^{n-1} p &= 2 \frac{(n-1)n(2n-1)}{6} + 3 \frac{n(n-1)}{2}
 \end{aligned}$$

operations are required, plus  $n^2$  operations for the resolution with the backward substitutions method of the triangular system  $U\mathbf{x} = \mathbf{b}^{(n)}$ . By keeping only the dominant elements (of order  $n^3$ ), we can say that the Gauss elimination method has a cost of around

$$\frac{2}{3}n^3 \text{ operations.}$$

The following table shows the estimated computation time to solve a system using the LU factorization in different computers:

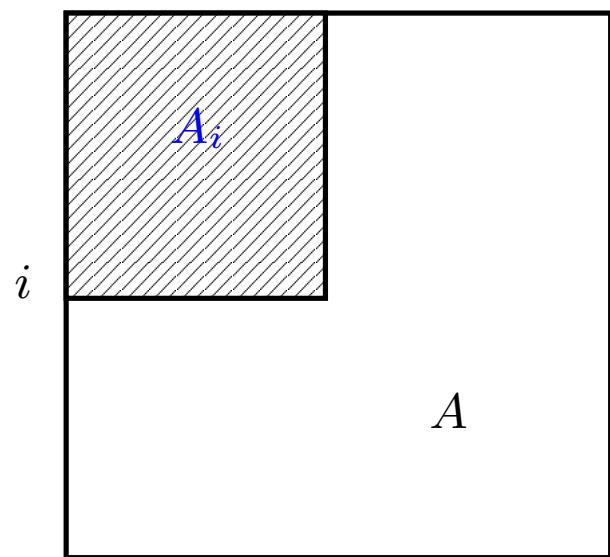
$n$	Number of flops of the computer		
	$10^9$ (Giga)	$10^{12}$ (Tera)	$10^{15}$ (Peta)
$10^2$	$7 \cdot 10^{-4}$ sec	negligible	negligible
$10^4$	11 min	0.7 sec	$7 \cdot 10^{-4}$ sec
$10^6$	21 years	7.7 months	11 min
$10^8$	o.r.	o.r.	21 years

The Gauss method is only properly defined if the pivots  $a_{kk}^{(k)}$  are non-zero for  $k = 1, \dots, n - 1$ . Unfortunately, knowing that the diagonal elements of  $A$  are not zero is not enough to avoid null pivots during the elimination phase. For example, the matrix  $A$  in (10) is invertible and its diagonal elements are non-zero

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 5 \\ 7 & 8 & 9 \end{bmatrix}, \text{ but we find } A^{(2)} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & \boxed{0} & -1 \\ 0 & -6 & -12 \end{bmatrix}. \quad (10)$$

Nevertheless, we have to stop the Gauss method at the second step, because  $a_{22}^{(2)} = 0$ .

Let  $A_i$  be the  $i$ -th main submatrix of  $A$  ( $i = 1, \dots, n - 1$ ), i.e. the submatrix made of the  $i$  first rows and columns of  $A$ :



and let  $d_i$  be the principal minor of  $A$  defined as  $d_i = \det(A_i)$ . We have the following result.

**Proposition 1.** (*Proposition 5.1 in the book*) For a given matrix  $A \in \mathbb{R}^{n \times n}$ , its Gauss factorization exists and is unique iff the principal submatrices  $A_i$  ( $i = 1, \dots, n - 1$ ) are non-singular (i.e. the principal minors  $d_i$  are non-zero:  $d_i \neq 0$ ).

Remark: If  $d_i \neq 0$  ( $i = 1, \dots, n - 1$ ), then the pivots  $a_{ii}^{(i)}$  are also non-zero.

The matrix of the previous example does not satisfy this condition because  $d_1 = 1$  but  $d_2 = 0$ .

There are some categories of matrices for which the hypothesis of the proposition (1) are fulfilled. In particular, we mention:

1. *(Strictly >) diagonal dominant by row* matrices. A matrix  $A$  is said diagonal dominant by row if
 

$|a_{ii}| \geq \sum_{j=1, \dots, n; j \neq i} |a_{ij}|, \quad i = 1, \dots, n.$

↳ for Gauss elimination, to be sure one could ask for strict case
2. *(Strictly >) diagonal dominant by column* matrices. A matrix  $A$  is said diagonal dominant by column if

$$|a_{jj}| \geq \sum_{i=1, \dots, n; i \neq j} |a_{ij}|, \quad j = 1, \dots, n.$$

Examples:  $\begin{bmatrix} -4 & 1 & 2 \\ 2 & 5 & 0 \\ -2 & 1 & 7 \end{bmatrix}$  is diagonal dominant by row and by column, whereas

$\begin{bmatrix} -3 & 1 & 2 \\ 2 & 5 & 0 \\ -2 & 1 & 7 \end{bmatrix}$  is only diagonal dominant by row.

↗ strictly      ↗ not strictly  
 because of  
 I column

3. *Symmetric definite positive* matrices. A matrix  $A$  is symmetric if  $A = A^T$ ; it is definite positive if all its eigenvalues are positive, i.e.:

$$\lambda_i(A) > 0, \quad i = 1, \dots, n.$$

$\geqslant \longrightarrow$  semidefinite positive

# Gauss $\sim LU$

We can show that the Gauss method is equivalent to the factorization  $A = LU$  of the matrix  $A$ , with  $L = \text{multiplier matrix}$  and  $U = A^{(n)}$ .

More exactly:

↖ now we know how to find  $L$

$$A = \underbrace{\begin{bmatrix} 1 & 0 & \dots & \dots & 0 \\ l_{21} & 1 & & & 0 \\ \vdots & l_{32} & \ddots & & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ l_{n1} & & & l_{n,n-1} & 1 \end{bmatrix}}_L \underbrace{\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & \dots & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & & & a_{2n}^{(2)} \\ \vdots & & \ddots & \ddots & \vdots \\ 0 & & & \ddots & \vdots \\ 0 & & & & a_{nn}^{(n)} \end{bmatrix}}_U.$$

The matrices  $L$  and  $U$  only depend on  $A$  (and not on  $\mathbf{b}$ ), the same factorization can be resused for solving several linear systems that share the same matrix  $A$  but **different vectors  $\mathbf{b}$** .

The number of operations is then considerably reduced, since most of the computational weight, around  $\frac{2}{3}n^3$  flops, is due to the Gaussian elimination process. Indeed, let us consider the  $M$  linear systems:

$$A\mathbf{x}_m = \mathbf{b}_m \quad m = 1, \dots, M$$

Donc:

- the cost of the factorization  $A = LU$  is  $\frac{2}{3}n^3$  flops;
- the cost of the resolution of both triangular systems,  $L\mathbf{y}_m = \mathbf{b}_m$  and  $U\mathbf{x}_m = \mathbf{y}_m$  ( $m = 1, \dots, M$ ) is  $2Mn^2$  flops,

for a total of  $\frac{2}{3}n^3 + 2Mn^2$  flops which is much smaller than  $\frac{2}{3}Mn^3$  flops required to solve all the systems Gauss elimination method.

# The pivoting technique

It has been already noted that the Gauss method fails if a pivot becomes zero. In that case, we can use a technique called **pivoting** that consists in exchanging the rows (or the columns) of the system in such a way that no pivot is zero.

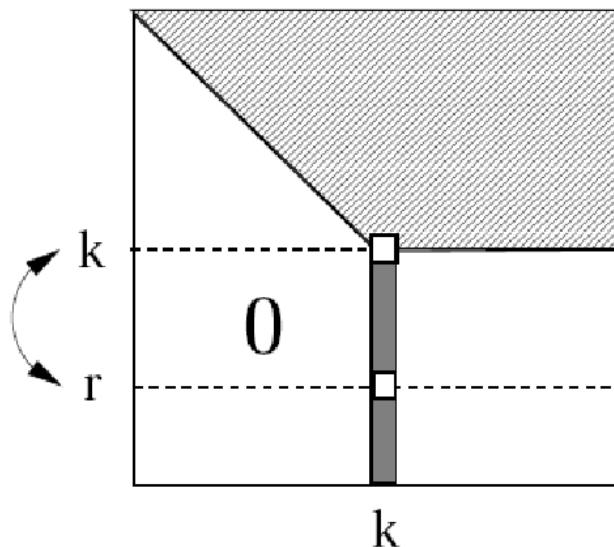
**Example 5.** Let us go back to the matrix (10) for which the Gauss method gives a null pivot at the second step. By just exchanging the second and the third rows, we get a non-zero pivot and can execute one step further. Indeed,

$$A^{(2)} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & \boxed{0} & -1 \\ 0 & -6 & -12 \end{bmatrix} \implies P_2 A^{(2)} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & \boxed{-6} & -12 \\ 0 & 0 & -1 \end{bmatrix},$$

where  $P_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$  is called **permutation** matrix.

↑ if you change rows  
 $\Rightarrow P_2 A : \text{PREMULT}$   
 ↑ if columns  
 $\Rightarrow A P_2 : \text{POSTMULTPLIC}$

The pivoting strategy used for the example 5 can be generalized by finding, at every step  $k$  of the elimination, a non-zero pivot among the elements of the subcolumn  $A^{(k)}(k : n, k)$ . This is called a *partial pivot change* (by row).

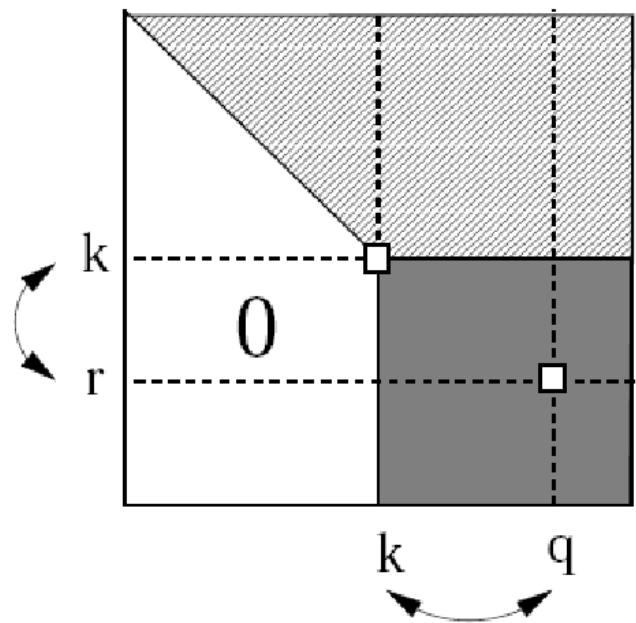


! Avoid big multipliers  
 $\Rightarrow$  avoid division by values close to zero  
 (you could have AMPLIFICATION OF ROUNDING ERRORS)

From (8) we know that a big value of  $l_{ik}$  (coming for instance from a small  $a_{kk}^{(k)}$ ) can amplify the rounding errors affecting the elements  $a_{kj}^{(k)}$ .

Consequently, in order to ensure a better stability, we choose as pivot the biggest element in module of the column  $A^{(k)}(k : n, k)$ , and the partial pivoting is performed at every step, even if it is not strictly necessary (i.e. even if the pivot is non-zero).

An alternative method consists looking for the pivot in the whole submatrix  $A^{(k)}(k : n, k : n)$ , performing what is called *complete pivoting*.



Remark that, whereas partial pivoting requires just an additional cost of  $n^2$  tests, complete pivoting needs some  $2n^3/3$ , what considerably increases the cost of the Gauss method.

*full permutation  $\Rightarrow$  doubling of cost*

In general, if at the step  $k$  we have to exchange the rows  $k$  and  $r$ , we will have to multiply  $A^{(k)}$  by the following permutation matrix  $P_k$  before continuing:

$$\begin{array}{c}
 k \rightarrow \\
 r \rightarrow
 \end{array}
 \left(
 \begin{array}{cccccc}
 1 & & & & & \\
 & \ddots & & & & \\
 & & 0 & \dots & 1 & \\
 & & & \vdots & & \\
 & & 1 & \dots & 0 & \\
 & & & & \ddots & \\
 & & & & & 1
 \end{array}
 \right) = P_k$$

↗ to use on  
 both sides  
 of the system!

This means we will consider  $P_k A^{(k)}$  instead of  $A^{(k)}$ .

We can prove that the result obtained is of the form:

$$PA = LU, \quad (11)$$

being  $P = P_{n-1}P_{n-2}\dots P_2P_1$  (global permutation matrix).  $L$  is the multiplier matrix (the new ones!) and  $U = A^{(n)}$ .

Once the matrices  $L$ ,  $U$  and  $P$  have been calculated, the resolution of the initial system is transformed into the resolution of the triangular systems

$$A\mathbf{x} = \mathbf{b} \implies PA\mathbf{x} = P\mathbf{b} \implies \begin{cases} Ly = P\mathbf{b}, \\ U\mathbf{x} = \mathbf{y}. \end{cases}$$

Remark that the coefficients of the matrix  $L$  have the same values as the multipliers calculated by a factorization  $LU$  of the matrix  $PA$  without pivoting.

If we use complete pivoting, it can be proved that we arrive to the following result:

$$PAQ = LU$$

where  $P = P_{n-1} \dots P_1$  is a permutation matrix that takes into account all permutations by row, and  $Q = Q_1 \dots Q_{n-1}$  is a permutation matrix that takes into account all permutations by column. By construction, the matrix  $L$  is still lower triangular, and its elements have a module lower or equal to 1.

As for the partial pivoting, the elements of  $L$  are the multipliers generated by the factorization  $LU$  of the matrix  $PAQ$  with no pivoting.

Once the matrices  $L$ ,  $U$ ,  $P$  and  $Q$  have been calculated, for solving the linear system we notice that we can write

$$A\mathbf{x} = \mathbf{b} \iff \underbrace{PAQ}_{LU} \underbrace{Q^{-1}\mathbf{x}}_{\mathbf{x}^*} = P\mathbf{b}.$$

What brings us to the resolution of triangular systems

$$\begin{cases} L\mathbf{y} = P\mathbf{b}, \\ U\mathbf{x}^* = \mathbf{y}. \end{cases}$$

and finally we compute

$$\mathbf{x} = Q\mathbf{x}^*.$$

**Remark 2.** In Matlab/Octave, we can get the factorization of a matrix  $A$  with the command

```
>> [L,U,P] = lu(A);
```

The matrix  $P$  is a permutation matrix. In the case where the matrix  $P$  is the identity, the matrices  $L$  and  $U$  are the matrices we are looking for (such that  $LU = A$ ). Otherwise, we have  $LU = PA$ .

*P is not needed*

**Example. 2 (continued)** In Matlab/Octave, we first need to define the matrix  $A$  and the vector  $b$  of the linear system:

```
>> A = [-0.37, 0.05, 0.05, 0.07; 0.05, -0.116, 0, 0.05; ...
          0.05, 0, -0.116, 0.05; 0.07, 0.05, 0.05, -0.202];
>> b = [-2; 0; 0; 0];
```

Then, we can use the command `\` as follows:

```
>> x = A\b
x =
    8.1172
    5.9893
    5.9893
    5.7779
```

This command computes the solution of the system with *ad hoc* algorithms (it tests the matrix to choose an optimal algorithm).

Remark: the pressures at the nodes 2 and 3 are the same (by symmetry).

If we wanted to use the  $LU$  factorization:

```
>> [L,U,P] = lu(A);  
>> y = L\ (P*b);  
>> x = U\y  
x =  
8.1172  
5.9893  
5.9893  
5.7779
```

The solution is the same. We can verify that, in this case, no permutation takes place ( $P$  is the identity matrix):

```
>> P  
P =  
1 0 0 0  
0 1 0 0  
0 0 1 0  
0 0 0 1
```

LU fact gives you the possibility to store LU in smarter ways

**Remark 3.** Using the LU factorization, obtained by fixing the value 1 for the  $n$  diagonal elements of  $L$ , we can calculate the determinant of a square matrix with  $O(n^3)$  operations, because

$$\det(A) \stackrel{\text{bunct theorem}}{=} \det(L) \det(U) = \det(U) = \prod_{k=1}^n u_{kk};$$

L has diagonal elements = 1  
 n operation

indeed, the determinant of a triangular matrix is the product of the diagonal elements. The Matlab/Octave command `det(A)` makes use of this.

# The inverse matrix

If  $A$  is a  $n \times n$  **non-singular** matrix, let us call  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$  the columns of its inverse matrix  $A^{-1}$  i.e.  $A^{-1} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)})$ .

The relation  $AA^{-1} = I$  can be expressed by the following  **$n$  systems** : for  $1 \leq k \leq n$ ,

$$A\mathbf{x}^{(k)} = \mathbf{e}^{(k)}, \quad \text{← } \begin{matrix} \text{\scriptsize $k$-th column} \\ \text{\scriptsize of identity} \end{matrix} \quad (12)$$

where  $\mathbf{e}^{(k)}$  is the column vector with all the elements equal to 0 except the one corresponding to the  $k$ -th row, which equals 1.

Once we know the matrices  $L$  and  $U$  that factorizes  $A$ , solving the  $n$  systems (12) defined by the same matrix  $A$  requires  $2n^3$  operations.

# The Cholesky factorization

In the case where the  $n \times n$  matrix  $A$  is symmetric and positive definite, there exists a unique upper triangular matrix  $R$  with positive diagonal elements such that

$$A = R^T R.$$

This factorization is called *Cholesky factorization*. In Matlab/Octave, the command

```
>> R = chol(A)
```

can be used to compute  $R$ .

↴  
 requires only  $1/2$   
 memory of LU  
 Factorization  
 (not 2 matrices)  
 but only 1  
 (matrix)  
 ↑  
 upper triangular matrix  
 and symmetric

The elements  $r_{ij}$  of  $R$  can be calculated using the expressions  $r_{11} = \sqrt{a_{11}}$  and for  $i = 2, \dots, n$ :

$$r_{ji} = \frac{1}{r_{jj}} \left( a_{ij} - \sum_{k=1}^{j-1} r_{ki} r_{kj} \right), \quad j = 1, \dots, i-1, \quad (13)$$

$$r_{ii} = \sqrt{a_{ii} - \sum_{k=1}^{i-1} r_{ki}^2} \quad (14)$$

The Cholesky factorization needs around  $\frac{n^3}{3}$  operations (half the operations for a LU factorization).

**Example 6.** In Matlab there are several families of predefined matrices (see `help gallery`). Let us consider the family of the *Lehmer* matrices, that are positive definite and thus candidates for a Cholesky factorization. We want to calculate the cost of the Cholesky factorization for matrices of size  $n = 10, 20, 30, 40$  and  $50$ .

*Remark :* In the version 5 of Matlab there was available the command `flops` which allowed to calculate the number of operations executed. Unfortunately, in the new version 6 of Matlab, this command has been removed (due to the integration in Matlab of certain linear algebra libraries). This example can't be executed on Matlab 6:

```
>> fl=[];
>> for n=10:10:50
    A=gallery('lehmer',n);
    flops(0)
    R=chol(A);
    fl = [fl, flops];
end
```

We get the number of *flops* for the different sizes

```
>> fl
fl =
      385      2870      9455     22140     42925
```

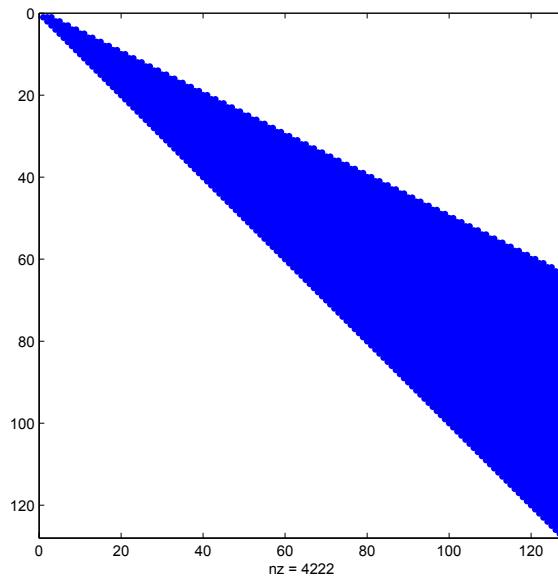
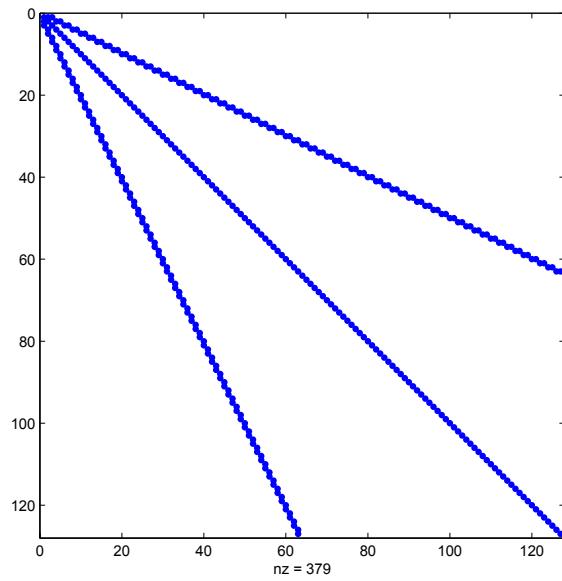
that grow proportionally to  $n^3$ . Indeed, if we calculate the ratio  $fl(n)/n^3$  we find a “quasi” constant value.

```
>> n=[10, 20, 30, 40, 50];
>> fl./(n.^3)
ans =
    0.3850    0.3588    0.3502    0.3459    0.3434
```

Remark that this value is approximately  $\frac{1}{3}$ . Thus, we can say that the computational cost of the Cholesky factorization is of order  $\frac{1}{3}n^3$ , or half the cost of the Gauss elimination method.

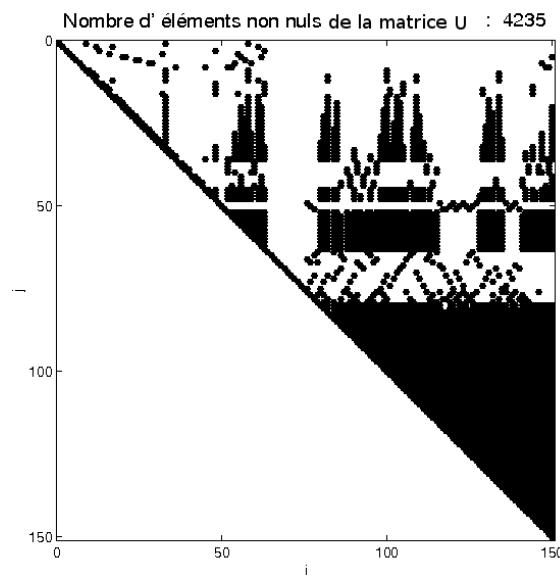
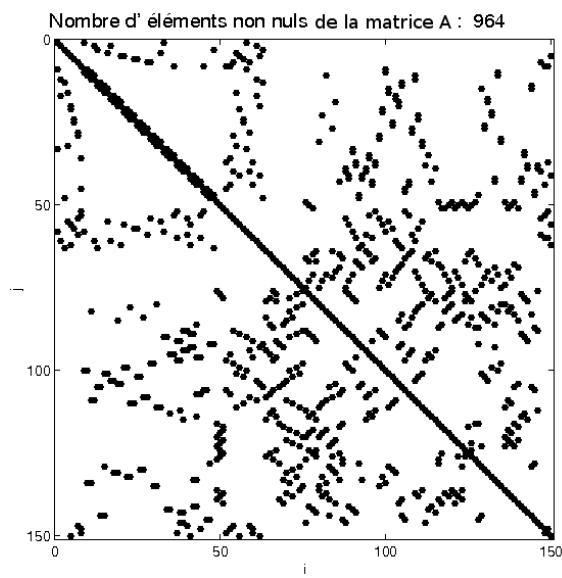
# Memory space limitations

**Example.** 1 (continued) Let  $A$  be a matrix of size  $127 \times 127$  corresponding to capillary bed with 8 bifurcation levels. This matrix is symmetric definite positive. The number of non-null entries of  $A$  is 379 and thus much smaller than  $(127)^2 = 16129$ . It is a *sparse* matrix. The figure on the left shows the disposition of the non-null entries of  $A$ , whereas the one on the right shows the non-null entries of the matrix  $R$ .



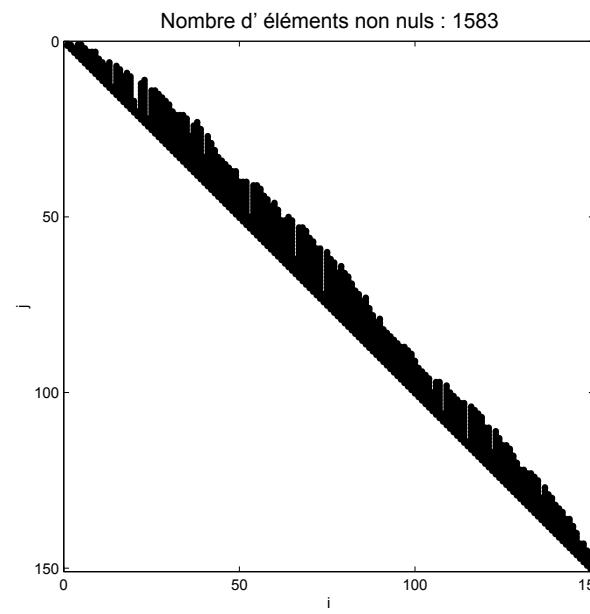
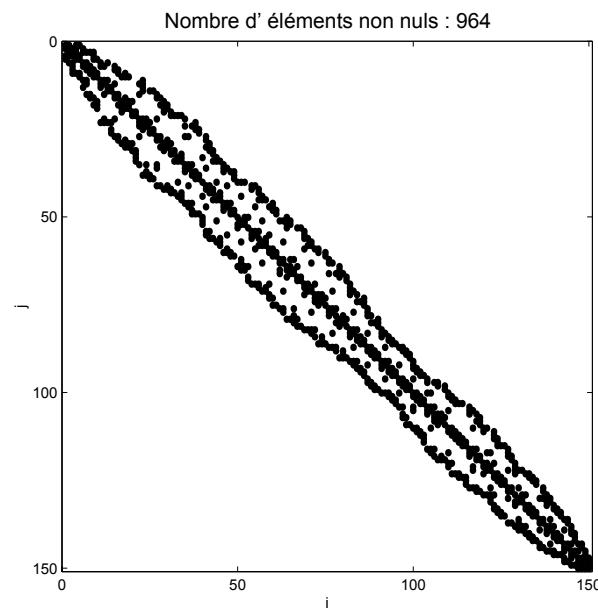
**Example 7.** Let us consider the problem of calculating the deformations in a structure subject to a given set of forces. The discretization using the finite elements method generates a matrix  $A$  of size  $150 \times 150$ . (The same matrix would have been produced by the approximation of an electric potential field.) This matrix is symmetric definite positive. The number of non-null entries of  $A$  is 964, and thus much smaller than  $(150)^2 = 22500$ . It is a *sparse* matrix.

The figure on the left shows the disposition of the non-null entries of  $A$ , whereas the one on the right shows the non-null entries of the matrix  $R$ .



We notice that the number of non-null entries of  $R$  is much bigger than those of  $A$  (*fill-in phenomenon*). This leads to a bigger memory usage. To reduce the fill-in phenomenon, we can re-order rows and columns of  $A$  in a particular fashion; this is called *re-ordering* of the matrix. There are several algorithms that allow us to do this (Matlab command: `symand`).

For example, the following figure shows, on the left, one possibility of reordering  $A$ , while the one on the right shows the disposition of the non-null entries of the Cholesky factorization of the reordered matrix  $A$ .



concentrating  
 elements  
 around diagonals  
 gives stability  
to be defined

# Precision limitations

**Example 8.** Rounding errors can induce important differences between the calculated solution using the Gauss elimination method (GEM) and the exact solution. This happens when the *conditioning* of the matrix of the system is very big.

The Hilbert matrix of size  $n \times n$  is a symmetric matrix defined by:

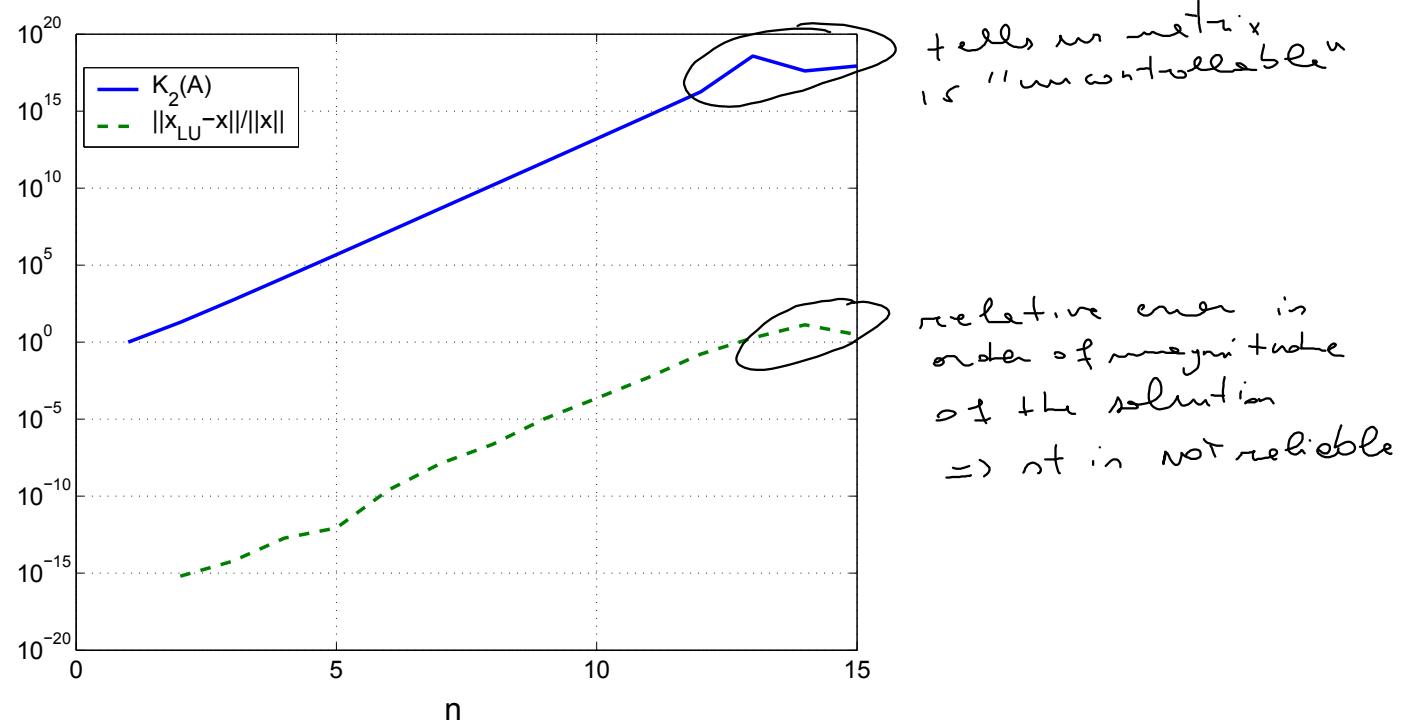
$$A_{ij} = \frac{1}{i+j-1}, \quad i, j = 1, \dots, n$$

In Matlab/Octave, we can build a Hilbert matrix of any size  $n$  with the command `A = hilb(n)`. For example, for  $n = 4$ , we get:

$$A = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} \\ \frac{1}{4} & \frac{1}{5} & \frac{1}{6} & \frac{1}{7} \end{bmatrix}$$

We consider the linear systems  $A_n \mathbf{x}_n = \mathbf{b}_n$  where  $A_n$  is the Hilbert matrix of size  $n$  with  $n = 4, 6, 8, 10, 12, 14, \dots$ , whereas  $\mathbf{b}_n$  is chosen such that the exact solution is  $\mathbf{x}_n = (1, 1, \dots, 1)^T$ .

For every  $n$ , we calculate the conditioning of the matrix, we solve the linear system by  $LU$  factorization and we get  $\mathbf{x}_n^{LU}$  as the found solution. The obtained conditioning as well as the error  $\|\mathbf{x}_n - \mathbf{x}_n^{LU}\|/\|\mathbf{x}_n\|$  (where  $\|\cdot\|$  is the euclidian norm of a vector,  $\|\mathbf{x}\| = \sqrt{\mathbf{x}^T \cdot \mathbf{x}}$ ) are shown in the figure below.



# Considerations on the precision (Sect. 5.5)

The methods we have seen until now allow us to find the solution of a linear system in a finite number of operations. That is why they are called *direct methods*. However, there are cases where these methods are not satisfactory.

**Definition 1.** We call *conditioning* of a matrix  $M$ , symmetric definite positive, the ratio between the maximum and minimum of its eigenvalues, i.e.

$$K(M) = \frac{\lambda_{\max}(M)}{\lambda_{\min}(M)}$$

If  $K$  is too big  $\Rightarrow$  got stability behavior

It can be shown that, the bigger the conditioning of a matrix, the worse the solution obtained by a direct method.

For example, let us consider a linear system  $Ax = b$ .

If we solve this system with a computer, due to rounding errors, we will not find the exact solution  $x$  but an approximate solution  $\hat{x}$ . The following relationship can be shown :

$$\frac{\|x - \hat{x}\|}{\|x\|} \leq K(A) \frac{\|r\|}{\|b\|} \quad (15)$$

needs to  
 be computed  
 because x is  
 generally not  
 known

where  $r$  is the residual  $r = b - A\hat{x}$ ; we write as  $\|\mathbf{v}\| = (\sum_{k=1}^n v_k^2)^{1/2}$  the Euclidean norm of a vector  $\mathbf{v}$ .

Remark that, if the conditioning of  $A$  is big, the distance  $\|x - \hat{x}\|$  between the exact solution and the numerically computed solution can be very big even if the residual is very small.

**Proof for (15) :** Let  $A$  be a symmetric definite positive matrix, we can consider the  $n$  eigenvalues  $\lambda_i > 0$  and the associated unitary eigenvectors  $\{\mathbf{v}_i\}$ ,  $i = 1, \dots, n$ :  $A\mathbf{v}_i = \lambda_i \mathbf{v}_i$ ,  $i = 1, \dots, n$ . These vectors form an orthonormal base of  $\mathbb{R}^n$ , what means  $\mathbf{v}_i^T \mathbf{v}_j = \delta_{ij}$  for  $i, j = 1, \dots, n$ . For any  $\mathbf{w} \in \mathbb{R}^n$ , if we write it as

$$\mathbf{w} = \sum_{i=1}^n w_i \mathbf{v}_i,$$

we have

$$\begin{aligned}
 \|A\mathbf{w}\|^2 &= (A\mathbf{w})^T (A\mathbf{w}) \\
 &\stackrel{\text{Av}_i \perp \mathbf{v}_i}{=} (\lambda_1 w_1 \mathbf{v}_1^T + \dots + \lambda_n w_n \mathbf{v}_n^T) (\lambda_1 w_1 \mathbf{v}_1 + \dots + \lambda_n w_n \mathbf{v}_n) \\
 &= \sum_{i,j=1}^n \lambda_i \lambda_j w_i w_j \mathbf{v}_i^T \mathbf{v}_j = \sum_{i,j=1}^n \lambda_i \lambda_j w_i w_j \delta_{ij} = \sum_{i=1}^n \lambda_i^2 w_i^2.
 \end{aligned}$$

And yet, as  $\|\mathbf{w}\|^2 = \sum_{i=1}^n w_i^2$ , we get  $\|A\mathbf{w}\|^2 \leq \lambda_{max}^2 \|\mathbf{w}\|^2$ , i.e.  
 $\|A\mathbf{w}\| \leq \lambda_{max} \|\mathbf{w}\|$  where  $\lambda_{max}$  is the biggest eigenvalue of  $A$ .

As the eigenvalues of  $A^{-1}$  are  $1/\lambda_i$ , we also get

$\|A^{-1}\mathbf{w}\| \leq \frac{1}{\lambda_{min}} \|\mathbf{w}\| \quad \forall \mathbf{w} \in \mathbb{R}^n$ , where  $\lambda_{min}$  is the smallest eigenvalue of  $A$ .

Thus, we have

$$\begin{aligned} \|\mathbf{x} - \hat{\mathbf{x}}\| &= \|A^{-1}\mathbf{r}\| \leq \frac{1}{\lambda_{min}} \|\mathbf{r}\|, \\ \|\mathbf{b}\| &= \|A\mathbf{x}\| \leq \lambda_{max} \|\mathbf{x}\|, \end{aligned}$$

$\mathbf{r} = \mathbf{b} - A\hat{\mathbf{x}} = A\mathbf{x} - A\hat{\mathbf{x}} = A(\mathbf{x} - \hat{\mathbf{x}})$   
 $\hookrightarrow A^{-1}\mathbf{r} = \mathbf{x} - \hat{\mathbf{x}}$

from where we directly find the inequality (15). ■

# Linear Systems - Iterative Methods



## Numerical Analysis

Profs. Gianluigi Rozza-Luca Heltai

2019-SISSA mathLab Trieste

(Sec. 5.9 of the book)

Solve the linear system  $Ax = b$  using an iterative method consists in building a series of vectors  $x^{(k)}$ ,  $k \geq 0$ , in  $\mathbb{R}^n$  that converge at the exact solution  $x$ , i.e.:

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

for any initial vector  $x^{(0)} \in \mathbb{R}^n$ .

We can consider the following recurrence relation:

$$x^{(k+1)} = Bx^{(k)} + g, \quad k \geq 0$$

↙ transformation matrix + let allow  $x^{(k)} \rightarrow x^{(k+1)}$   
 ↘ needed for respecting consistency of (1)

where  $B$  is a well chosen matrix (depending on  $A$ ) and  $g$  is a vector (that depends on  $A$  and  $b$ ), satisfying the relation (of consistence)

$$x = Bx + g. \quad (2)$$

↙ if doesn't hold,  
 our method is NOT  
 CONSISTENT

some Spectral properties needed

Given  $\mathbf{x} = A^{-1}\mathbf{b}$ , we get  $\mathbf{g} = (I - B)A^{-1}\mathbf{b}$ ; the iterative method is therefore completely defined by the matrix  $B$ , known as *iteration matrix*.  
By defining the error at step  $k$  as

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

*error for each iteration  
(2) - (12)*

we obtain the following recurrence relation:

$$\mathbf{e}^{(k+1)} = B\mathbf{e}^{(k)} \quad \text{and thus} \quad \mathbf{e}^{(k+1)} = B^{k+1}\mathbf{e}^{(0)}, \quad k = 0, 1, \dots$$

We can show that  $\lim_{k \rightarrow \infty} \mathbf{e}^{(k)} = \mathbf{0}$  for all  $\mathbf{e}^{(0)}$  (and thus for all  $\mathbf{x}^{(0)}$ ) if and only if

$$\rho(B) < 1,$$

or  $\rho(B)$  is the *spectral radius* of the matrix  $B$ , defined by

$$\rho(B) = \max |\lambda_i(B)|$$

and  $\lambda_i(B)$  are the eigenvalues of the matrix  $B$ .

The smaller the value of  $\rho(B)$ , the less iterations are needed to reduce the initial error of a given factor.

# Construction of an iterative method

A general way of setting up an iterative method is based on the decomposition of the matrix  $A$ :

$$A = P - (P - A)$$

where  $P$  is an invertible matrix called *preconditioner* of  $A$ .

Hence,

$$Ax = \mathbf{b} \Leftrightarrow Px = (P - A)x + \mathbf{b}$$

which is of the form (2) leaving

$$\begin{matrix} \downarrow \\ x^{(k+1)} \end{matrix} \quad \begin{matrix} \downarrow \\ x^{(k)} \end{matrix}$$

*as far the due to the splitting we basically rewrite 1, but that can help to better face the problem*

$$B = P^{-1}(P - A) = I - P^{-1}A \quad \text{and} \quad \mathbf{g} = P^{-1}\mathbf{b}.$$

C! Condition number is a measure of stability

Preconditioner helps for stability

We can define the corresponding iterative method

$$P(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = \mathbf{r}^{(k)} \quad k \geq 0$$

*brace* substitute to believe!

where  $\mathbf{r}^{(k)}$  represents the *residual* at the iteration  $k$ :  $\boxed{\mathbf{r}^{(k)} = \mathbf{b} - A\mathbf{x}^{(k)}}$

We can generalise this method as follows:

*Nasty if  $b$  and  $Ax^{(k)}$  are close at machine precision*

$$P(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = \alpha_k \mathbf{r}^{(k)} \quad k \geq 0 \quad (3)$$

*brace* Helper in convergence

where  $\alpha_k \neq 0$  is a parameter that improves the convergence of the series  $\mathbf{x}^{(k)}$ .  
 The method (3) is called *Richardson's method*.

The matrix  $P$  has to be chosen in such a way that renders the cost of solving (3) small enough. For example a diagonal or triangular  $P$  matrix would comply with this criterion.

# Jacobi method

If the elements of the diagonal of  $A$  are non-zero, we can write

$$P = D = \text{diag}(a_{11}, a_{22}, \dots, a_{nn})$$

$D$  with the diagonal part of  $A$  being:

$$D_{ij} = \begin{cases} 0 & \text{si } i \neq j \\ a_{ij} & \text{if } i = j. \end{cases}$$

The Jacobi method corresponds to this choice with  $\alpha_k = 1$  for all  $k$ .

We deduce:

$$D\mathbf{x}^{(k+1)} = \mathbf{b} - (A - D)\mathbf{x}^{(k)} \quad k \geq 0.$$

By components:

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left( b_i - \sum_{j=1, j \neq i}^n a_{ij} x_j^{(k)} \right), \quad i = 1, \dots, n. \quad (4)$$

NO DIAGONAL!

The Jacobi method can be written under the general form

$$\mathbf{x}^{(k+1)} = B\mathbf{x}^{(k)} + \mathbf{g},$$

with

$$B = B_J = D^{-1}(D - A) = I - D^{-1}A, \quad \mathbf{g} = \mathbf{g}_J = D^{-1}\mathbf{b}.$$

# Gauss-Seidel method

This method is defined as follows:

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

*Take consideration else of partial yet at  
current iteration it will already computed*

This method corresponds to (1) with  $P = D - E$  and  $\alpha_k = 1$  ( $\forall k \geq 0$ ) where  $E$  is the lower triangular matrix

↑      ↗  
 diagonal      lower triangular

$$\begin{cases} E_{ij} = -a_{ij} & \text{if } i > j \\ E_{ij} = 0 & \text{if } i \leq j \end{cases}$$

(lower triangular part of  $A$  without the diagonal and with its elements' sign inverted).

We can write this method under the form (3), with the iteration matrix  $B = B_{GS}$  given by

$$B_{GS} = (D - E)^{-1}(D - E - A)$$

and

$$\mathbf{g}_{GS} = (D - E)^{-1}\mathbf{b}.$$



why not only  $\varepsilon$ ?  
 Library (historically)  
 builds ~~is~~ matrix by block,  
 using all the diagonal

**Example 1.** Given the matrix

$$A = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{pmatrix}.$$

We have then

$$D = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 6 & 0 & 0 \\ 0 & 0 & 11 & 0 \\ 0 & 0 & 0 & 16 \end{pmatrix};$$

Thus, the iteration matrix for the Jacobi method is

Not symmetric

$$B_J = D^{-1}(D - A) = I - D^{-1}A = \begin{pmatrix} 0 & -2 & -3 & -4 \\ -5/6 & 0 & -7/6 & -4/3 \\ -9/11 & -10/11 & 0 & -12/11 \\ -13/16 & -14/16 & -15/16 & 0 \end{pmatrix}.$$

For defining the matrix  $A$  and extracting its diagonal  $D$  and its lower triangular part  $E$  (without the diagonal and the sign inverted) with Matlab/Octave, we use the commands

```
>> A = [1,2,3,4;5,6,7,8;9,10,11,12;13,14,15,16];
>> D = diag(diag(A));
>> E = - tril(A,-1);
```

diag twice = put diag(A) into a diagonal  
 matrix

extracted diagonal

These allow us, for example, to compute the iteration matrix  $B_{GS}$  for the Gauss-Seidel method in the following way:

```
>> B_GS = (D-E)\(D-E-A);
```

We find:

$$B_{GS} = \begin{pmatrix} 0.0000 & -2.0000 & -3.0000 & -4.000 \\ 0.0000 & 1.6667 & 1.3333 & 2.0000 \\ 0.0000 & 0.1212 & 1.2424 & 0.3636 \\ 0.0000 & 0.0530 & 0.1061 & 1.1591 \end{pmatrix}.$$

# Convergence

We have the following convergence results:

- (Prop 5.3) If the matrix  $A$  is strictly diagonally dominant by row, i.e.,

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

then the Jacobi and the Gauss-Seidel methods converge.

- If  $A$  is **symmetric positive definite**, then the Gauss-Seidel method converges (Jacobi maybe not). *Non zero element only on the 3 biggest "diagonals"*
- (Prop 5.4) Let  $A$  be a tridiagonal non-singular matrix whose diagonal elements are all non-null. Then the Jacobi and the Gauss-Seidel methods are either **both divergent or both convergent**. In the latter case,  $\rho(B_{GS}) = \rho(B_J)^2$ .

↳ For Jacobi and Gauss-Seidel  
 $\alpha_k = 1$  ( $\Rightarrow$  we are NOT accelerating)

# Richardson method

(Sec. 5.10)

Let consider the following iterative method:

$$P(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = \alpha_k \mathbf{r}^{(k)}, \quad k \geq 0. \quad (5)$$

If  $\alpha_k = \alpha$  (a constant) this method is called **stationary preconditioned Richardson method**; otherwise **dynamic preconditioned Richardson method** when  $\alpha_k$  varies during the iterations.

The matrix  $P$  is called **preconditioner** of  $A$ .

If  $A$  and  $P$  are **symmetric positive definite**, then there are two optimal criteria to choose  $\alpha_k$ :

1. *Stationary case:*

$$\alpha_k = \alpha_{opt} = \frac{2}{\lambda_{min} + \lambda_{max}}, \quad k \geq 0,$$

← from algebraic  
and geometrical  
properties

where  $\lambda_{min}$  and  $\lambda_{max}$  represent the smaller and the larger eigenvalue of the matrix  $P^{-1}A$ .

2. *Dynamic case:*

$$\alpha_k = \frac{(\mathbf{z}^{(k)})^T \mathbf{r}^{(k)}}{(\mathbf{z}^{(k)})^T A \mathbf{z}^{(k)}}, \quad k \geq 0,$$

This is a lin system

where  $\underbrace{\mathbf{z}^{(k)} = P^{-1}\mathbf{r}^{(k)}}$  is the preconditioned residual. ← it is a TRF of the residual vector  
 This method is also called **preconditioned gradient method**.

If  $P = I$  and  $A$  is symmetric definite positive, we get the following methods:

- the **Stationary Richardson** if we choose:

$$\alpha_k = \alpha_{opt} = \frac{2}{\lambda_{min}(A) + \lambda_{max}(A)}. \quad (6)$$

$P^{-1}A = A$

- the **Gradient** method if :

$$\alpha_k = \frac{(\mathbf{r}^{(k)})^T \mathbf{r}^{(k)}}{(\mathbf{r}^{(k)})^T A \mathbf{r}^{(k)}}, \quad k \geq 0. \quad (7)$$

$P^{-1} r = r$

The gradient method can be written as:

Let  $\mathbf{x}^{(0)}$  be given, set  $\mathbf{r}^{(0)} = \mathbf{b} - A\mathbf{x}^{(0)}$ , then for  $k \geq 0$ ,

I residual

$$\begin{aligned}
 P\mathbf{z}^{(k)} &= \mathbf{r}^{(k)} \\
 \alpha_k &= \frac{(\mathbf{z}^{(k)})^T \mathbf{r}^{(k)}}{(\mathbf{z}^{(k)})^T A \mathbf{z}^{(k)}} \\
 \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \alpha_k \mathbf{z}^{(k)} \\
 \mathbf{r}^{(k+1)} &= \mathbf{r}^{(k)} - \alpha_k A \mathbf{z}^{(k)}.
 \end{aligned}$$

gradient method  
 short vec direction  
 correction (only  
 fixd direction along  
 the gradient); only  
 mod. filtration are on the  
 amplification of  
 parameter

We have to apply once  $A$  and inverse  $P$  at each iteration.  $P$  should then be such that the resolution of the associated system results easy (i.e. it requires a reasonable amount of computing cost). For example, we can choose a diagonal  $P$  (Like in the gradient or the stationary Richardson cases) or triangular.

# Convergence of Richardson method

When  $A$  and  $P$  are s.p.d. and with the two optimal choices for  $\alpha$ , we can show that the preconditioned Richardson Method converges to  $x$  when  $k \rightarrow \infty$ , and that

$$\|\mathbf{x}^{(k)} - \mathbf{x}\|_A \leq \left( \frac{K(P^{-1}A) - 1}{K(P^{-1}A) + 1} \right)^k \|\mathbf{x}^{(0)} - \mathbf{x}\|_A, \quad k \geq 0, \quad (8)$$

$\underbrace{\phantom{\left( \frac{K(P^{-1}A) - 1}{K(P^{-1}A) + 1} \right)^k}}_{< 1}$

where  $\|\mathbf{v}\|_A = \sqrt{\mathbf{v}^T A \mathbf{v}}$  and  $K(P^{-1}A)$  is the condition number of  $P^{-1}A$ .

*we can never  
see that error  
is converging*

**Remark** If  $A$  et  $P$  are s.p.d., we have that

$$K(P^{-1}A) = \frac{\lambda_{max}}{\lambda_{min}}.$$

ENERGY NORM  
associated with  
MATRIX A

**Demonstration** The iteration matrix of the method is given by

$R_\alpha = I - \alpha P^{-1} A$ , where the eigenvalues of  $R_\alpha$  are of the form  $1 - \alpha \lambda_i$ . The method is convergent if and only if  $|1 - \alpha \lambda_i| < 1$  for  $i = 1, \dots, n$ , therefore  $-1 < 1 - \alpha \lambda_i < 1$  for  $i = 1, \dots, n$ . As  $\alpha > 0$ , this is the equivalent to  $-1 < 1 - \alpha \lambda_{max}$ , from where the necessary and sufficient condition for convergence remains  $\alpha < 2/\lambda_{max}$ . Consequently,  $\rho(R_\alpha)$  is minimal if  $1 - \alpha \lambda_{min} = \alpha \lambda_{max} - 1$ , i.e., for  $\alpha_{opt} = 2/(\lambda_{min} + \lambda_{max})$ . By substitution, we obtain

$$\rho_{opt} = \rho(R_{opt}) = 1 - \alpha_{opt} \lambda_{min} = 1 - \frac{2\lambda_{min}}{\lambda_{min} + \lambda_{max}} = \frac{\lambda_{max} - \lambda_{min}}{\lambda_{min} + \lambda_{max}}$$

what allows us to complete the proof. □

In the dynamic case, we get a result that allows us to optimally choose the iteration parameters at each step, if the matrix  $A$  is symmetric definite positive:

**Theorem 1** (Dynamic case). *If  $A$  is symmetric definite positive, the optimal choice for  $\alpha_k$  is given by*

$$\alpha_k = \frac{(\mathbf{r}^{(k)}, \mathbf{z}^{(k)})}{(A\mathbf{z}^{(k)}, \mathbf{z}^{(k)})}, \quad k \geq 0 \quad (9)$$

where

$$\mathbf{z}^{(k)} = P^{-1}\mathbf{r}^{(k)}. \quad (10)$$

**Demonstration** On the one hand we have

$$\mathbf{r}^{(k)} = \mathbf{b} - A\mathbf{x}^{(k)} = A(\mathbf{x} - \mathbf{x}^{(k)}) = \cancel{-A}\mathbf{e}^{(k)}, \quad (11)$$

sometimes not present in literature

and thus, using (10),

$$P^{-1}A\mathbf{e}^{(k)} = -\mathbf{z}^{(k)}, \quad (12)$$

where  $\mathbf{e}^{(k)}$  represents the error at the step  $k$ . On the other hand

$$\mathbf{e}^{(k+1)} = \mathbf{e}^{(k+1)}(\alpha) = \underbrace{(I - \alpha P^{-1}A)}_{R_\alpha} \mathbf{e}^{(k)}.$$

some kind of  $\alpha$   
 (constant or dynamic)

We notice that, in order to update the residual, we have the relation

$$\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} - \alpha A \mathbf{z}^{(k)} = \mathbf{r}^{(k)} - \alpha A P^{-1} \mathbf{r}^{(k)}.$$

Thus, expressing as  $\|\cdot\|_A$  the vector norm associated to the scalar product  $(\mathbf{x}, \mathbf{y})_A = (\mathbf{A}\mathbf{x}, \mathbf{y})$ , what means,  $\|\mathbf{x}\|_A = (\mathbf{A}\mathbf{x}, \mathbf{x})^{1/2}$  we can write

*Remember: A is symmetric*

$$\begin{aligned}
 \|\mathbf{e}^{(k+1)}\|_A^2 &= (A\mathbf{e}^{(k+1)}, \mathbf{e}^{(k+1)}) = -(\mathbf{r}^{(k+1)}, \mathbf{e}^{(k+1)}) \\
 &= -(\mathbf{r}^{(k)} - \alpha A P^{-1} \mathbf{r}^{(k)}, \mathbf{e}^{(k)} - \alpha P^{-1} A \mathbf{e}^{(k)}) \\
 &= -(\mathbf{r}^{(k)}, \mathbf{e}^{(k)}) + \alpha [(\mathbf{r}^{(k)}, P^{-1} A \mathbf{e}^{(k)}) + (A \mathbf{z}^{(k)}, \mathbf{e}^{(k)})] \\
 &\quad - \alpha^2 (A \mathbf{z}^{(k)}, P^{-1} A \mathbf{e}^{(k)})
 \end{aligned}$$

Now we choose  $\alpha$  as the  $\alpha_k$  that minimises  $\|\mathbf{e}^{(k+1)}(\alpha)\|_A$ :

$$\frac{d}{d\alpha} \|\mathbf{e}^{(k+1)}(\alpha)\|_A \Big|_{\alpha=\alpha_k} = 0$$

We then obtain

$$\alpha_k = \frac{1}{2} \frac{(\mathbf{r}^{(k)}, P^{-1} A \mathbf{e}^{(k)}) + (A \mathbf{z}^{(k)}, \mathbf{e}^{(k)})}{(A \mathbf{z}^{(k)}, P^{-1} A \mathbf{e}^{(k)})} = \frac{1}{2} \frac{-(\mathbf{r}^{(k)}, \mathbf{z}^{(k)}) + (A \mathbf{z}^{(k)}, \mathbf{e}^{(k)})}{-(A \mathbf{z}^{(k)}, \mathbf{z}^{(k)})}$$

and using the equality  $(A \mathbf{z}^{(k)}, \mathbf{e}^{(k)}) = (\mathbf{z}^{(k)}, A \mathbf{e}^{(k)})$  knowing that  $A$  is symmetric definite positive, and noting that  $A \mathbf{e}^{(k)} = -\mathbf{r}^{(k)}$ , we find

$$\alpha_k = \frac{(\mathbf{r}^{(k)}, \mathbf{z}^{(k)})}{(A \mathbf{z}^{(k)}, \mathbf{z}^{(k)})}$$

For the stationary case and for the dynamic one we can prove that, if  $A$  and  $P$  are symmetric definite positive, the series  $\{\mathbf{x}^{(k)}\}$  given by the Richardson method (stationary and dynamic) converges towards  $\mathbf{x}$  when  $k \rightarrow \infty$ , and

$$\|\mathbf{x}^{(k)} - \mathbf{x}\|_A \leq \left( \frac{K(P^{-1}A) - 1}{K(P^{-1}A) + 1} \right)^k \|\mathbf{x}^{(0)} - \mathbf{x}\|_A, \quad k \geq 0, \quad (13)$$

where  $\|\mathbf{v}\|_A = \sqrt{\mathbf{v}^T A \mathbf{v}}$  and  $K(P^{-1}A)$  is the conditioning of the matrix  $P^{-1}A$ .

↗ Again,  $K$  should make sense  $\leftrightarrow A$  should be stable

**Remark.** In the case of the gradient method or the Richardson stationary method the error estimation becomes

$$\|\mathbf{x}^{(k)} - \mathbf{x}\|_A \leq \left( \frac{K(A) - 1}{K(A) + 1} \right)^k \|\mathbf{x}^{(0)} - \mathbf{x}\|_A, \quad k \geq 0. \quad (14)$$

**Remark.** If  $A$  and  $P$  are symmetric definite positive, we have

$$K(P^{-1}A) = \frac{\lambda_{\max}(P^{-1}A)}{\lambda_{\min}(P^{-1}A)}.$$

# The conjugate gradient method

(Sec. 5.11)

When  $A$  and  $P$  are s.p.d, there exists a very efficient and effective method to iteratively solve the system: the conjugate gradient method

Let  $\mathbf{x}^{(0)}$  be given; we compute  $\mathbf{r}^{(0)} = \mathbf{b} - A\mathbf{x}^{(0)}$ ,  $\mathbf{z}^{(0)} = P^{-1}\mathbf{r}^{(0)}$ ,  
 $\mathbf{p}^{(0)} = \mathbf{z}^{(0)}$ , then for  $k \geq 0$ ,

$$\begin{aligned}\alpha_k &= \frac{\mathbf{p}^{(k)T} \mathbf{r}^{(k)}}{\mathbf{p}^{(k)T} A \mathbf{p}^{(k)}} \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)} \\ \mathbf{r}^{(k+1)} &= \mathbf{r}^{(k)} - \alpha_k A \mathbf{p}^{(k)} \\ P\mathbf{z}^{(k+1)} &= \mathbf{r}^{(k+1)} \\ \beta_k &= \frac{(A\mathbf{p}^{(k)})^T \mathbf{z}^{(k+1)}}{(A\mathbf{p}^{(k)})^T \mathbf{p}^{(k)}} \\ \mathbf{p}^{(k+1)} &= \mathbf{z}^{(k+1)} - \beta_k \mathbf{p}^{(k)}.\end{aligned}$$

Now we correct  
else the DIRECTION.  
wrong given by  $\mathbf{p}^{(k)}$

gives an

A conjugate (A  
orthogonal) direction  
with the direction  
 $\mathbf{p}^{(k-1)}$

special directions are

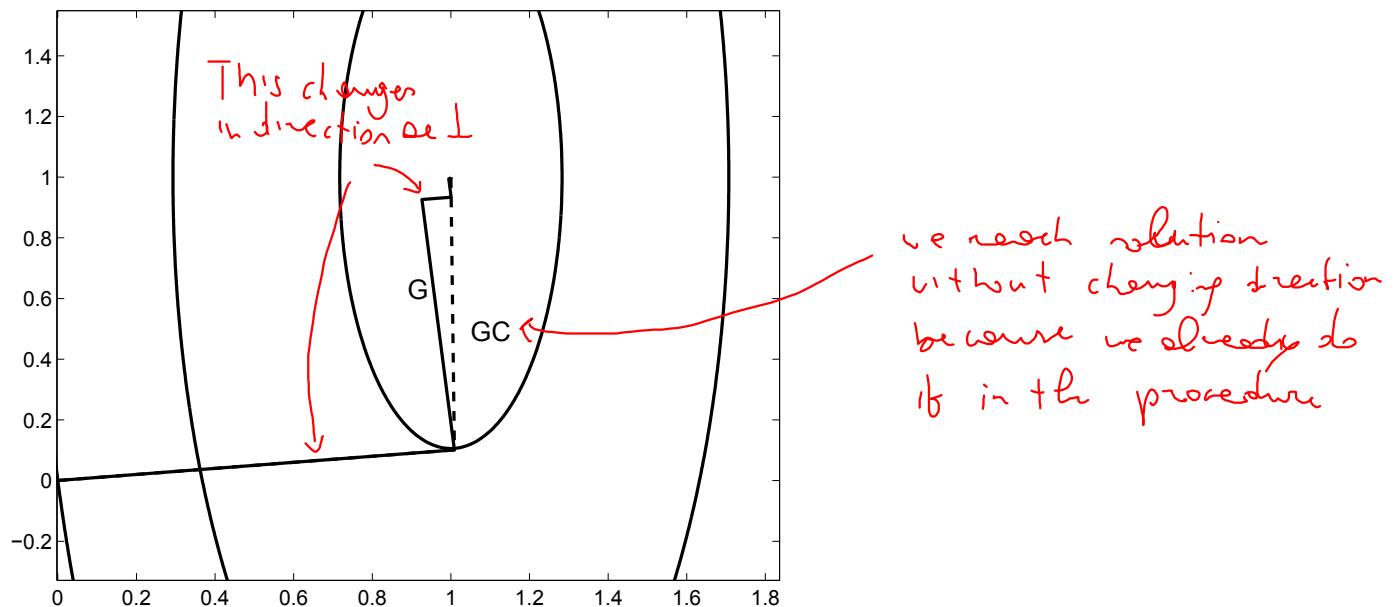
- Not so, but min number = to order of  $A$  matrix

The error estimate is given by

$$\|\mathbf{x}^{(k)} - \mathbf{x}\|_A \leq \frac{2c^k}{1 + c^{2k}} \|\mathbf{x}^{(0)} - \mathbf{x}\|_A, \quad k \geq 0$$

$$\text{ou } c = \frac{\sqrt{K_2(P^{-1}A)} - 1}{\sqrt{K_2(P^{-1}A)} + 1}.$$

(15)



**Example 2.** Let consider the following linear system:

$$\begin{cases} 2x_1 + x_2 = 1 \\ x_1 + 3x_2 = 0 \end{cases} \quad (16)$$

whose matrix is  $A = \begin{pmatrix} 2 & 1 \\ 1 & 3 \end{pmatrix}$  is s.p.d. The solution to this system is  $x_1 = 3/5 = 0.6$  et  $x_2 = -1/5 = -0.2$ .

## Preliminary convergence studies

- $A$  is strictly diagonal dominant by row. Hence Jacobi and Gauss-Seidel methods converge.
- $A$  is regular, tridiagonal with non-zero diagonal elements. Then  $\rho(B_{GS}) = \rho(B_J)^2$ . Therefore we expect a quicker convergence of Gauss-Seidel w.r.t. Jacobi.
- $A$  is s.p.d., hence the gradient and the conjugate gradient methods converge. Moreover (see error estimates), the CG shall converge faster.

We want to approximate the solution with an iterative method starting with

$$\mathbf{x}^{(0)} = \begin{pmatrix} x_1^{(0)} \\ x_2^{(0)} \end{pmatrix} = \begin{pmatrix} 1 \\ \frac{1}{2} \end{pmatrix}.$$

We can see that

$$\mathbf{r}^{(0)} = \mathbf{b} - A\mathbf{x}^{(0)} = \begin{pmatrix} -\frac{3}{2} \\ -\frac{5}{2} \end{pmatrix}$$

and

$$\|\mathbf{r}^{(0)}\|_2 = \sqrt{(\mathbf{r}^{(0)})^T \mathbf{r}^{(0)}} = \frac{\sqrt{34}}{2} \approx 2.9155.$$

## Jacobi method

$$\mathbf{x}^{(k+1)} = B_J \mathbf{x}^{(k)} + \mathbf{g}_J, \quad k \geq 0, \quad \text{where } B_J = I - D^{-1}A \text{ and } \mathbf{g}_J = D^{-1}\mathbf{b}.$$

We have

$$B_J = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} - \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{3} \end{pmatrix} \begin{pmatrix} 2 & 1 \\ 1 & 3 \end{pmatrix} = \begin{pmatrix} 0 & -\frac{1}{2} \\ -\frac{1}{3} & 0 \end{pmatrix}$$

$$\mathbf{g}_J = \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{3} \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} \frac{1}{2} \\ 0 \end{pmatrix}$$

and  $\rho(B_J) = \max|\lambda_i(B_J)| = \max(\text{abs}(\text{eig}(B_J))) = 0.4082$ .

For  $k = 0$  (first iteration) we find:

$$\mathbf{x}^{(1)} = B_J \mathbf{x}^{(0)} + \mathbf{g}_J = \begin{pmatrix} 0 & -\frac{1}{2} \\ -\frac{1}{3} & 0 \end{pmatrix} \begin{pmatrix} 1 \\ \frac{1}{2} \end{pmatrix} + \begin{pmatrix} \frac{1}{2} \\ 0 \end{pmatrix} = \begin{pmatrix} \frac{1}{4} \\ -\frac{1}{3} \end{pmatrix} \approx \begin{pmatrix} 0.25 \\ -0.3333 \end{pmatrix}.$$

Notice that

$$\mathbf{r}^{(1)} = \mathbf{b} - A\mathbf{x}^{(1)} = \begin{pmatrix} 0.8333 \\ 0.75 \end{pmatrix} \quad \text{and } \|\mathbf{r}^{(1)}\|_2 = 1.1211.$$

## Gauss-Seidel method

$$\mathbf{x}^{(k+1)} = B_{GS}\mathbf{x}^{(k)} + \mathbf{g}_{GS}, \quad k \geq 0, \quad \text{where } B_{GS} = (D - E)^{-1}(D - E - A)$$

$$\text{and } \mathbf{g}_{GS} = (D - E)^{-1}\mathbf{b}.$$

We have

$$B_{GS} = \begin{pmatrix} 2 & 0 \\ 1 & 3 \end{pmatrix}^{-1} \begin{pmatrix} 0 & -1 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} \frac{1}{2} & 0 \\ -\frac{1}{6} & \frac{1}{3} \end{pmatrix} \begin{pmatrix} 0 & -1 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & -\frac{1}{2} \\ 0 & \frac{1}{6} \end{pmatrix}$$

$$\mathbf{g}_{GS} = \begin{pmatrix} \frac{1}{2} & 0 \\ -\frac{1}{6} & \frac{1}{3} \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} \frac{1}{2} \\ -\frac{1}{6} \end{pmatrix}$$

In this case  $\rho(B_{GS}) = \max|\lambda_i(B_{GS})| = \max(\text{abs}(\text{eig}(B_{GS}))) = 0.1667$ .

We can verify that  $\rho(B_{GS}) = \rho(B_J)^2$ .

For  $k = 0$  (first iteration) we find:

$$\mathbf{x}^{(1)} = B_{GS}\mathbf{x}^{(0)} + \mathbf{g}_{GS} = \begin{pmatrix} 0 & -\frac{1}{2} \\ 0 & \frac{1}{6} \end{pmatrix} \begin{pmatrix} 1 \\ \frac{1}{2} \end{pmatrix} + \begin{pmatrix} \frac{1}{2} \\ -\frac{1}{6} \end{pmatrix} = \begin{pmatrix} \frac{1}{4} \\ -\frac{1}{12} \end{pmatrix} \approx \begin{pmatrix} 0.25 \\ -0.0833 \end{pmatrix}.$$

We have

$$\mathbf{r}^{(1)} = \mathbf{b} - A\mathbf{x}^{(1)} = \begin{pmatrix} 0.5833 \\ 0 \end{pmatrix} \quad \text{and } \|\mathbf{r}^{(1)}\|_2 = 0.5833.$$

## Preconditioned gradient method with $P = D$

We set  $\mathbf{r}^{(0)} = \mathbf{b} - A\mathbf{x}^{(0)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} - \begin{pmatrix} 2 & 1 \\ 1 & 3 \end{pmatrix} \begin{pmatrix} 1 \\ \frac{1}{2} \end{pmatrix} = \begin{pmatrix} -\frac{3}{2} \\ -\frac{5}{2} \end{pmatrix}$ .

For  $k = 0$ , we have:

$$P\mathbf{z}^{(0)} = \mathbf{r}^{(0)} \Leftrightarrow \mathbf{z}^{(0)} = P^{-1}\mathbf{r}^{(0)} = \begin{pmatrix} -\frac{3}{4} \\ -\frac{5}{6} \end{pmatrix}$$

$$\alpha_0 = \frac{(\mathbf{z}^{(0)})^T \mathbf{r}^{(0)}}{(\mathbf{z}^{(0)})^T A \mathbf{z}^{(0)}} = \frac{77}{107}$$

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{z}^{(0)} = \begin{pmatrix} 0.4603 \\ -0.0997 \end{pmatrix}$$

$$\mathbf{r}^{(1)} = \mathbf{r}^{(0)} - \alpha_0 A \mathbf{z}^{(0)} = \begin{pmatrix} 0.1791 \\ -0.1612 \end{pmatrix} \quad \text{and} \quad \|\mathbf{r}^{(1)}\|_2 = 0.2410.$$

## Conjugated preconditioned gradient method with $P = D$

We set  $\mathbf{r}^{(0)} = \mathbf{b} - A\mathbf{x}^{(0)}$ ,  $\mathbf{z}^{(0)} = P^{-1}\mathbf{r}^{(0)}$  and  $\mathbf{p}^{(0)} = \mathbf{z}^{(0)}$ . For  $k = 0$ , we have:

$$\alpha_0 = \frac{(\mathbf{p}^{(0)})^T \mathbf{r}^{(0)}}{(\mathbf{p}^{(0)})^T A \mathbf{p}^{(0)}} = \frac{(\mathbf{z}^{(0)})^T \mathbf{r}^{(0)}}{(\mathbf{z}^{(0)})^T A \mathbf{z}^{(0)}}$$

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{p}^{(0)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{z}^{(0)}$$

$$\mathbf{r}^{(1)} = \mathbf{r}^{(0)} - \alpha_0 A \mathbf{p}^{(0)} = \mathbf{r}^{(0)} - \alpha_0 A \mathbf{z}^{(0)}.$$

We see that the first iteration  $\mathbf{x}^{(1)}$  matches with the one obtained by the preconditioned gradient method.

We then complete the first iteration of the preconditioned conjugate gradient method:

$$P\mathbf{z}^{(1)} = \mathbf{r}^{(1)} \iff \mathbf{z}^{(1)} = P^{-1}\mathbf{r}^{(1)} = \begin{pmatrix} 0.0896 \\ -0.0537 \end{pmatrix}$$

$$\beta_0 = \frac{(A\mathbf{p}^{(0)})^T \mathbf{z}^{(1)}}{(A\mathbf{p}^{(0)})^T A\mathbf{p}^{(0)}} = \frac{(A\mathbf{z}^{(0)})^T \mathbf{z}^{(1)}}{(A\mathbf{z}^{(0)})^T \mathbf{z}^{(0)}} = -0.0077$$

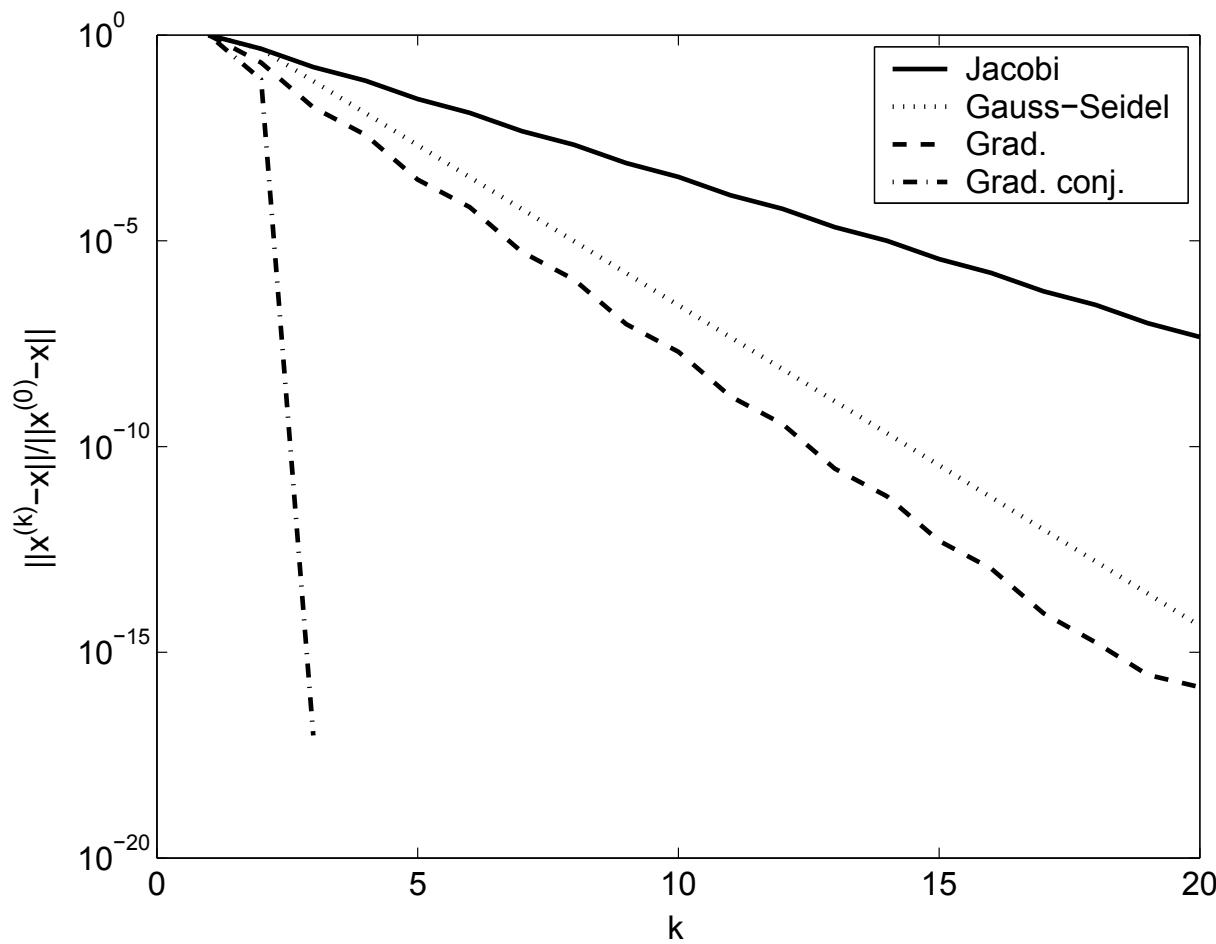
$$\mathbf{p}^{(1)} = \mathbf{z}^{(1)} - \beta_0 \mathbf{p}^{(0)} = \mathbf{z}^{(1)} - \beta_0 \mathbf{z}^{(0)} = \begin{pmatrix} 0.0838 \\ -0.0602 \end{pmatrix}.$$

At the second iteration, with the four different methods, we have:

Method	$\mathbf{x}^{(2)}$	$\mathbf{r}^{(2)}$	$\ \mathbf{r}^{(2)}\ _2$
Jacobi	$\begin{pmatrix} 0.6667 \\ -0.0833 \end{pmatrix}$	$\begin{pmatrix} -0.2500 \\ -0.4167 \end{pmatrix}$	0.4859
Gauss-Seidel	$\begin{pmatrix} 0.5417 \\ -0.1806 \end{pmatrix}$	$\begin{pmatrix} 0.0972 \\ 0 \end{pmatrix}$	0.0972
PG	$\begin{pmatrix} 0.6070 \\ -0.1877 \end{pmatrix}$	$\begin{pmatrix} -0.0263 \\ -0.0438 \end{pmatrix}$	0.0511
PCG	$\begin{pmatrix} 0.60000 \\ -0.2000 \end{pmatrix}$	$\begin{pmatrix} -0.2220 \\ -0.3886 \end{pmatrix} \cdot 10^{-15}$	$4.4755 \cdot 10^{-16}$

Behavior of the relative error applied to the system (16) :

Graph important  
for the exam



**Example 3.** Let now consider another example:

$$\begin{cases} 2x_1 + x_2 &= 1 \\ -x_1 + 3x_2 &= 0 \end{cases} \quad (17)$$

whose solution is  $x_1 = 3/7$ ,  $x_2 = 1/7$ .

## Preliminary convergence studies

The associated matrix is  $A = \begin{pmatrix} 2 & 1 \\ -1 & 3 \end{pmatrix}$ . ← not symmetric  
=> not positive definite

- $A$  is strictly diagonal dominant by row. Hence Jacobi and Gauss-Seidel methods converge.
- $A$  is regular, tridiagonal with non-zero diagonal elements. Then  $\rho(B_{GS}) = \rho(B_J)^2$ . Therefore we expect a quicker convergence of Gauss-Seidel w.r.t. Jacobi.
- $A$  is **not s.p.d.**, therefore we have no idea if the gradient or the conjugate gradient converge.

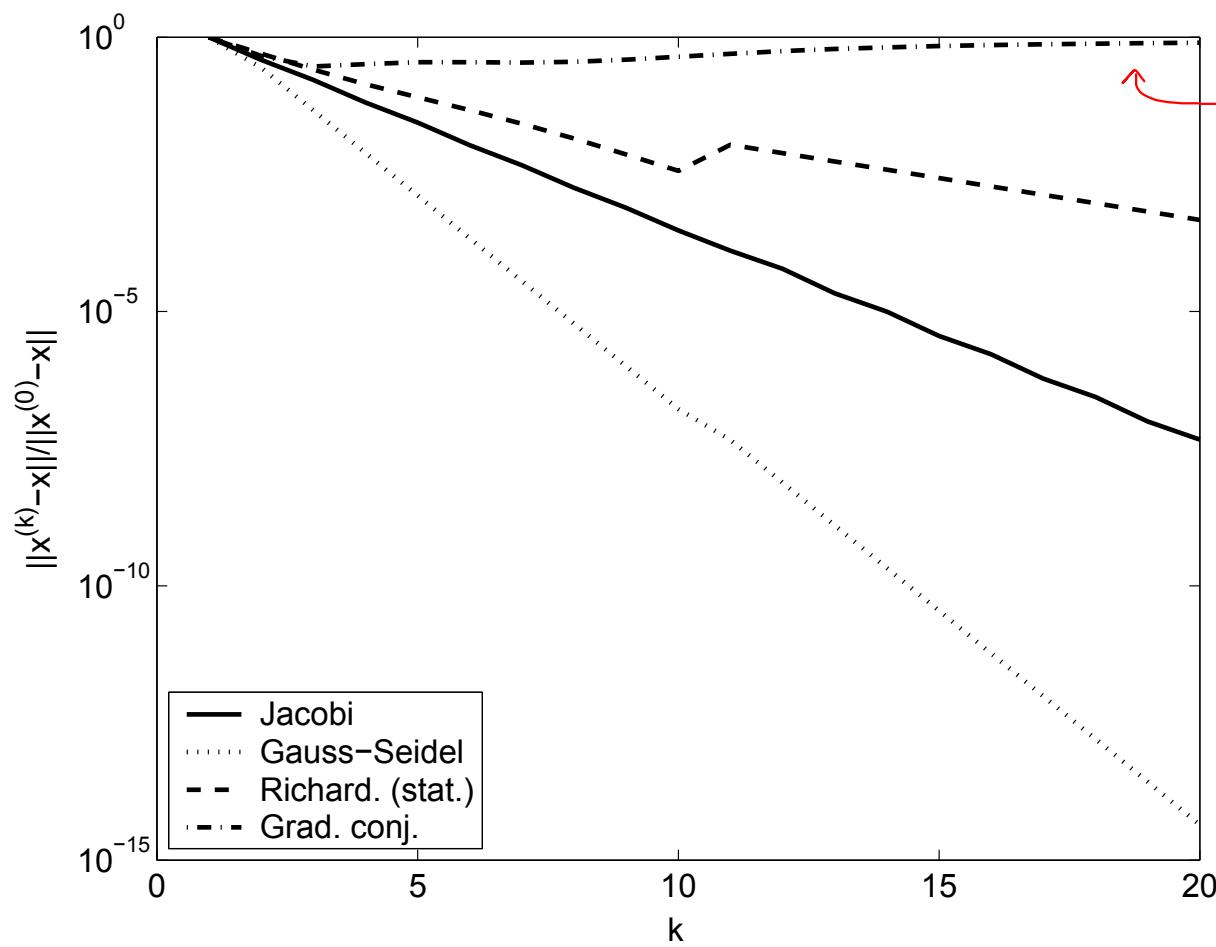
We approximate the solution with an iterative method starting from

$$\mathbf{x}^{(0)} = \begin{pmatrix} x_1^{(0)} \\ x_2^{(0)} \end{pmatrix} = \begin{pmatrix} 1 \\ \frac{1}{2} \end{pmatrix}.$$

The following figure shows the value of  $\frac{\|\mathbf{x}^{(k)} - \mathbf{x}\|}{\|\mathbf{x}^{(0)} - \mathbf{x}\|}$  for the Jacobi, Gauss-Seidel, Richardson stationary (preconditioned with  $\alpha = 0.5$  and  $P = D = \begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix}$ ), and the preconditioned (with  $P = D$ ) conjugate gradient methods.

Remark that this time the preconditioned conjugate gradient method doesn't converge.

Behavior of the relative error applied to the system (17) :



why? we are applying wrong methodology, since A is NOT symmetric

when this happen  
 use other method,  
 i.e. conjugate gradient

# Convergence Criteria

(Sec. 5.12)

We have the following error bound:

*If  $A$  is s.p.d, then*

$$\frac{\|\mathbf{x}^{(k)} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq K(A) \frac{\|\mathbf{r}^{(k)}\|}{\|\mathbf{b}\|}. \quad (18)$$

The relative error at the iteration  $k$  is bounded by the condition number of  $A$  times the residual scaled with the right hand side.

We can also use another relation in case of a preconditioned system:

$$\frac{\|\mathbf{x}^{(k)} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq K(P^{-1}A) \frac{\|P^{-1}\mathbf{r}^{(k)}\|}{\|P^{-1}\mathbf{b}\|}.$$

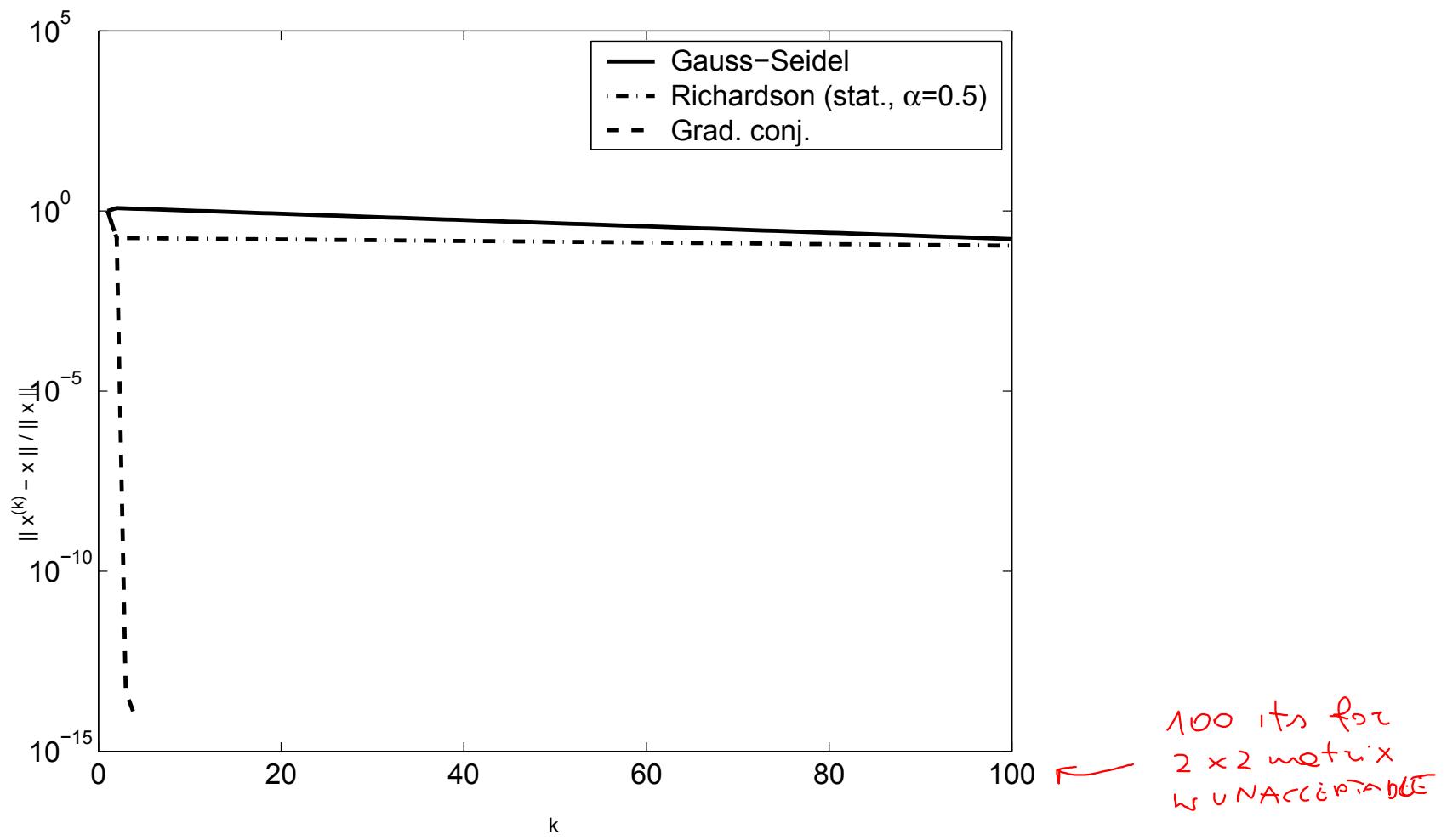
Some examples of convergence of iterative methods applied to linear systems of the kind  $A\mathbf{x} = \mathbf{b}$ .

**Example 4.** Lets start with the matrix

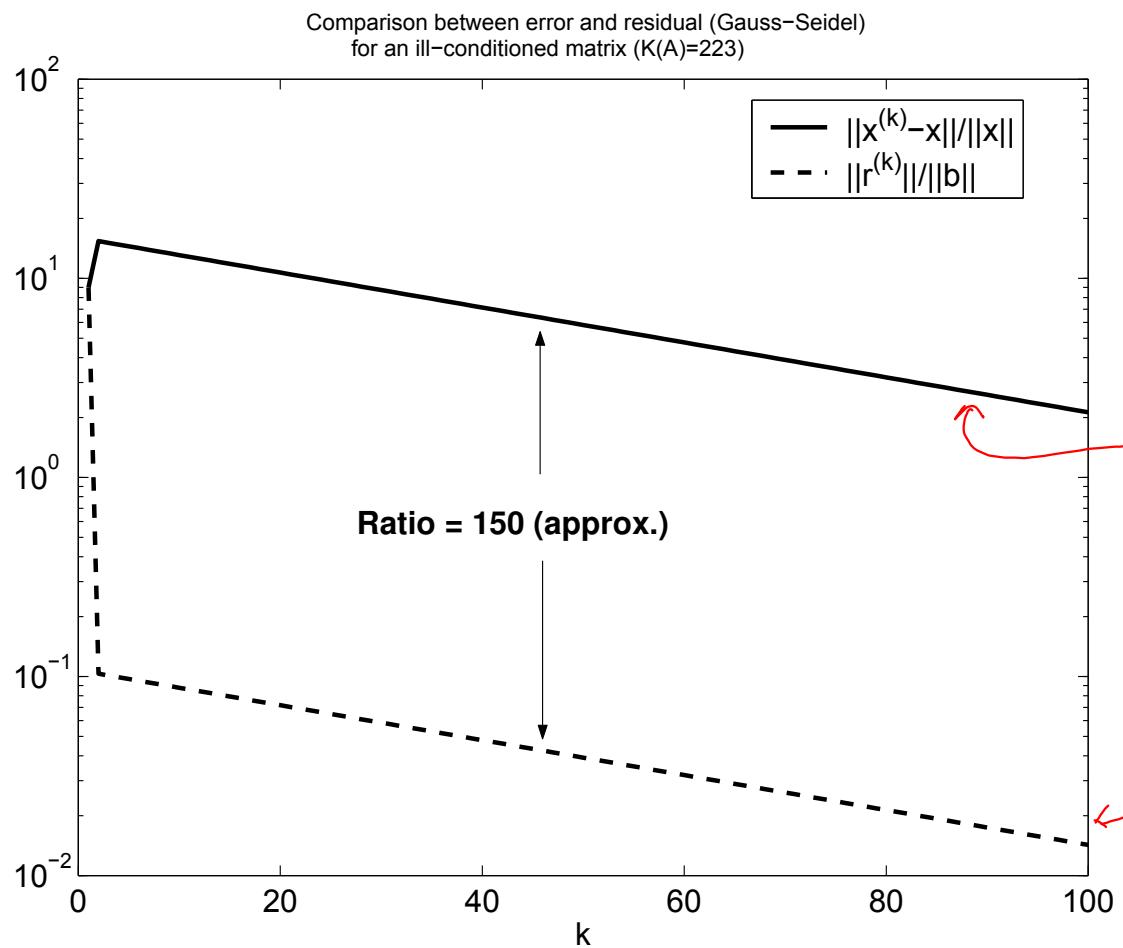
$$A = \begin{pmatrix} 5 & 7 \\ 7 & 10 \end{pmatrix} \quad (19)$$

The condition number of this matrix is  $K(A) \approx 223$ . We consider the Gauss-Seidel method and stationary preconditioned Richardson method with  $\alpha = 0.5$  and  $P = \text{diag}(A)$ . We have that  $\rho(B_{GS}) = 0.98$  and  $\rho(B_{Rich}) = 0.98$ , where  $B_{Rich} = I - \alpha P^{-1} A$  is the iterative matrix for the stationary Richardson method.

This figure shows the relative error behavior for Gauss-Seidel method and stationary preconditioned Richardson method and the conjugate gradient preconditioned with the same matrix



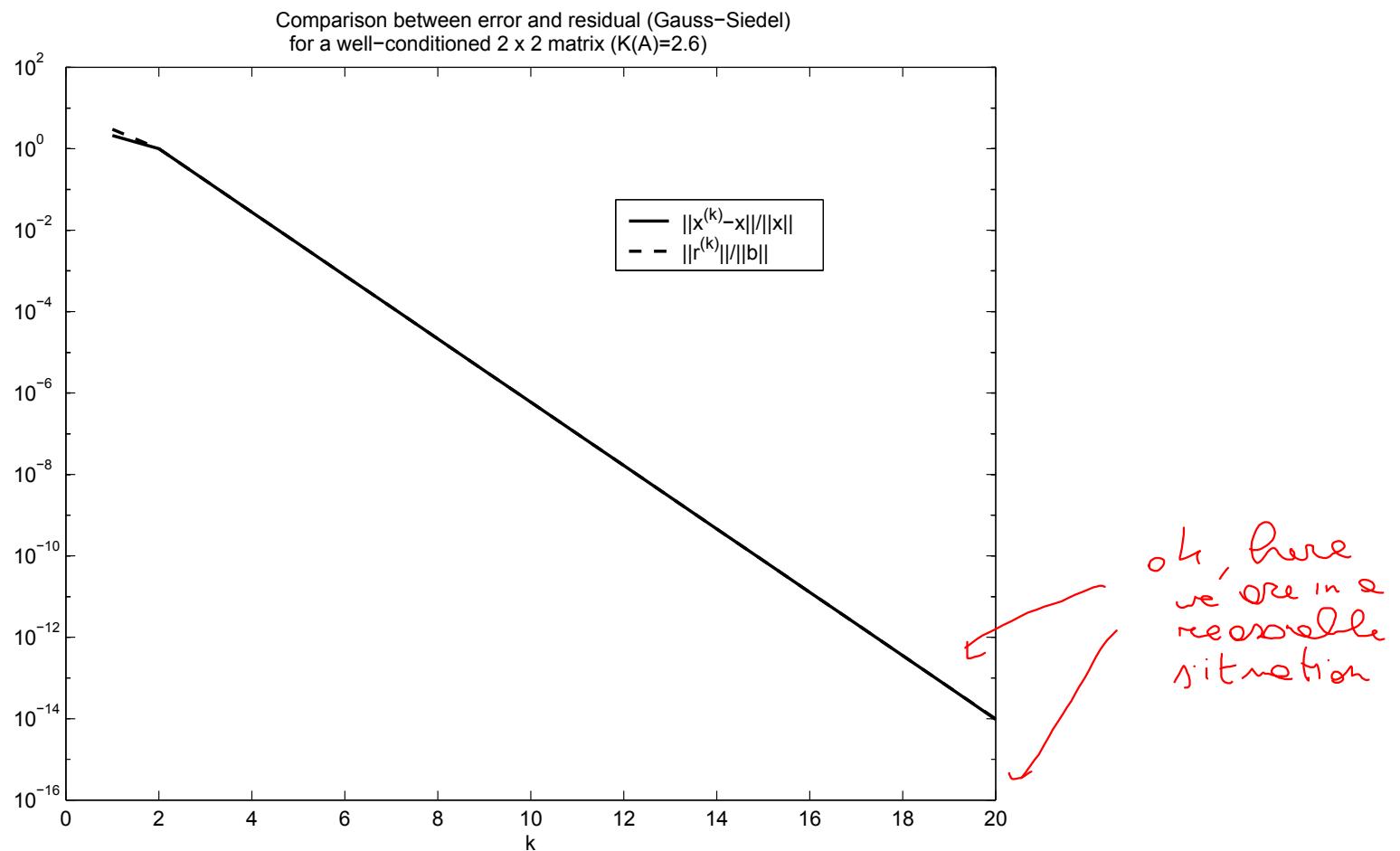
Comparison between error and residual (Gauss-Seidel) (19). Recall that  $K(A) \approx 223$



Look AT THE DB  
S P E M A N T I U D E !  
But often this  
line is not available

If a matrix  $2 \times 2$   
requires 100 its  
should tell me  
sth is wrong

We can compare the previous figure with the same curves for well conditioned matrix  $2 \times 2$  ( $K(A) \approx 2.6$ ):



*Comparison between the relative error and the residual, for a well conditioned matrix.*

**Example 5.** We take the Hilbert matrix and we solve a linear system for different size  $n$ . We set  $P = \text{diag}(A)$  and use the preconditioned gradient method. The stoping tolerance is set to  $10^{-6}$  on the relative residual. We take  $\mathbf{x}^{(0)} = \mathbf{0}$ .

We evaluate the relative error

```
for j=2:7;
    n=2*j; i=j-1; nn(i)=n;
    A=hilb(n);
    x_ex=ones(n,1); b=A*x_ex;
    Acond(i)=cond(A);
    tol=1.e-6; maxit=10000;
    R=diag(diag(A));
    x0=zeros(n,1);
    [x,iter_gr(i)]=gradient(A,b,x0,maxit,tol,R);
    error_gr(i)=norm(x-x_ex)/norm(x_ex);
end
```

For the iterative method, we set a tolerance of  $10^{-6}$  on the relative residual. Because the matrix is ill conditioned, it is to be expected to get a relative error in the solution greater than  $10^{-6}$  (see inequality (18)).

		Gradient method		
$n$	$K(A)$	Error	Iterations	Residual
4	1.55e+04	8.72e-03	995	1.00e-06
6	1.50e+07	3.60e-03	1813	9.99e-07
8	1.53e+10	6.30e-03	1089	9.96e-07
10	1.60e+13	7.99e-03	875	9.99e-07
12	1.67e+16	5.09e-03	1355	9.99e-07
14	2.04e+17	3.91e-03	1379	9.98e-07

↑  
 Already a WARNING  
 A 0D of magnitude  
 greater than  
 residual → could lead to very wrong  
 assumption on the problem you  
 are trying to solve

# Some observations

# Memory and computational costs

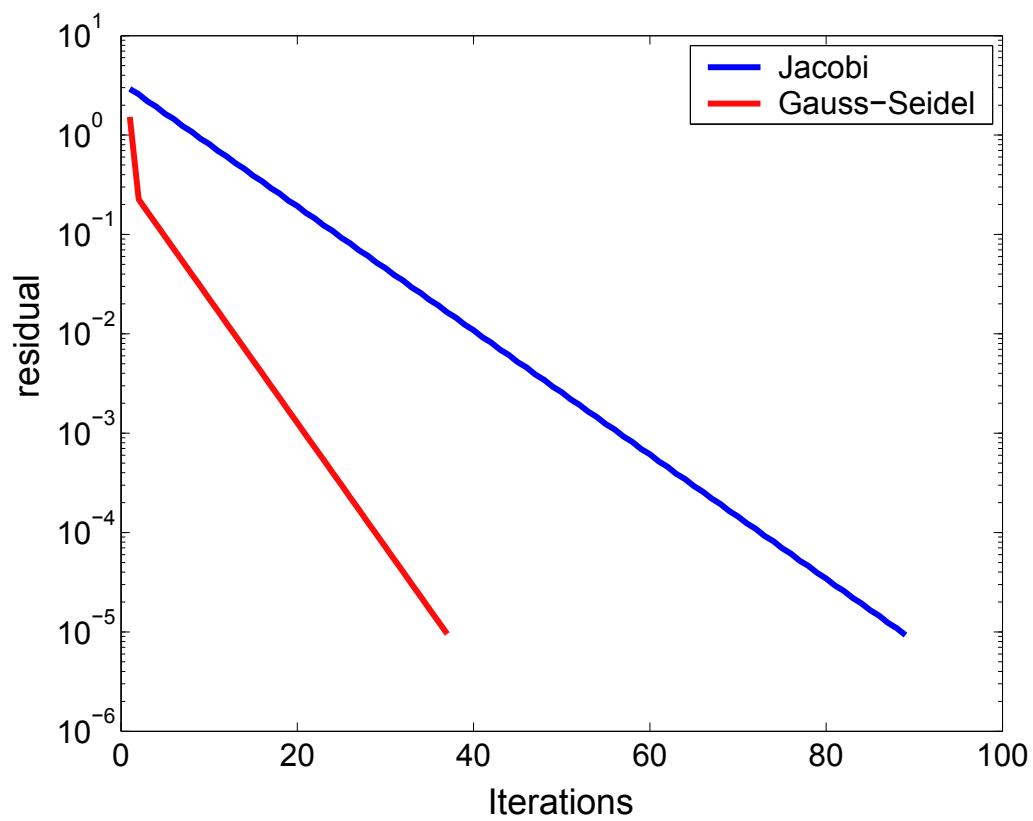
Computational cost (*flops*) and used memory (*bytes*): we consider the Cholesky and the conjugate gradient methods for sparse matrices of size  $n$  (originated in the approximation of solutions to the Poisson equation in the finite elements method) with  $m$  non zero elements.

		Cholesky		Conjugate gradient		flops(Chol.)/ flops(GC)	Mem(Chol.)/ Mem(GC)
$n$	$m/n^2$	flops	Memory	flops	Memory		
47	0.12	8.05e+03	464	1.26e+04	228	0.64	2.04
83	0.07	3.96e+04	1406	3.03e+04	533	1.31	2.64
150	0.04	2.01e+05	4235	8.86e+04	1245	2.26	3.4
225	0.03	6.39e+05	9260	1.95e+05	2073	3.27	4.47
329	0.02	1.74e+06	17974	3.39e+05	3330	5.15	5.39
424	0.02	3.78e+06	30815	5.49e+05	4513	6.88	6.83
530	0.01	8.31e+06	50785	8.61e+05	5981	9.65	8.49
661	0.01	1.19e+07	68468	1.11e+06	7421	10.66	9.23

Conjugate gradient solves 10 times the system and needs then the Cholesky method

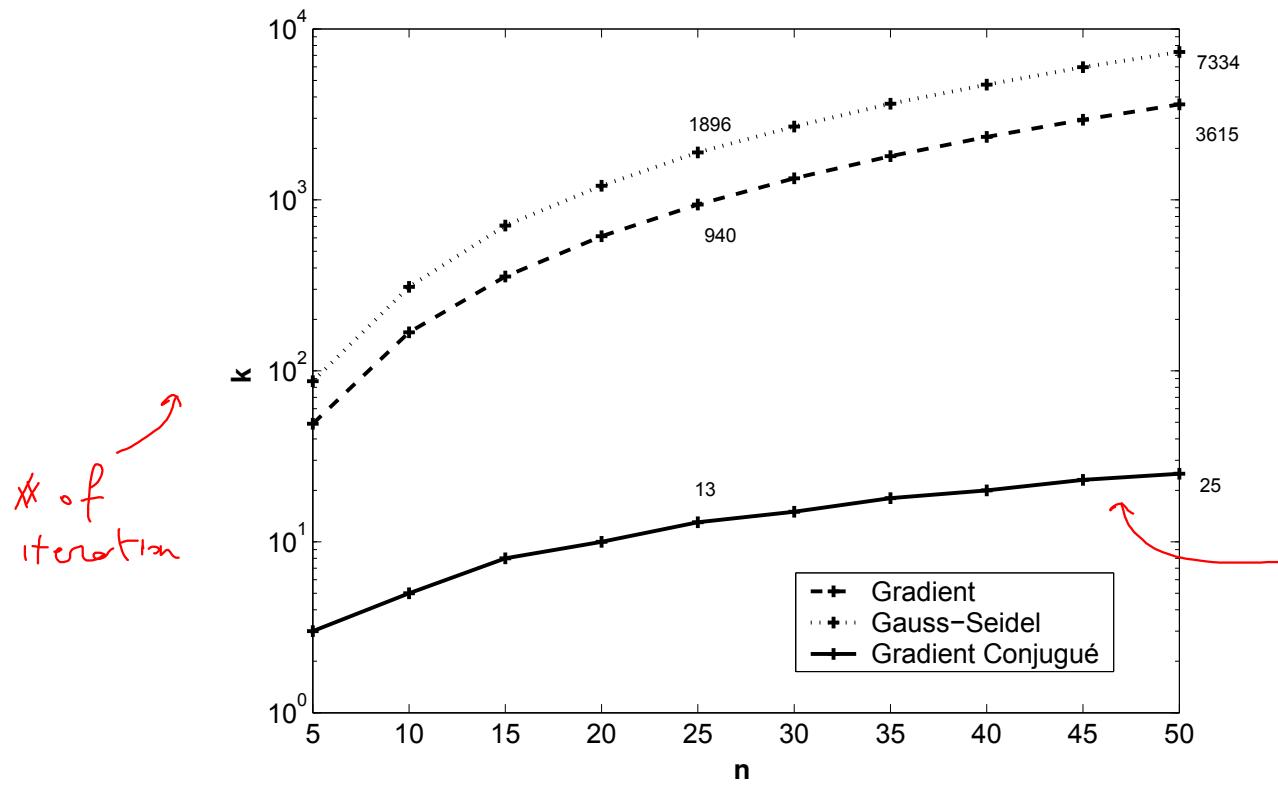
# Iterative methods: Jacobi, Gauss-Seidel

Behaviour of the error for a well conditioned matrix ( $K = 20$ )



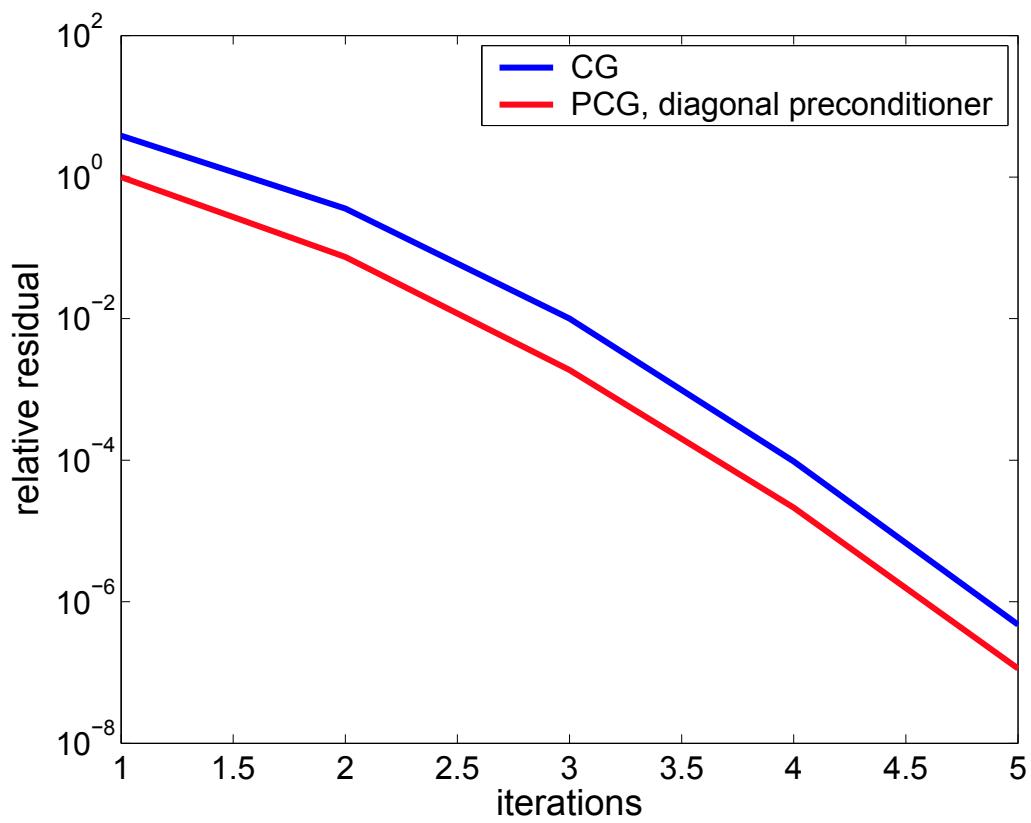
# Iterative methods: G-S, Gradient, Gradient Conjugate

Number of iterations as function of the size  $n$  of a stiffness matrix, tolerance  $10^{-6}$



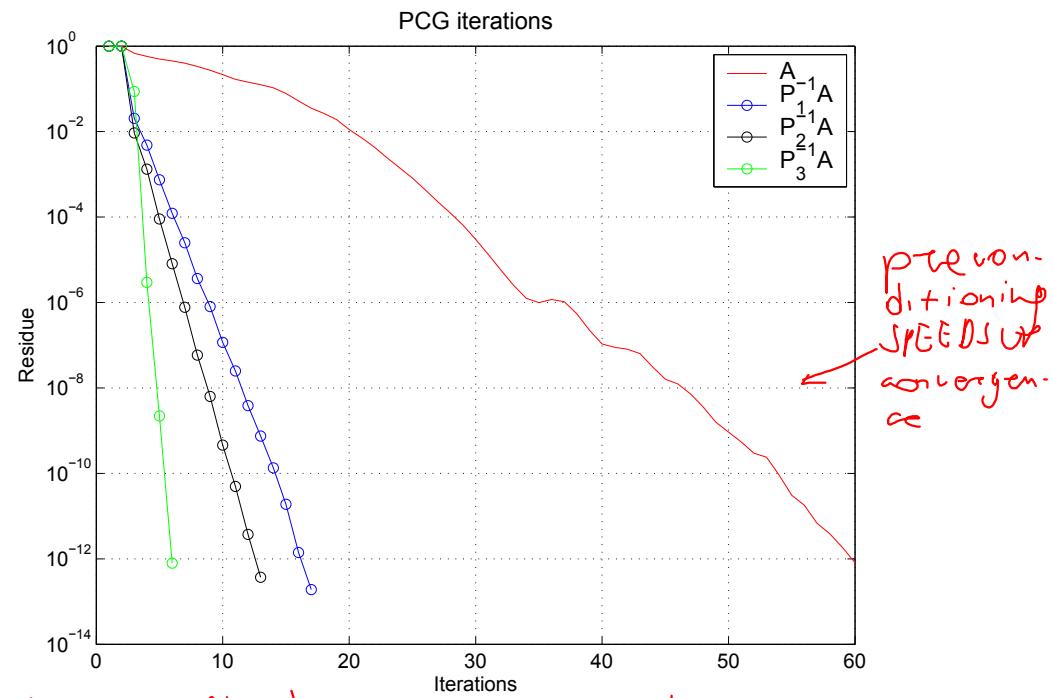
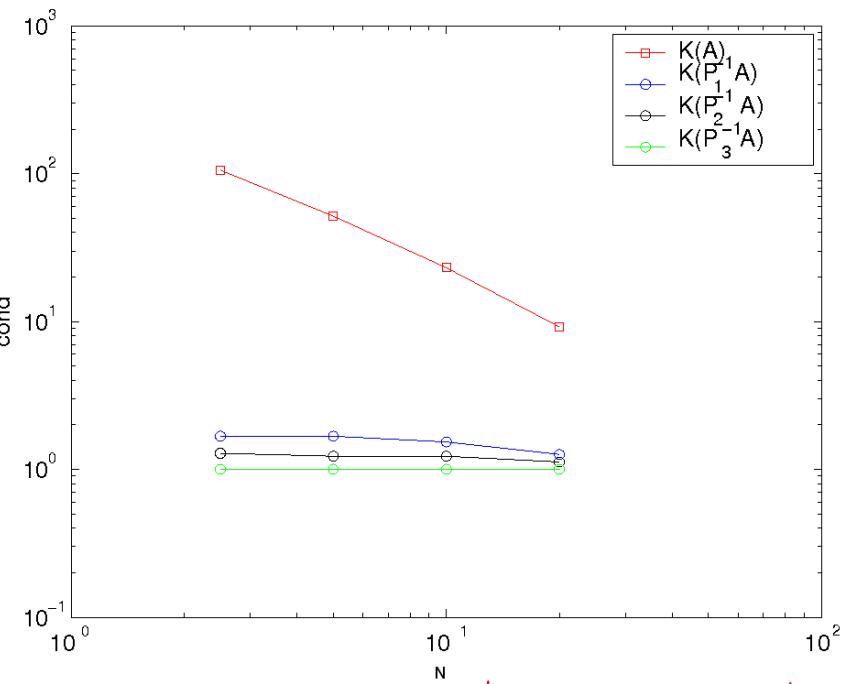
# Preconditioning

The convergence conjugate gradient and the preconditioned conjugate gradient methods with a diagonal preconditioner ( $K = 4 \cdot 10^8$ )



# Preconditioning

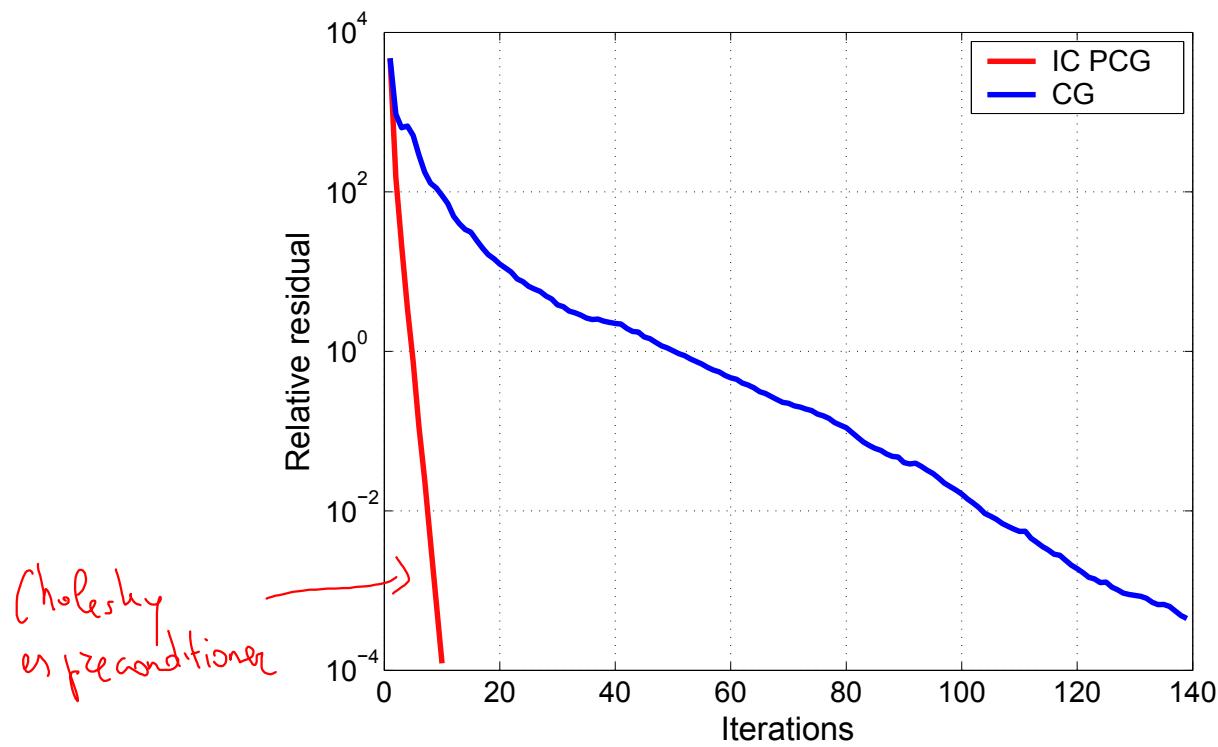
Role of the preconditioning over the condition number of a matrix  
 (approximation by the finite element method of the Laplacian operator)



Question: do you know which matrix could be preconditioned and which  
 can't? If  $K \sim 1$  could be preconditioned. If  $K$  is high, is  
 likely not preconditioned

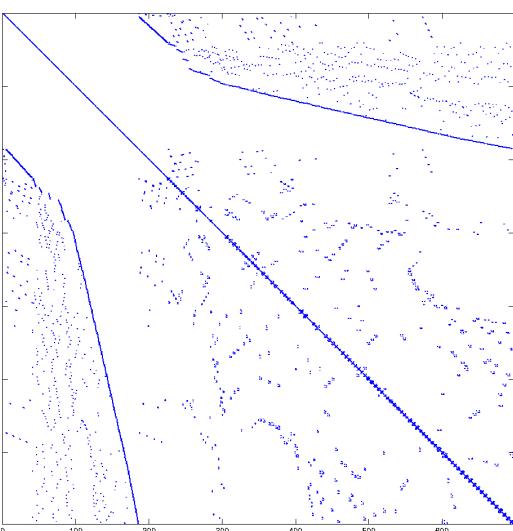
# Preconditioning

The convergence of the conjugate gradient method and the preconditioned conjugate gradient method with a preconditioner based on an Cholesky incomplete decomposition for a sparse matrix originated from the finite element method ( $K = 1.5 \cdot 10^3$ )

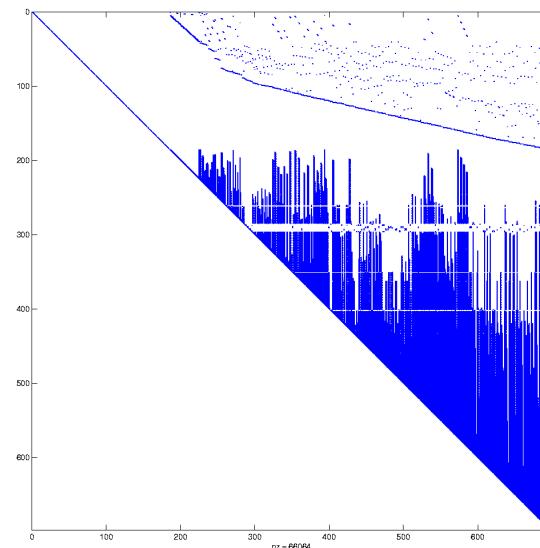


# Pre-conditioning

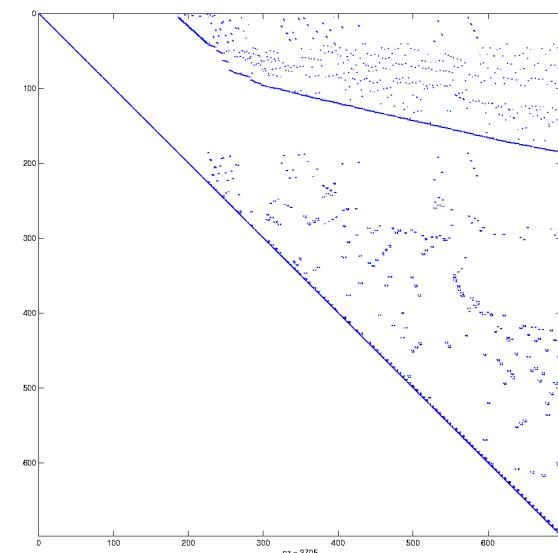
Comparison between the non zero elements of the sparse matrix  $A$  from the previous example, its Cholesky factor  $R$  and the matrix  $\tilde{R}$  obtained by the Cholesky incomplete decomposition:



$A$



$R$



$\tilde{R}$

# Non-linear systems

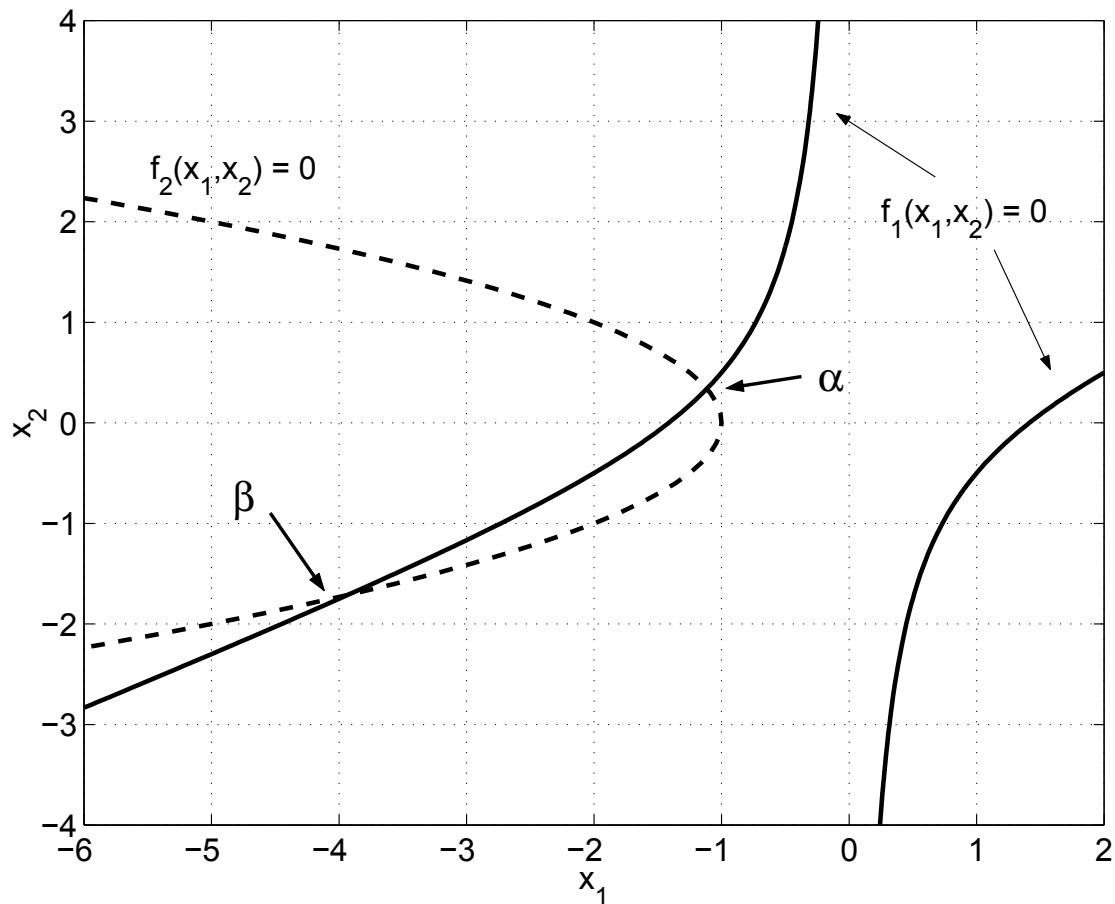
**Example 6.** Lets consider the following system of non-linear equations:

$$\begin{cases} x_1^2 - 2x_1x_2 = 2 \\ x_1 + x_2^2 = -1. \end{cases} \quad (20)$$

This system can be written in the form:

$$\mathbf{f}(\mathbf{x}) = \mathbf{0} \quad \text{i.e.} \quad \begin{cases} f_1(x_1, x_2) = 0 \\ f_2(x_1, x_2) = 0 \end{cases}$$

where  $\mathbf{f} = (f_1, f_2)$ ,  $f_1(x_1, x_2) = x_1^2 - 2x_1x_2 - 2$  and  $f_2(x_1, x_2) = x_1 + x_2^2 + 1$ .



*Curves  $f_1 = 0$  and  $f_2 = 0$  in the square  $-6 \leq x_1 \leq 2, -4 \leq x_2 \leq 4$*

We want to generalise the [Newton method](#) for the case of non-linear systems. To do this, we define the *Jacobian* matrix of the vector  $\mathbf{f}$ :

$$J_{\mathbf{f}}(\mathbf{x} = (x_1, x_2)) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1 - 2x_2 & -2x_1 \\ 1 & 2x_2 \end{bmatrix}.$$

If  $J_{\mathbf{f}}(\mathbf{x}^{(k)})$  is invertible, the Newton method for non linear systems is written : Let  $\mathbf{x}^{(0)} = (x_1^{(0)}, x_2^{(0)})$ , we compute for  $k = 0, 1, 2, \dots$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - [J_{\mathbf{f}}(\mathbf{x}^{(k)})]^{-1} \mathbf{f}(\mathbf{x}^{(k)}), \quad k = 0, 1, 2, \dots \quad (21)$$

We can write (21) as

C in place of  $1/p^1$  of Newton method, we put plus 1 each iteration -1 (inverse of Jacobian)

$$[J_{\mathbf{f}}(\mathbf{x}^{(k)})](\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = -\mathbf{f}(\mathbf{x}^{(k)}), \quad k = 0, 1, 2, \dots \quad (22)$$

↑ quasi-newton method update Jacobian every m iterations with m > 1

Description of the first step of the algorithm: given the vector  $\mathbf{x}^{(0)} = [1, 1]^T$ .  
We calculate:

$$J_{\mathbf{f}}(\mathbf{x}^{(0)}) = \begin{bmatrix} 0 & -2 \\ 1 & 2 \end{bmatrix}.$$

We determine  $\mathbf{x}^{(1)}$  as solution of the equation :

$$[J_{\mathbf{f}}(\mathbf{x}^{(0)})](\mathbf{x}^{(1)} - \mathbf{x}^{(0)}) = -\mathbf{f}(\mathbf{x}^{(0)}).$$

And we continue with (22).

## Lecture 21

## Solution of prototypical PDE in dimension 1.

10/12/2020

### Prototypical problem

$$\begin{cases} -u'' = f & \text{in } [0,1] \\ u(0) = u(1) = 0 \end{cases}$$

Unknown of the problem is a ft

Boundary cns

Given  $f$  (regular enough) find  $u$  s.t.

Problems like this can be found:

• Heat conduction

• Elasticity

• Maxwell equations

Requirement for  
WELL POSEDNESS

Strong formulation,  $f \in C^0([0,1])$ , and  $u \in C^2([0,1])$

First approach

This prob will be well posed if  
the final linear system is invertible  
(det of matrix is  $\neq 0$ )

## FINITE DIFFERENCES.

Approximate the differential operators using the incremental definition of the derivative:

Can do  
obs for  
 $u''$

$$u'(x) := \lim_{h \rightarrow 0} \frac{u(x+h) - u(x)}{h}$$

Forward difference  
FD  
for  $C^1$   
fcts there  
are ( $\Rightarrow$   
for  $h \rightarrow 0$ )

$$u'(x) := \lim_{h \rightarrow 0} \frac{u(x) - u(x-h)}{h}$$

backward  
BD

$$u'(x) := \lim_{h \rightarrow 0} \frac{u(x+h) - u(x-h)}{2h}$$

centered  
CD

Then the FINITE DIFFERENCE scheme is:

Fix  $h$ , and define  $D_u$  as

$$D_u^F = \frac{u(x+h) - u(x)}{h} \quad FD$$

$$D_u^B = \frac{u(x) - u(x-h)}{h} \quad BD$$

$$D_u^C = \frac{u(x+h) - u(x-h)}{2h} \quad CD$$

Assume  $u \in C^\infty(I(x))$  in an open set of  $x \in I(x)$

then a Taylor expansion

$$u(x+h) = u(x) + u'(x)h + \frac{u''(x)h^2}{2} + \frac{u'''(x)h^3}{6} + \frac{u^{(4)}(x)h^4}{24} + \dots + \frac{u^{(n)}(x)h^n}{n!}$$

$$u(x-h) = u(x) - u'(x)h + \frac{u''(x)h^2}{2} - \frac{u'''(x)h^3}{6} + \frac{u^{(4)}(x)h^4}{24} + \dots + \frac{u^{(n)}(-h)h^n}{n!}$$

Theorem Given  $u \in C^\infty([a,b])$ ,  $x \in (a,b)$ , then  $\forall n > 0$

equality EXACT  $\exists \xi \in (x, x+h)$

$$u(x+h) = \left( \sum_{i=0}^{k-1} \frac{u^{(i)}(x)h^i}{i!} \right) + \frac{u^{(k)}(\xi)h^k}{k!}$$

You can  
take a taylor  
expansion  
at any point  
and what you  
find is the value  
of the ft  
at one point  
in  $(x, x+h)$

Truncation Error of Finite Difference schemes:

FD  $\frac{u(x+h) - u(x)}{h} = u'(x) + \frac{u''(\xi)h}{2}$  1st order scheme

BD  $\frac{u(x) - u(x-h)}{h} = u'(x) - \frac{u''(\xi)h}{2}$  1st order scheme

CD  $\frac{u(x+h) - u(x-h)}{2h} = u'(x) + \frac{u'''(\xi)h^2}{6}$  2nd order scheme

Use low order (1st) if ft is not regular (maybe  $u''$  doesn't exist) and we h low, while high order if u is regular, and h high

CFD II  $\frac{u(x-h) - 2u(x) + u(x+h)}{h^2} = u''(x) + \frac{u''''(\xi)h^2}{12}$  2nd order scheme for 2nd derivative

Apply CFD<sup>II</sup> to  $-u'' = f$ ,  $u(0) = u(1) = 0$  in  $[0, 1]$

Split the interval  $(0, 1)$  into  $N-1$  subintervals ( $N$  points)

$$\{x_i\}_{i=0}^{N-1} = \left\{ 0, \frac{1}{N-1}, \frac{2}{N-1}, \dots, \frac{N-1}{N-1} = 1 \right\} = \left\{ \frac{i}{(N-1)} \right\}_{i=0}^{N-1} = \{i \cdot h\}_{i=0}^{N-1}$$

$$h = \frac{1}{(N-1)}$$

$$\text{set } u^i := u(x_i)$$

$$\{u_i\}_{i=0}^{N-1} \in \mathbb{R}^N$$

For  $i = 1$  to  $N-2$  write: For all:

$$\frac{1}{h^2} \left[ -u^{i-1} + 2u^i - u^{i+1} \right] \underset{\approx -u''}{\sim} f_i \equiv A \underline{u} = \underline{f}$$

system  
of  
equations

$$A \in \mathbb{R}^{N \times N}, \underline{u}, \underline{f} \in \mathbb{R}^N$$

$$\sum_{j=0}^{N-1} A_{ij} u^j = f_i$$

$$A_{ij} = \frac{1}{h^2} \begin{cases} -1 & \text{when } j = i-1 \text{ or } i+1 \\ 2 & \text{when } j = i \\ 0 & \text{otherwise} \end{cases} \quad \text{and } i \in \{1, \dots, N-2\}$$

$$\begin{pmatrix} 0 & -1 & 2 & -1 & & & & \\ -1 & 2 & -1 & & & & & \\ & \ddots & \ddots & \ddots & & & & \\ & & -1 & 2 & -1 & & & \\ & & & & \ddots & \ddots & & \\ & & & & & 0 & & \\ & & & & & & \ddots & \\ N-2 & & & & & & & 0 \end{pmatrix} \cdot \frac{1}{h^2}$$

$$f_i = f(x_i) \quad i = 1, \dots, N-2$$

For  $i=0, i=N-1$  we impose

$$u^0 = 0$$

$$u^{N-1} = 0$$



$$A_{00} = 1$$

$$A_{N-1, N-1} = 1$$

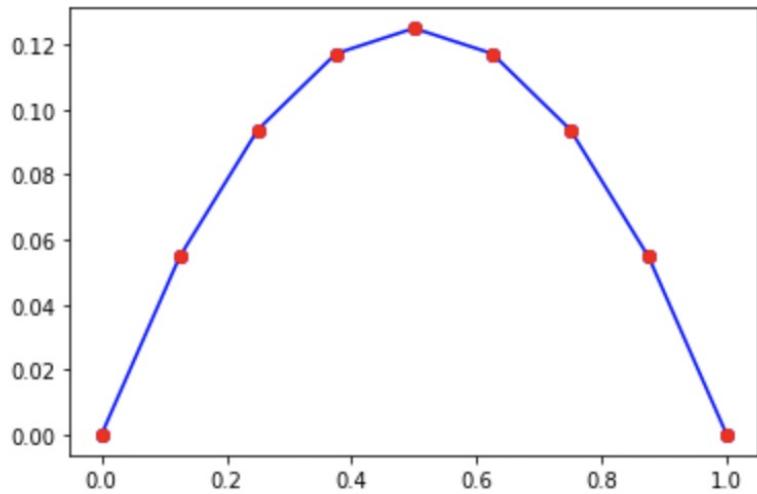
$$f_0 = 0$$

$$f_{N-1} = 0$$

We set the BOUNDARY CONDITION

! The more the points the higher the order NUMERIC PRECONDITIONER

$\Rightarrow$  choose good PRECONDITIONER



Example of solution  
to  
 $-u'' = 1$   
 $u(0) = u(1) = 0$   
with  $N = 9$

the red dots are the exact solution.

$$u_{\text{exact}} = \frac{x(1-x)}{2}$$

exact function

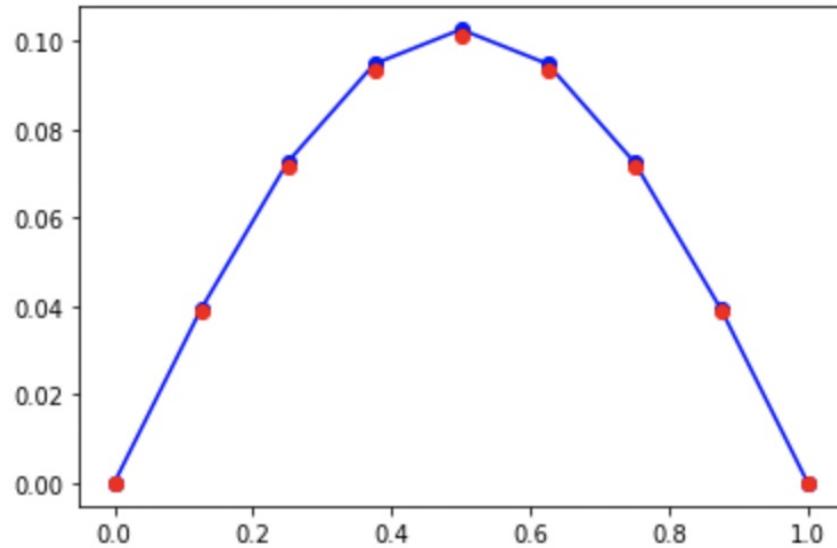
We have errors (blue overlaps red)  
because error is up to  $u'''$ , but  $u'''$  of  $e$   
parallel is zero!

Same thing with

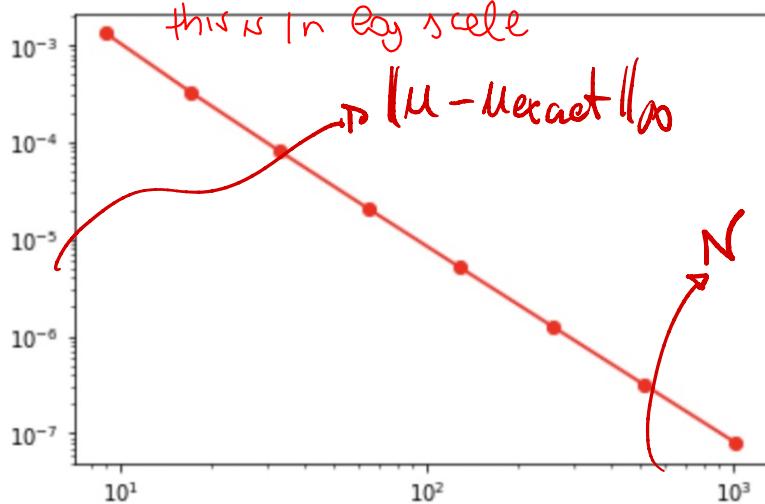
$$f(x) = \sin(\pi x)$$

Here the error is  
not zero, because

$$u'''' = -f'' \neq 0$$



Plot the error  $\|u - u_{\text{exact}}\|_\infty = \max_{i \in \{0, N-1\}} |u^i - u(x_i)|$



Straight line :

$$E = c h^\kappa$$

$$= c \left(\frac{1}{N-1}\right)^\kappa$$

what is  $\kappa$ ?

inclina-  
tion of  
curve (with  
minim line)  
( $\rightarrow$  4th and 2nd E)

# How to measure the rate $k$ ?

- Two options:
- 1) we know the exact solution
  - 2) we do not know the exact solution

## Case 1

Hypothesis:  $C$  is indep. of  $w$  / and  $h$   $\rightarrow$  target

$$E_1 = \|\mu_{h_1} - u_{\text{exact}}\|_\infty \simeq C h_1^k$$
$$E_2 = \|\mu_{h_2} - u_{\text{exact}}\|_\infty \simeq C h_2^k$$

some for the 2 axis

$$\frac{E_1}{E_2} \simeq \left( \frac{h_1}{h_2} \right)^k$$

We need at least 2 measurements

$$k = \frac{\log(E_1/E_2)}{\log(h_1/h_2)} = \frac{\log(E_1) - \log(E_2)}{\log(h_1) - \log(h_2)}$$

## Case 2

Make sure we reduce  $h$  by a constant factor  $\theta$ , and use at least 3 approx solutions.

$$\|\mu_{h_2} - \mu_{h_1}\|_\infty \leq \|\mu_{h_2} - u_{\text{exact}}\|_\infty + \|\mu_{h_1} - u_{\text{exact}}\|_\infty$$

$$\simeq C h_1^k + C h_2^k$$

$$\simeq C h_1^k + C \theta^k h_2^k$$

$$\simeq \overline{C(1+\theta^k)} h_1^k$$

$$\frac{\|\mu_{h_1} - \mu_{h_2}\|_\infty}{\|\mu_{h_2} - \mu_{h_3}\|_\infty} \sim \frac{C_2 h_1^{-k}}{C_2 h_2^{-k}} = \bar{\theta}^{-k}$$

$$h_2 = \theta h_1$$

$$h_3 = \theta h_2$$

$\leftarrow$  I do a third measurement  
and I do the same  
by the fact that  
2 consecutive steps go  
the same in terms of

$$k = \frac{\log(\|\mu_{i+2} - \mu_i\|_\infty) - \log(\|\mu_{i+1} - \mu_i\|_\infty)}{\log(\bar{\theta})}$$

I remember  
that we had  
 $m_i, m_{i+1}, m_{i+2}$ , that  
are 3 consecutive  
approximate  
solutions

## Lecture 22 LH

### PDE - 1D - FEM

Finite Element Method

15.12.2020

$$\textcircled{1} \quad \begin{cases} -u'' = f & \text{in } [0,1] \\ u(0) = u(1) = 0 \end{cases}$$

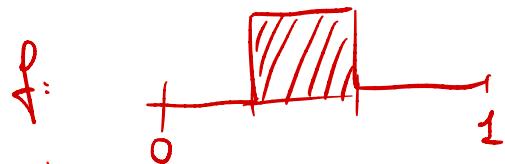
We tried to:  
approximate  $u''$

- split  $[0,1]$  in  $N-1$  intervals with size  $h = \frac{1}{N-1}$
- define approx of  $u''$  on each  $x_i = ih \quad i = 0, \dots, N-1$

$$-u^{i-1} + 2u^i - u^{i+1} = f_i$$

- issues: i)  $u''$  is defined if  $u \in C^2([0,1])$
- ii) ( $f \in C^0([0,1])$ )

excluded situations:



- iii) truncation error: defines the order by which we approximate  $u''$

$$u'' = CD^2(u) + O(\|u'''\|_\infty h^2)$$

SOLUTION: rewrite  $\textcircled{1}$  in weak form

Let's see a special case

Let  $v$  be a smooth function:  $V: [0,1] \rightarrow \mathbb{R}$

Multiply  $\textcircled{1}$  by  $v$  and integrate on  $[0,1]$

$$\int_0^1 -u'' \cdot v = \int_0^1 f \cdot v \quad \nabla v \text{ "smooth enough"}$$

Let's further assume that

→ integrate by parts:  
so the problem is now:

$$v(0) = V(1) = 0 \quad \text{because } V(0) = V(1) = 0$$

$$(wv)' = w'v + wv' \quad w = u'$$

Find  $u$  in  $V$  s.t.

$\nabla v \in V$

$$\textcircled{2} \quad \int_0^1 u' v' = \int_0^1 f v$$

Weak formulation

of  $\textcircled{1}$  ( $-u'' = f$ )

This is weaker: we are asking  $u$  to have only one degree of continuity (first derivative not necessarily continuous)

What's the "right" space  $V$  for ② to make sense?

i)  $V$  must incorporate that  $v(0) = v(1) = 0 \quad \forall v \in V$

ii) ② must make sense if  $u, v \in V$

$$\xrightarrow{\text{in particular}} v = u \Rightarrow \int_0^1 (u')^2 = \int_0^1 f u$$

$\Rightarrow$  It only makes sense if  $\int_0^1 (u')^2 < +\infty$

So it must be

$$\Rightarrow V = \{v \text{ s.t. } \int_0^1 (v')^2 < +\infty, \quad v(0) = v(1) = 0\}$$

( $V$  is the Sobolev space  $H_0^1([0, 1])$ )

space of  $f$  + but one in  $L^2$  and whose first derivatives are in  $L^2$

$V$  is the space of functions whose first derivative is in  $L^2$ , and that are zero on  $\{0, 1\}$

A solution to ① is a solution to ② (always true!)

$$(Wv)' = W'v + Wv'$$

$$\Rightarrow \int_0^1 (Wv)' = \int_0^1 W'v + \int_0^1 Wv'$$

it's not in weak st but not the contrary in general

$$\int_0^1 Wv' = \int_0^1 W'v + \int_0^1 Wv'$$

$$\Rightarrow \int_0^1 W'v = - \int_0^1 Wv' + \int_0^1 Wv'$$

$$\boxed{\int_0^1 -u''v = \int_0^1 u'v' + \int_0^1 u'v'} \quad \text{for the because } v(0) = v(1) = 0$$

$$-u'' = \int_0^1 f v \Rightarrow \text{false: } \int_0^1 u'v' = \int_0^1 f v$$

$\forall v \in V$

You def the WEAK sol of PDE as the solution of this problem

$$V \equiv H_0^1([0,1]) := \{ v \in L^2, v' \in L^2, v(0) = v(1) = 0 \}$$

f must be s.t.  $\int_0^1 f v < +\infty \quad \forall v \in H_0^1([0,1])$

true if  $f \in L^2$

but also if f is less reg. than  $L^2$ , but

such that  $\int_0^1 f \cdot v < +\infty \quad \forall v \in H_0^1([0,1])$

$$f \in V^* \xleftarrow{\text{dual space of } V} \xleftarrow{\text{what is that?}}$$

for example,  $Dne f$  is not in  $L^2$ , but this works anyway

Weak derivative:

$$Du(x) = u'(x) \quad \forall u \in C^1([0,1]) \quad , \quad \forall x \in [0,1]$$

$$\underline{Du(x)}$$
 is s.t.  $- \int_0^1 u \varphi' = \int_0^1 Du \varphi \quad \forall \varphi \in C_0^\infty([0,1])$

By construction, if  $u \in C^1([0,1]) \Rightarrow$

$$- \int_0^1 u \varphi' = \underbrace{\int_0^1 u' \varphi}_{\substack{\text{integration by part}}} = \int_0^1 Du \varphi \Rightarrow Du = u'$$

$$Du = u^{(\alpha)} \quad \text{if } u \in C^\alpha([0,1])$$

$$\boxed{\int_0^1 D u, \varphi = (-1)^\alpha \int_0^1 \varphi^{(\alpha)} Du \quad \forall \varphi \in C_0^\infty([0,1])}$$

↑ true only if  $u \in C^1$ , but then we can be def whenever  $(u^{(\alpha)})$  exist

$D u$  is weak derivative of  $u$ , if this is satisfied

$$\text{where } C_0^\infty([0,1]) := \{ \varphi \in C^\infty([0,1]) \text{ s.t. } \varphi^{(\alpha)}(0) = \varphi^{(\alpha)}(1) = 0 \quad \forall \alpha \in \mathbb{N}_0 \}$$

The final weak form of our problem is

Given  $f \in L^2([0,1])$  (could be  $(H_0^1([0,1]))^*$ )  
dual of  $H_0^1([0,1])$

Find  $u \in H_0^1([0,1])$  s.t.

$$\int_0^1 u' v' = \int_0^1 f v \quad \forall v \in H_0^1([0,1])$$

$$H_0^1([0,1]): \{v \in L^2([0,1]), v' \in L^2([0,1]), v(0) = v(1) = 0\}$$

Assume that we have  $V_h \subset V \equiv H_0^1(\Omega)$

$$V_h = \text{Span} \{ v_i \}_{i=0}^{N-1} \quad \dim(V_h) = N \leftarrow \text{finite dim}$$

Write a approximate problem:

$$\text{find } u_h \in V_h \quad \text{s.t.} \quad (u_h(x) = \sum_{j=0}^{N-1} u_j^j v_j(x))$$

$$\int_0^1 u_h' v_i' = \int_0^1 v_i f \quad i = 0, \dots, N-1$$

$\leftarrow$  that is for ALL basis functions

$$\rightarrow \sum_{j=0}^{N-1} \int_0^1 u_j^j v_j' v_i' = \int_0^1 v_i f$$

$$A u = F$$

$$A_{ij} := \int_0^1 v_j' v_i' \quad F_i = \int_0^1 v_i f$$

| DO NOT CHANGE WAY IN WHICH  
| COMPUTE DERIVATIVES:

| well choose basis so I can compute  
| derivative exactly

Define  $V_h$  piecewise polynomial and  $C^0([0,1])$

Select  $N$  points that define  $N-1$  intervals;

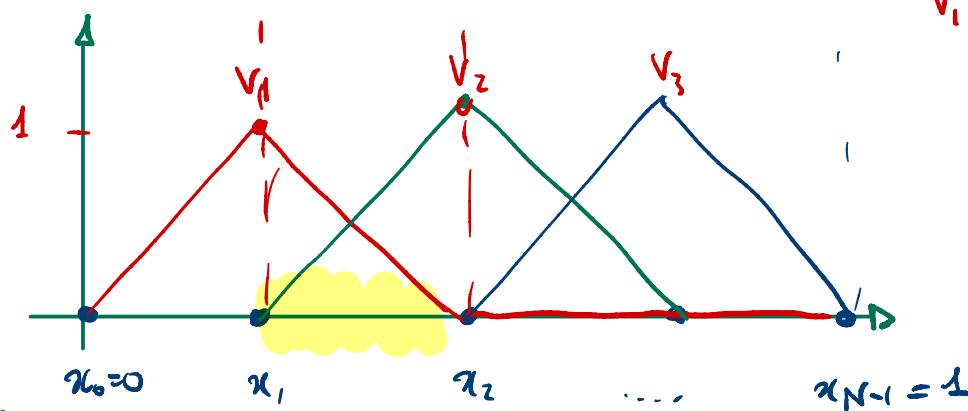
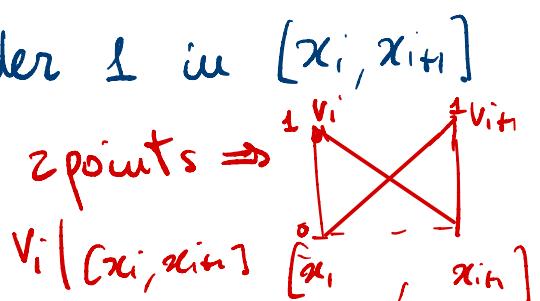
$$\{x_i\}_{i=0}^{N-1} = ih = \frac{i}{(N-1)} \quad \text{the intervals are } [x_i, x_{i+1}] \quad i=0, \dots, N-2$$

$$V_h^k := \left\{ v \in \mathbb{R}^k([x_i, x_{i+1}]) \mid \begin{array}{l} v(0) = v(1) = 0 \\ v \in C^0([0, 1]) \end{array} \right\}$$

Continuous piecewise polynomial of order  $k$ .

Example  $k=1 \Rightarrow P^1([x_i, x_{i+1}])$

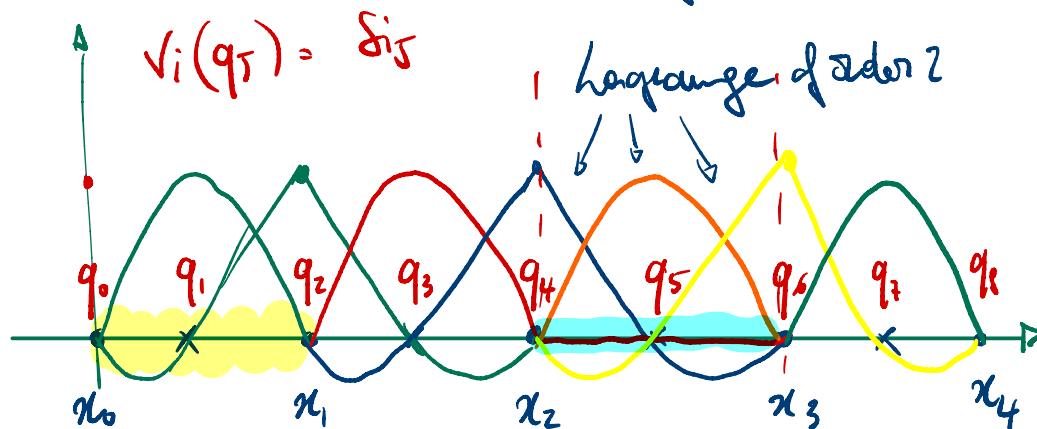
We choose Lagrange basis of order 1 in  $[x_i, x_{i+1}]$   
with support points  $\{x_i, x_{i+1}\}$



$q_i = x_i$   
for linear  
polynomials.

For example, choose

$$v_i(x_j) = \delta_{ij} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}, \quad v_i \in V_h^1$$

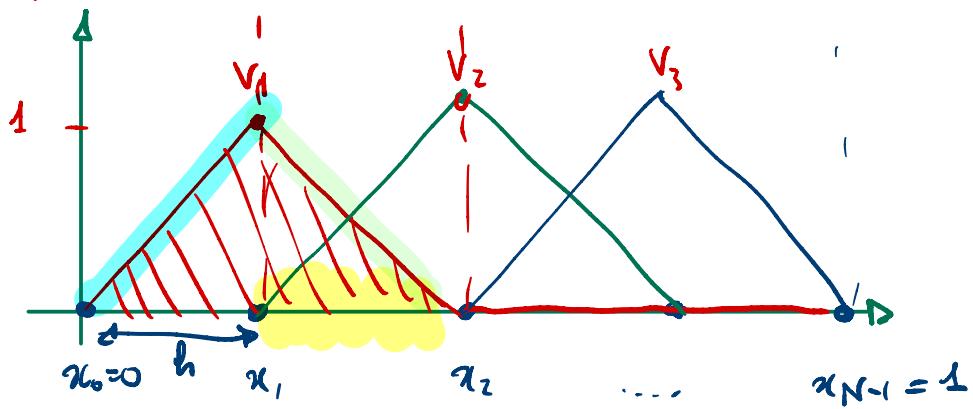


For polynomials of

order 2, choose

$$\{q_i\}_{i=0}^{(N-1)(N-2)} = \{x_i\}_{i=0}^{N-1} \cup \left\{ \frac{x_i + x_{i+1}}{2} \right\}_{i=0}^{N-2}$$

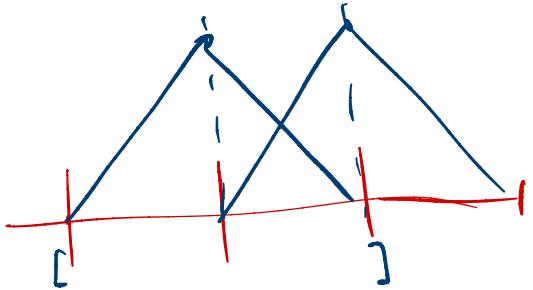
1D case linear



$$\dim(V_h) = N-2$$

$$v_i = \begin{cases} \frac{(x - x_{i-1})}{h} & x \in [x_{i-1}, x_i] \\ \frac{(x_{i+1} - x)}{h} & x \in [x_i, x_{i+1}] \\ 0 & \text{elsewhere} \end{cases} \quad i = 1, \dots, N-2$$

$$v_i^l = \begin{cases} \frac{1}{h} & \text{in } [x_{i-1}, x_i] \\ -\frac{1}{h} & \text{in } [x_i, x_{i+1}] \\ 0 & \text{elsewhere} \end{cases}$$



$$A_{ij} := \int_0^1 v_i^l v_j^l \neq 0 \quad \begin{matrix} \text{if } i = j, j-1, j+1 \\ \text{elsewhere} \end{matrix}$$

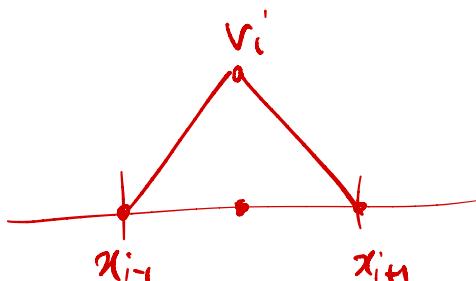
$$A_{jj} = \int_{x_{i-1}}^{x_{i+1}} \left(\frac{1}{h}\right)^2 = 2 \frac{1}{h} \quad \text{if } i=j$$

$$A_{ij} = \int_{x_i}^{x_{i+1}} \left(-\frac{1}{h}\right)^2 = -\frac{1}{h} \quad \begin{matrix} \text{if } i = j-1 \\ \text{or } i = j+1 \end{matrix}$$

$$A_{ij} := \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & & -1 & 2 & -1 \end{pmatrix} \frac{1}{h}$$

Same matrix as FD ( $\frac{1}{h}$ )

$$F_i = \int_{x_{i-1}}^{x_{i+1}} f(x) v_i(x) dx$$



$\approx f(x_i) \cdot h$  using Trapez rule (same as FD with  $h$ )

proof

$$\int_{x_{i-1}}^{x_i} v_i(x) f(x) dx + \int_{x_i}^{x_{i+1}} v_i(x) f(x) dx$$

since integral of  $v_i$  between  $x_{i-1}, x_{i+1}$  is  $1/2$  of the whole rectangle of height  $1/h$

$$h \left[ \frac{1}{2} v_i(x_{i+1}) f(x_{i+1}) + \frac{1}{2} v_i(x_i) f(x_i) \right] + h \left[ \frac{1}{2} v_i(x_i) f(x_i) + \frac{1}{2} (v_i(x_{i+1}) f(x_{i+1})) \right]$$

$$= f(x_i) h$$

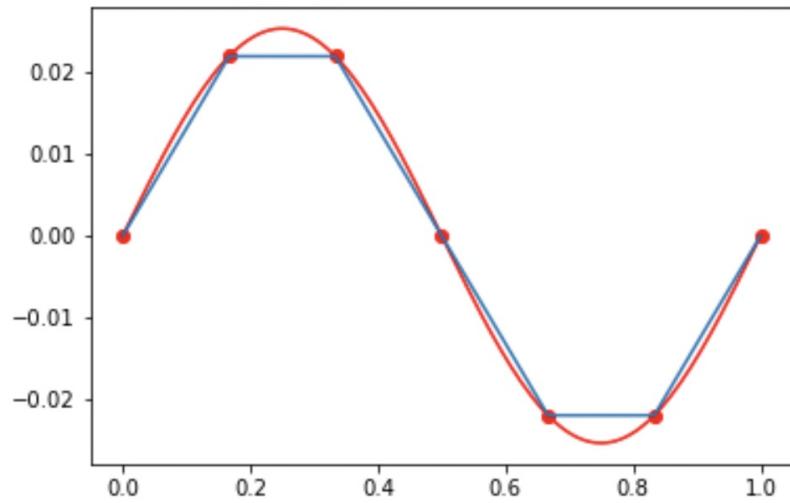
Approximation property of FEM:

$$\int_0^1 u' v_h = \int_0^1 f v_h \quad \forall v_h \in V_h$$

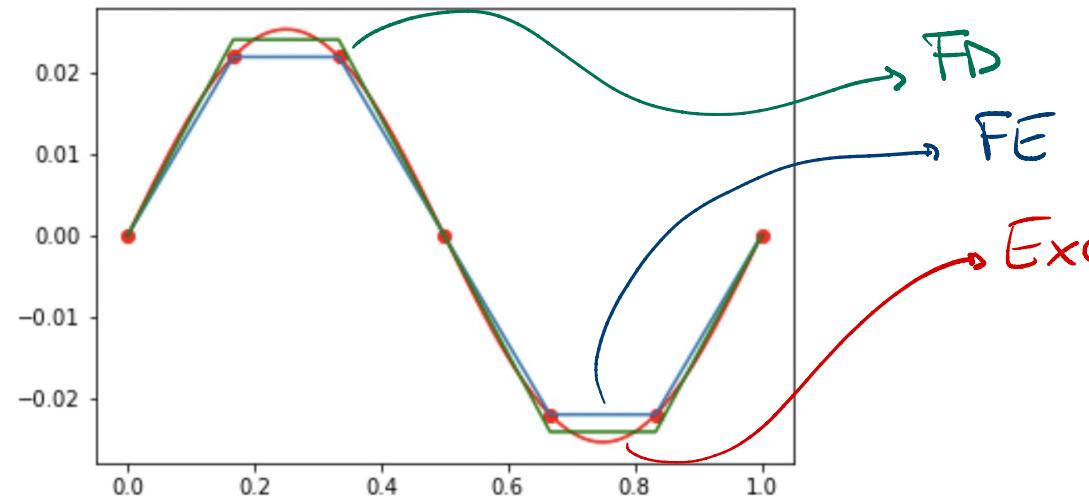
$$\int_0^1 u_h' v_h = \int_0^1 f v_h \quad \forall v_h \in V_h \subset V$$

$$\int_0^1 (u' - u_h') v_h = 0 \quad \forall v_h \in V_h \subset V$$

orthogonality of the error



In 1D, the FEM solution coincides with the exact solution at the nodes (while FD does not)



but it depends on the partial diff equation used  
IT IS EXACT FOR ANY SYMMETRIC POSITIVE DEFINITE SYSTEM  
(when you can construct norm with differential operators)

FD is a special case of FEM, but it is the "petals of a rose" of FEM

# Small recap of necessary concepts of functional analysis

Given  $V, W$  normed vector spaces (on the real field  $\mathbb{R}$ )

Define  $L(V, W)$  the space of linear functions

from  $V$  to  $W$ :

$$f \in L(V, W) \quad f: V \longrightarrow W$$

$$\text{i)} \forall u, v \in V, \forall \alpha, \beta \in \mathbb{R} \quad W \ni f(\alpha u + \beta v) = \alpha f(u) + \beta f(v)$$

$$\text{ii)} \forall u \in V \quad \|f(u)\|_W \leq \|f\|_* \|u\|_V$$

$$\|f\|_* := \sup_{0 \neq v \in V} \frac{\|f(v)\|_W}{\|v\|_V}$$

**Special cases:**

$$L(V, \mathbb{R}) = V^* \quad \text{the dual space of } V$$

$$f \in V^* \rightarrow \begin{aligned} f: V &\longrightarrow \mathbb{R} \\ v &\mapsto f(v) \equiv \underbrace{\langle f, v \rangle}_V \end{aligned} \quad \begin{matrix} \text{just} \\ \text{notation for now} \end{matrix} \quad v \in \mathbb{R}$$

you define

$$\|f\|_{V^*} = \|f\|_* = \sup_{0 \neq v \in V} \frac{|f(v)|}{\|v\|_V} = \sup_{0 \neq v \in V} \frac{|\langle f, v \rangle_V|}{\|v\|_V}$$

Bilinear forms:  $a: V \times V \rightarrow \mathbb{R}$

$$a(u, v) = \langle Au, v \rangle \quad \forall u, v \in V$$

where

$$A \in L(V, V^*)$$

$$A: V \rightarrow V^*$$

$$(Au)(v) \in \mathbb{R}$$

This is an object in  $V^*$

$a$ : is a bilinear form if it is a function from  $V \times V$  to  $\mathbb{R}$  which is linear in both arguments separately

$$\forall u, v, w, z \in V \quad \alpha, \beta, \gamma, \delta \in \mathbb{R}$$

Linearity in both arguments:

$$a(\alpha u + \beta v, \gamma w + \delta z) = \gamma \alpha a(u, w) + \gamma \delta a(u, z) + \dots$$

$$|a(u, v)| \leq \|A\|_* \|u\|_V \|v\|_V \quad \forall u, v \in V$$

Where

$$\|A\|_* := \sup_{0 \neq u \in V} \frac{\|Au\|_{V^*}}{\|u\|_V} = \sup_{0 \neq u, v \in V} \frac{|\langle Au, v \rangle|}{\|u\|_V \|v\|_V} = \sup_{u, v \in V} \frac{|a(u, v)|}{\|u\|_V \|v\|_V}$$

If  $\|A\|_*$  is  $< +\infty \rightarrow$  then  $a(u, v)$  is bounded

## Finite element method:

- Take the strong form of the PDE (1)
- Look for a solution in a vector space  $V$
- multiply (1) by test functions in  $V$  and integrate on the domain  $\Omega$  of the PDE
- integrate by parts until you cannot reduce any more the order of the differential op.

$\Rightarrow$  you end up with:

- i) a bilinear form in  $V \times V$
- ii) a linear operator in  $V^*$

## Example

$$\left\{ \begin{array}{l} -\Delta u = f \text{ in } \Omega \\ u = 0 \text{ on } \partial\Omega \end{array} \right. \quad \begin{array}{l} \text{(1)} \\ \text{new notation for II derivatives} \end{array}$$

(2) Take  $u$  in  $V$   
 $(V$  contains the condition  $u=0|_{\partial\Omega})$

(3)  $\int_{\Omega} \Delta u \cdot v = \int_{\Omega} f v - \operatorname{div}(\nabla u)$

(4) Integration by part

$$v = 0 \Big|_{\partial\Omega}$$

$$\int_{\Omega} \nabla u \cdot \nabla v - \int_{\partial\Omega} \nabla u \cdot n v = \int_{\Omega} f v$$

$$\int_{\Omega} \operatorname{div}(w) = \int_{\partial\Omega} w \cdot n$$

$$a(u, v) = \int_{\Omega} \nabla u \cdot \nabla v$$

$$f(v) := \int_{\Omega} f v$$

(5)

$$f \in V^*$$

Consider the case where  $V$  is Hilbert,

i)  $(u, v)_V$  is its scalar product

ii)  $\|u\|_V^2 := (u, u)_V$  of Poisson Problem

Example for FEM of Poisson Problem  
 $V = H_0(\Omega) := \{v \in L^2(\Omega), \nabla v \in L^2(\Omega), v|_{\partial\Omega} = 0\}$

$$(u, v)_V = \int_{\Omega} uv + \int_{\Omega} \nabla u \cdot \nabla v$$

$$\Rightarrow \|u\|_V = \|u\|_1 := \left( \int_{\Omega} u^2 + \int_{\Omega} |\nabla u|^2 \right)^{\frac{1}{2}}$$

A scalar product is a BILINEAR FORM

Abstract setting -  $V$  Hilbert (not only  $H_0$ , but any Hilbert)

Lax Milgram Lemma:

i) <sup>let</sup>  $a$ : bilinear, bounded operator

ii)  $a$ : coercive

$$1) a(u, v) \leq \|A\|_* \|u\|_V \|v\|_V \quad \forall u, v \in V$$

$$2) a(u, u) \geq \alpha \|u\|_V^2 \quad \forall u \in V$$

Then  $\nexists f \in V^*$ ,  $\exists!$  solution  $u \in V$  to

$$① a(u, v) = f(v) = \langle f, v \rangle \quad \forall v \in V$$

this it has the property:

$$\text{And } \|u\| \leq \frac{1}{\alpha} \|f\|_*$$

solution is BOUNDED  
with respect to "state"

Since ① is true for every  $v \in V$ , and  $u \in V$   
then:

$$d\|u\|^2 \leq a(u, u) = \langle f, u \rangle = f(u) \leq \|f\|_* \|u\|$$

$$\rightarrow \|u\| \leq \frac{1}{2} \|f\|_*$$

### Proof of Lax Milgram

2 ingredients

1) Rietz Theorem

2) Contraction theorem (fixed points)

Rietz theorem(1): (Rietz operator  $\mathcal{T}: V^* \rightarrow V$ )

$\forall f \in V^*, \exists! \tilde{f} \in V$  s.t.

$$i) (\tilde{f}, v)_V = \langle f, v \rangle = f(v)$$

$$ii) \|\tilde{f}\|_V = \|f\|_* \quad \tilde{f} = \mathcal{T}f$$

### Contraction Theorem

Let  $T: V \rightarrow V$  be a contraction, i.e:

$\exists L < 1$  s.t.

$$\|T(u) - T(v)\|_V \leq L \|u - v\|$$

Then  $\exists! \varphi$  fixed point of  $T$ :  $T(\varphi) = \varphi$

Let  $u^0 \in V$ , and define  $u^{k+1} = T(u^k)$

$$\rightarrow \|u^{k+1} - u^k\| = \|T(u^k) - T(u^{k-1})\| \leq L \|u^k - u^{k-1}\|$$

$$\rightarrow \|u^{k+1} - u^k\| \leq L^k \|T(u^0) - u^0\|$$

$$|L| < 1 \rightarrow \|u^{k+1} - u^k\| \rightarrow 0$$

$$\Rightarrow \exists! u \text{ s.t. } u^k \rightarrow u$$

### Proof of Lax Milgram

$$\langle Au, v \rangle = \langle f, v \rangle \quad \forall v \in V$$

$$a(u, v) = f(v)$$

$$\langle \underbrace{Au - f}_{V^*}, v \rangle \quad \forall v \in V$$

Rietz repr. theorem

$$(z(Au - f), v)_V \quad \forall v \in V$$

Construct a fixed point iteration (for  $\rho \in \mathbb{R}$ )

$$T_\rho(v) = v - \rho z(Av - f)$$

$$T_\rho: V \rightarrow V$$

How we see this is a contraction!

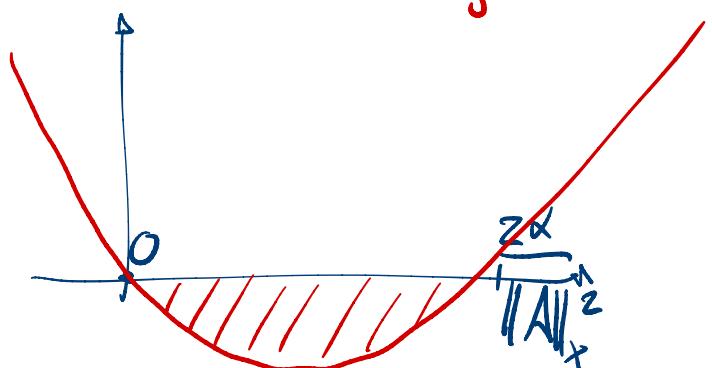
$$\|T_\rho(v) - T_\rho(w)\|_V^2 = \|v-w - \rho z A(v-w)\|_V^2$$

$$\begin{aligned}
 \|T_p(v) - T_p(w)\|_V^2 &= \|v-w\|_V^2 + g^2 \left\| z A(v-w) \right\|_V^2 - 2g \langle z A(v-w), v-w \rangle_V \\
 &= \|v-w\|_V^2 + g^2 \|A(v-w)\|_{V^*}^2 - 2g \langle A(v-w), v-w \rangle \\
 &\leq \|v-w\|_V^2 + g^2 \|A\|_*^2 \|v-w\|_V^2 - 2g \alpha \|v-w\|_V^2
 \end{aligned}$$

$$\|T_g(v) - T_g(w)\|^2 \leq \left( g^2 \|A\|_*^2 - 2g\alpha + 1 \right) \|v-w\|_V^2$$

$\underbrace{\phantom{\dots}}_{< 1}$

Can we choose  $g$  s.t.  $\underbrace{g^2 \|A\|_*^2 - 2g\alpha + 1}_{g \left( g - \frac{2\alpha}{\|A\|_*^2} \right)} < 1$ ?



If  $g \in (0, \frac{2\alpha}{\|A\|_*^2})$  then  $T$  is a contraction

$$\Rightarrow \exists ! u \mid T(u) = u \quad \Rightarrow$$

$$\begin{aligned}
 u &= u - g z(Au - f) \\
 \Rightarrow Au &= f
 \end{aligned}$$

## Galerkin methods:

Choose  $V_h = \text{span} \{v_i\}_{i=0}^{N-1}$  st.  $V_h \subset V$

$\rightarrow$  Find  $u_h \in V_h \Rightarrow u_h = \sum_{j=0}^{N-1} u_j^j v_j$   
st.

$$\langle A u_h, v_i \rangle = \langle f, v_i \rangle = f(v_i) \quad i=0, \dots, N-1$$

$$\sum_j \langle A u_j^j v_j, v_i \rangle = \langle f, v_i \rangle = f(v_i)$$

$$A_u = f \quad A_{ij} := \langle A v_j, v_i \rangle \quad f_i := f(v_i)$$

$$1): \quad \langle A u_h, v_h \rangle = \langle f, v_h \rangle \quad \forall v_h \in V_h$$

$$\langle A u, v_h \rangle = \langle f, v_h \rangle \quad \forall v_h \in V_h$$

(implying that the error is a conjugate w.r.t  $v_h$ )

$$\langle A(u - u_h), v_h \rangle = 0 \quad \forall v_h \in V_h$$

The proof above uses Cea's Lemma (orthogonality of error)

Coercivity:

$$\begin{aligned} \alpha \|u - u_h\|^2 &\leq \langle A(u - u_h), u - u_h \rangle = \\ &= \langle A(u - u_h), u - v_h \rangle \leq \|A\|^* \|u - u_h\| \|u - v_h\| \end{aligned} \quad \forall v_h \in V_h$$

$$\|u - u_h\|_V \leq \frac{\|A\|_*}{\alpha} \inf_{v_h \in V_h} \|u - v_h\|_V$$

Infimum

A priori error estimate

Trick: Make  $\text{dist}(V_h, V)$  small!

$$\text{dist}(V_h, V) := \sup_{u \in V} \inf_{v_h \in V_h} \|u - v_h\|$$

if you can  
 $\text{dist}(V_h, V) \leq \epsilon$   
 $\Rightarrow V_h$  is  $\epsilon$ -approximable  
 by  $V$

For Lagrangian Finite elements of order  $k$

We bound  $u - I(u)$  in terms of Sobolev norms of  $u$ :

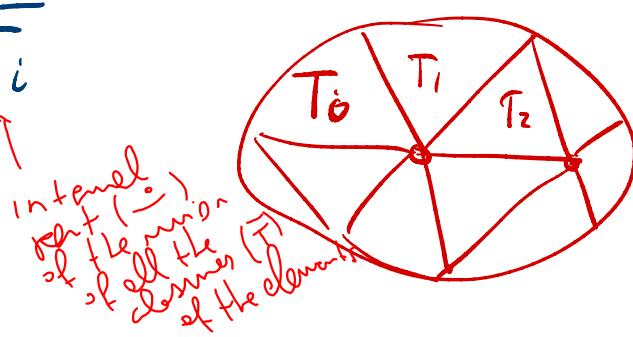
If  $I(u)$  is Lagrange interpolation of  $u$ , then

$$\|u - I(u)\|_{m,T} \leq C h_T^{k+1-m} |u|_{k+1,T}$$

measure  
of the  
element

on each "element"  $T$  of  $\mathcal{S}_h$ , discretization of  $\mathcal{S}$ . 1D:  $\mathcal{S} = [0,1]$   $T = [x_i, x_{i+1}]$

$$\mathcal{S}_h = \bigcup_{i=1}^n T_i$$



$$h_T = (x_{i+1} - x_i)$$

$$\|u\|_{\alpha,T} := \left( \sum_{|\alpha|=d} \int_T ((D^\alpha u)^2) \right)^{\frac{1}{2}}$$

weak derivative  
of  $\alpha$  in  $m$

$\alpha$  is multi index

$\alpha_1, \alpha_2, \dots, \alpha_d$

$$\|u\|_{\alpha,T} := \left( \sum_{|\alpha| \leq d} |u|_{\alpha,T}^2 \right)^{\frac{1}{2}}$$

in dimension

$$\text{and } |\alpha| = \sum_{i=1}^d |\alpha_i|$$

$$\text{in 1D: } |u|_{\alpha,T} := \left( \int_T (u^{(\alpha)})^2 \right)^{\frac{1}{2}}$$

$$V = \mathbb{R}^n \quad V_h = \text{span} \left\{ v_i \right\}_{i=0}^{m-1}$$

Exercise

We try to approx a finite dim space  $V$ , with the smaller finite dim space  $V_h$

$$m \leq n \quad v_i \in \mathbb{R}^n (\equiv V)$$

Apply Galerkin to this problem:

$$(u, v) := \sum_{i=0}^{n-1} u^i v^i$$

$$a : V \times V \rightarrow \mathbb{R}$$

$$\|u\| := \left( \sum_{i=0}^{n-1} (u^i)^2 \right)^{\frac{1}{2}}$$

$$\underline{A} : V = \mathbb{R}^n \longrightarrow V^* (\equiv V) = \mathbb{R}^n$$

A is a  $\mathbb{R}^{n \times n}$  matrix,  $v_i$  are  $\mathbb{R}^n$  vectors

Galerkin approx.

$$\langle \underline{A} u^j v_j, v_i \rangle = \langle f, v_i \rangle \quad i = 0, \dots, n-1$$

Should give me a  $\mathbb{R}^{n \times n}$  matrix A

$$A_{ij} = \langle \underline{A} v_j, v_i \rangle = \underline{v}_i^T \underline{A} \underline{v}_j$$

$$F_i = \langle f, v_i \rangle = \underline{v}_i^T \underline{f}$$

$$V_{id} := (v_d)_i$$

$$V \in \mathbb{R}^{n \times m}$$

Matrix containing the basis  $\{v_i\}_{i=0}^{m-1}$

$$V^T A V v_0 = V^T f$$

$$A_0 = F$$

$$V \in \mathbb{R}^{n \times m} \quad A \in \mathbb{R}^{n \times n} \quad f \in \mathbb{R}^n$$

$$v \in \mathbb{R}^m$$

If we manage to construct  $V^T A V = I_d$   
 Then  $v = V^T f$

nothing  
left to  
be inserted  
here

$$m \leq V V^T f$$

$$Vv \approx u \quad m=0$$

$$\text{Start with } q_0 = r_0 = A u_0 - f = -f$$

$$q_0^T A q_0 := \beta_0 \rightarrow v_0 = \frac{f_0}{\beta_0}$$

$$\rightarrow r_0^T A v_0 = 1$$

construct  $r_{k+1}$  s.t.

$$v_{k+1}^T A v_j = 0 \quad \forall j \leq k$$

and

$$v_{k+1}^T A v_{k+1} = 1$$

And you  
get exactly  
the

$\Rightarrow$  at most after  $n$  iterations  $u = V^T V f$

CG  
method

# Lecture 29 - LH

## FEM / FDM 2D

14.01.2021

One problem formulation is

$$1) \begin{aligned} -\Delta u &= f && \text{in } \Omega \\ u &= 0 && \text{on } \partial\Omega \end{aligned}$$

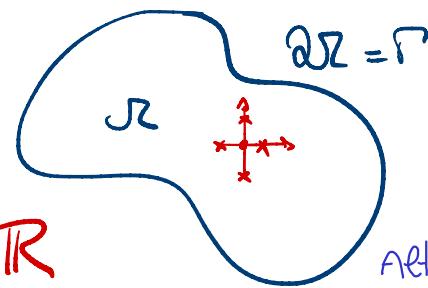
LAPLACIAN

$$\downarrow \quad -\Delta u := -\frac{\partial^2 u}{\partial x_0^2} - \frac{\partial^2 u}{\partial x_1^2} = \sum_{d=0}^{N-1} -\frac{\partial^2 u}{\partial x_d^2}$$

$$u : \Omega \rightarrow \mathbb{R}$$

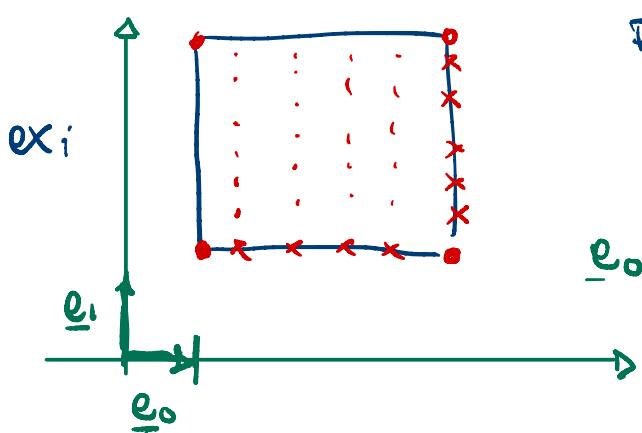
sometimes:

$$\begin{aligned} \text{Alternatively,} \\ \text{another problem} \\ \text{that can appear is} \\ -\Delta u + u = f \\ 2) \frac{\partial u}{\partial n} = 0 \quad \text{on } \partial\Omega \end{aligned}$$



$N=2$

For FD  $\Omega$  has to be "simple" and "regular"



Fix a grid of points  $X_i \in \Omega \subset \mathbb{R}^2$

$$\Omega = [a, b] \times [c, d]$$

$$\underline{e}_0 := (1, 0) \quad \underline{e}_1 := (0, 1)$$

$$x = \underline{e}_0 x^0 + \underline{e}_1 x^1 = \underline{e}_i x^i$$

FD: along  $e_0$ : choose  $n$  points as  $\{\varphi_i\}_{i=0}^{n-1} := [a, a+\dots, \dots, b]$

along  $e_1$ : choose  $m$  points as  $\{\psi_d\}_{d=0}^{m-1} := [c, \dots, \dots, d]$

construct  $\underline{X}_{di} \in \mathbb{R}^2$  ( $n \times m$  points)  $d \in \{0, \dots, m-1\}$

$(\underline{X}_{di})_d \rightarrow d$  comp. of  $\underline{X}_{di}$   $i \in \{0, \dots, n-1\}$

$$\underline{X}_{di} := (\varphi_i, \psi_d)$$

You want same # of pts on both X and Y direction, for 2 errors:  
1) multiply 2)  
If I flatten the things, I get

$$\text{flat}(\underline{X}) \in \mathbb{R}^{m \cdot n \cdot 2}$$

The X object has elements

$$\underline{X}_{did} = \begin{cases} \varphi_i & \text{if } d=0 \\ \psi_d & \text{if } d=1 \end{cases}$$

$$\begin{matrix} \{\underline{X}_{did}\} \\ d=0, \dots, m-1 \\ i=0, \dots, n-1 \\ d=0, 1 \end{matrix}$$

at  $m \times n$  points  
in directions (2)

$$X \in \mathbb{R}^m \times \mathbb{R}^n \times \mathbb{R}^N \quad \text{shape}(X) = (m, n, N)$$

$$X_{d,j} \in \mathbb{R}^N \quad \underline{x}$$

$$X_{d,j,d} \in \mathbb{R} \quad x_i$$

$$X_d \in \mathbb{R}^{n \times N}$$

*✓ net is  
only on the  
grid*

Now we assume we know  $u$  on  $X$  only  
(on all  $n \cdot m$  points  $X_{di} \in S^2$ )

$$u_{di} := u(X_{di}) \quad \leftarrow \begin{array}{l} u \text{ is a ft/that if you pass it} \\ \text{a "matrix" (matrix element)} \\ \text{split into a matrix (matrix ele-} \\ \text{ment)} \end{array}$$

We apply CFD to  $u_{di}$  on both directions:

$$\text{CFD}^2 - 1D: \quad -v_i'' \sim \frac{-v_{i-1} + 2v_i - v_{i+1}}{h^2}$$

(along  $i \rightarrow e_0$ )

$$\text{CFD}^2 - 1D \quad -w_d'' \sim \frac{-w_{d-1} + 2w_d - w_{d+1}}{h^2}$$

along  $d \rightarrow e_1$

$$v_i := u_{di} \quad -\Delta u_{di} \approx \frac{1}{h^2} [4u_{di} - u_{d-1,i} - u_{d+1,i} - u_{d-1} - u_{d+1}]$$

$$w_d := u_{di}$$

$$-\tilde{\Delta} u_{di} = f_{di} = f(X_{di})$$

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

*because I have  
2 directions  
 $d = 0$  and  $1$*

$$Au = f$$

Construct

$$A_{\alpha i \beta j} \text{ s.t. } \sum_{\beta, j} A_{\alpha i \beta j} \mu_{\beta j} = -\tilde{\Delta} u_{\alpha i}$$

$$A_{\alpha i \beta j} = \begin{cases} 4 & \beta = \alpha \\ -1 & \beta = \alpha - 1 \\ -1 & \beta = \alpha + 1 \\ -1 & \beta = \alpha \\ -1 & \beta = \alpha \\ 0 & \text{all other cases} \end{cases} \quad \begin{array}{lll} \beta = \alpha & j = i \\ \beta = \alpha - 1 & j = i \\ \beta = \alpha + 1 & j = i \\ \beta = \alpha & j = i - 1 \\ \beta = \alpha & j = i + 1 \end{array}$$

$$A_{\alpha i \beta j} \mu_{\beta j} = f_{\alpha i}$$

↓ reshape

$$\sum_j \tilde{A}_{Ij} \tilde{\mu}_j = \tilde{f}_I$$

$$I, j \in [0, \dots, m \cdot n - 1]$$

$$\tilde{\mu}_j = \sum_i (\tilde{A}^{-1})_{ji} \tilde{f}_i$$