

# Least Squares



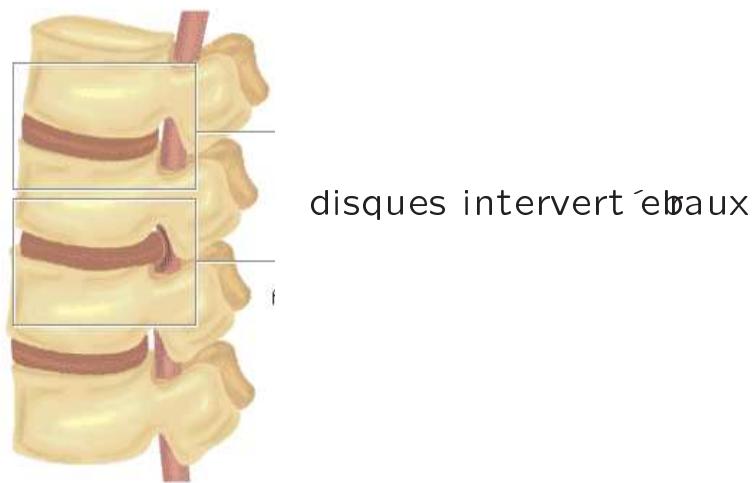
## Numerical Analysis

Profs. Gianluigi Rozza- Luca Heltai

2019-SISSA mathLab Trieste

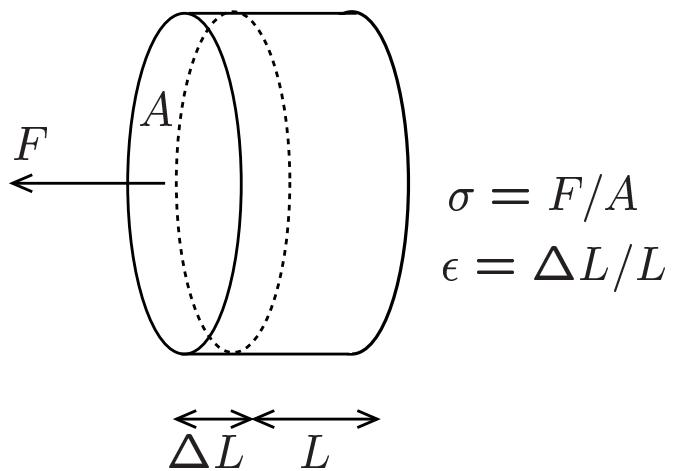
# Examples and motivations

**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).



disques intervertébraux

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

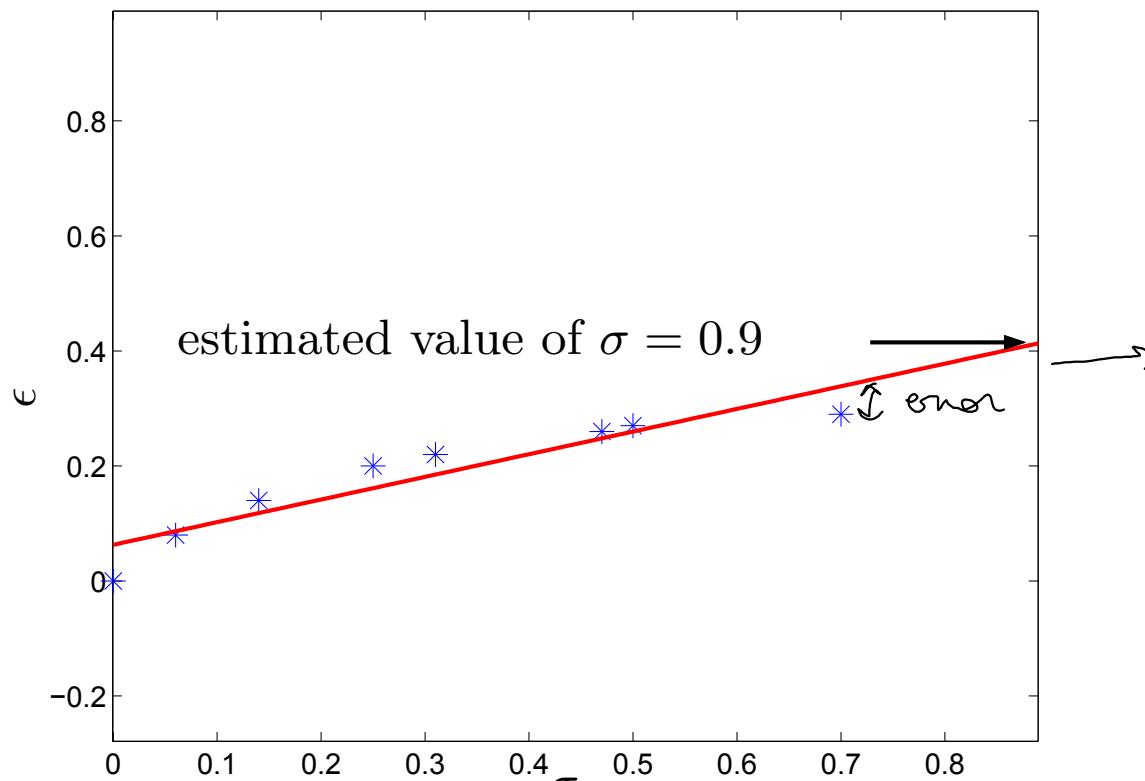


$$\sigma = F/A$$

$$\epsilon = \Delta L/L$$

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

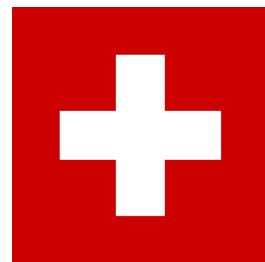
By the method of least squares, we get that best approximation for the data is  $p(x) = 0.3938x - 0.0629$ . The approximation can be used (called *a linear regression*) to estimate  $\epsilon$  when  $\sigma = 0.9$  MPa: We get  $p(0.9) \simeq 0.4$ .



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

! If you want to approx this data (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):

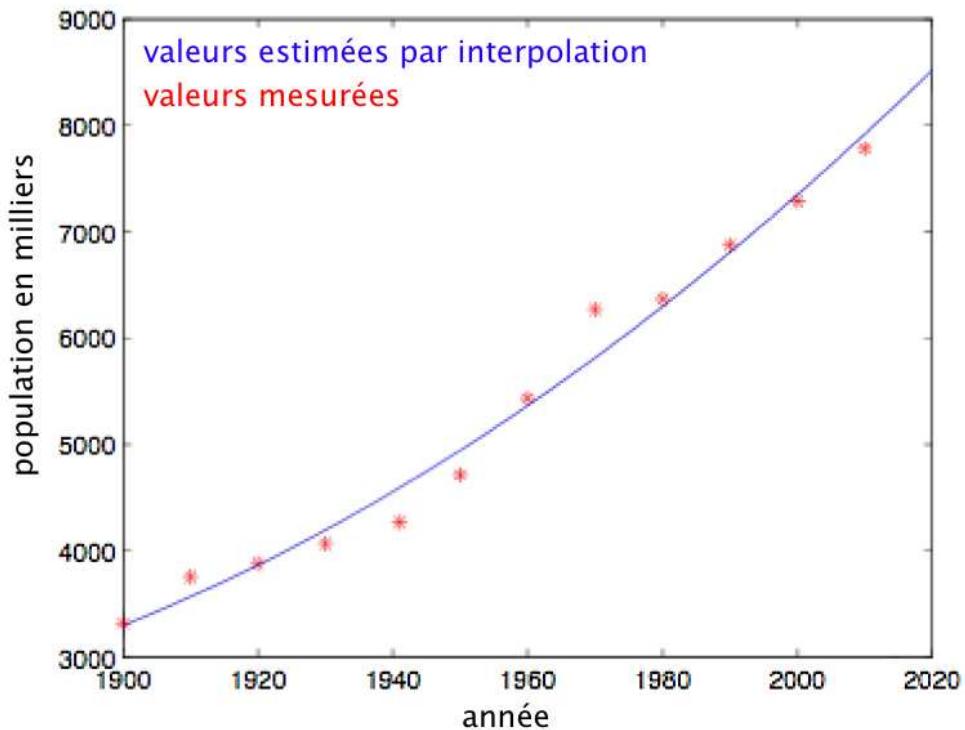


year	1900	1910	1920	1930	1941	1950
population	3315	3753	3880	4066	4266	4715
year	1960	1970	1980	1990	2000	2010
population	5429	6270	6366	6874	7288	7783

- 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  $p(x) = 0.15x^2 - 549.9x + 501600$ .

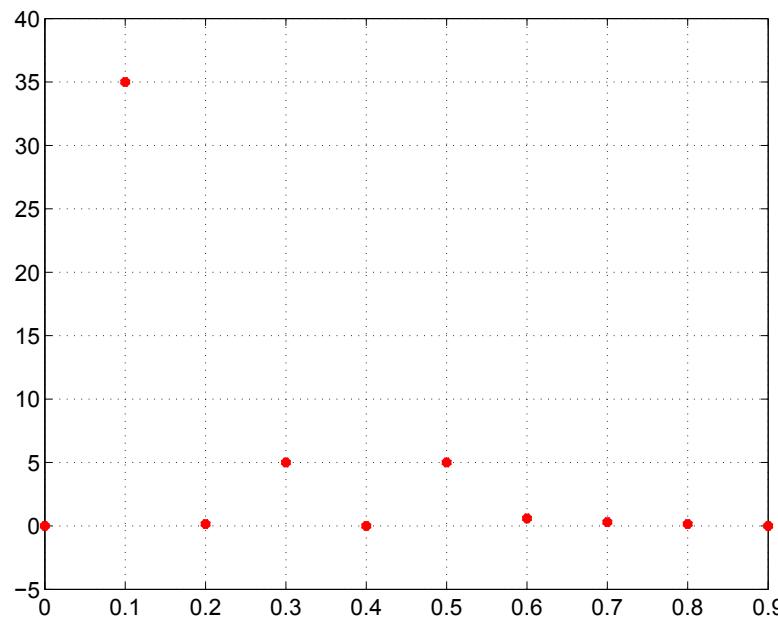
↑ II od poly. LEAST SQUARES



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

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

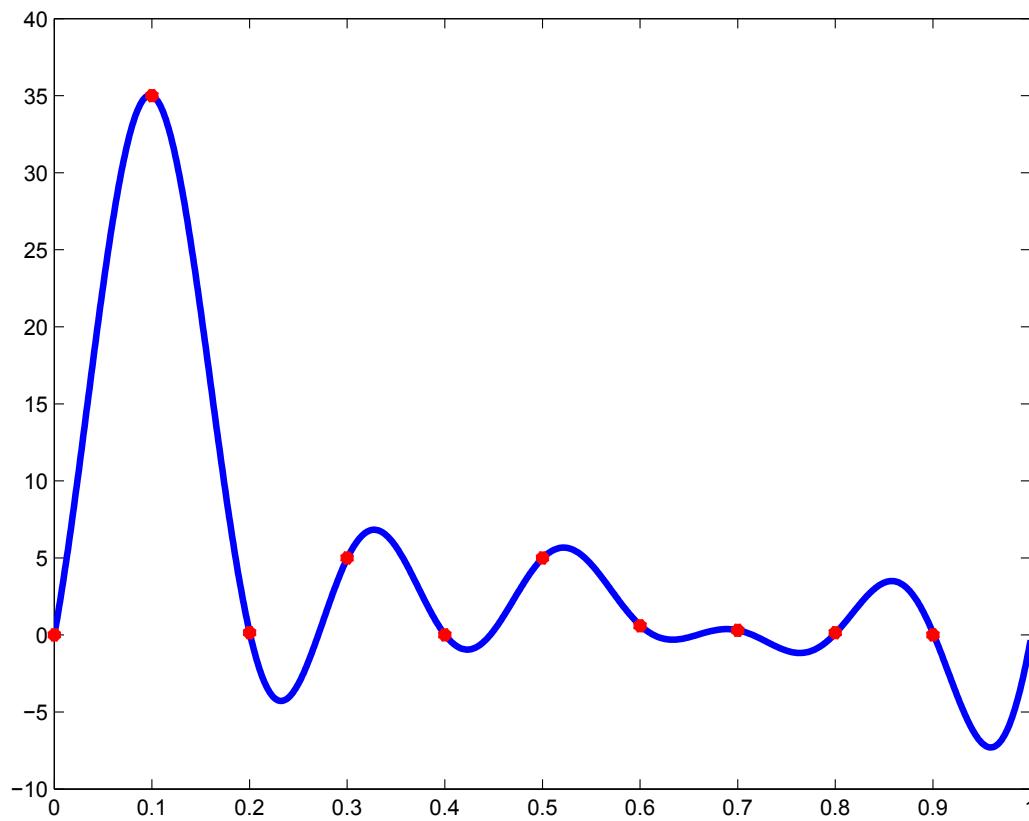
$i$  is the imaginary unit. We can calculate a vector of the coefficients  $\mathbf{c} = [c_1, \dots, c_{10}]^T$  by the FFT algorithm (Fast Fourier Transform, Cooley & Tuckey, 1965), implemented in Matlab/Octave by the **fft** command.

# 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];
```

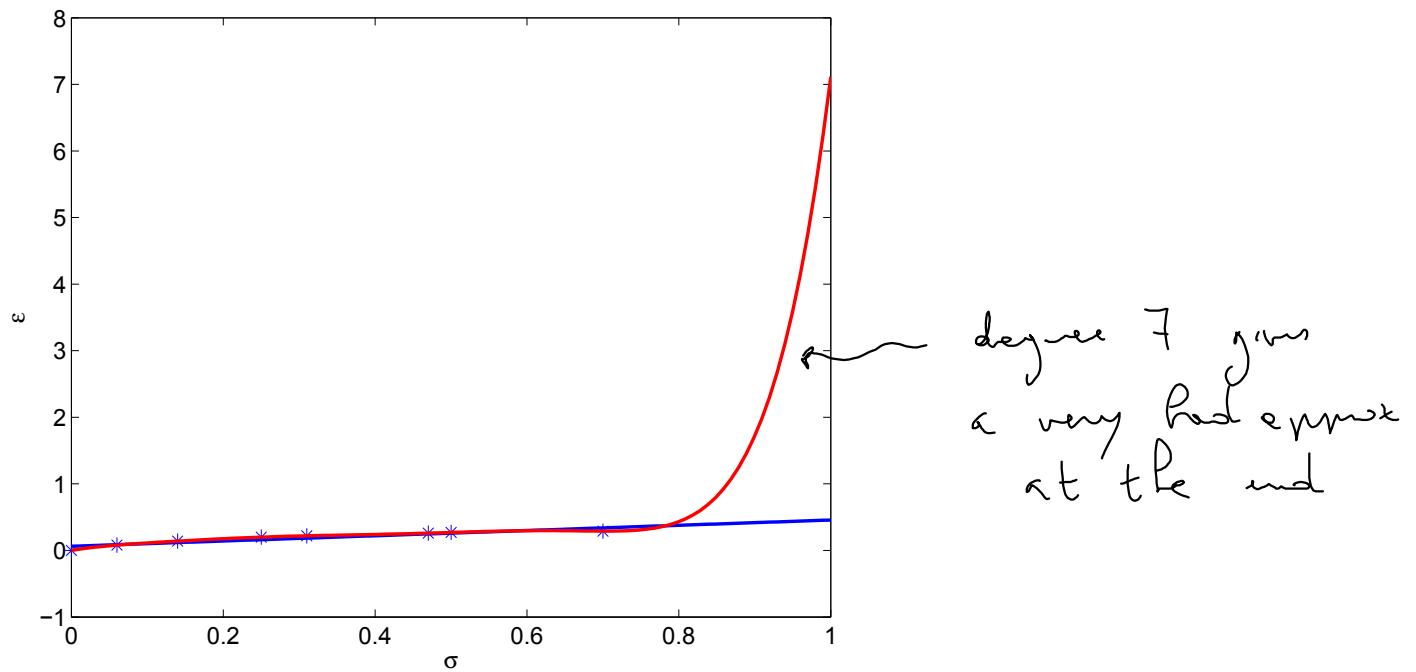
We want to extrapolate the value of  $\epsilon$  for  $\sigma = 0.4$ . We consider two ways:

- 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

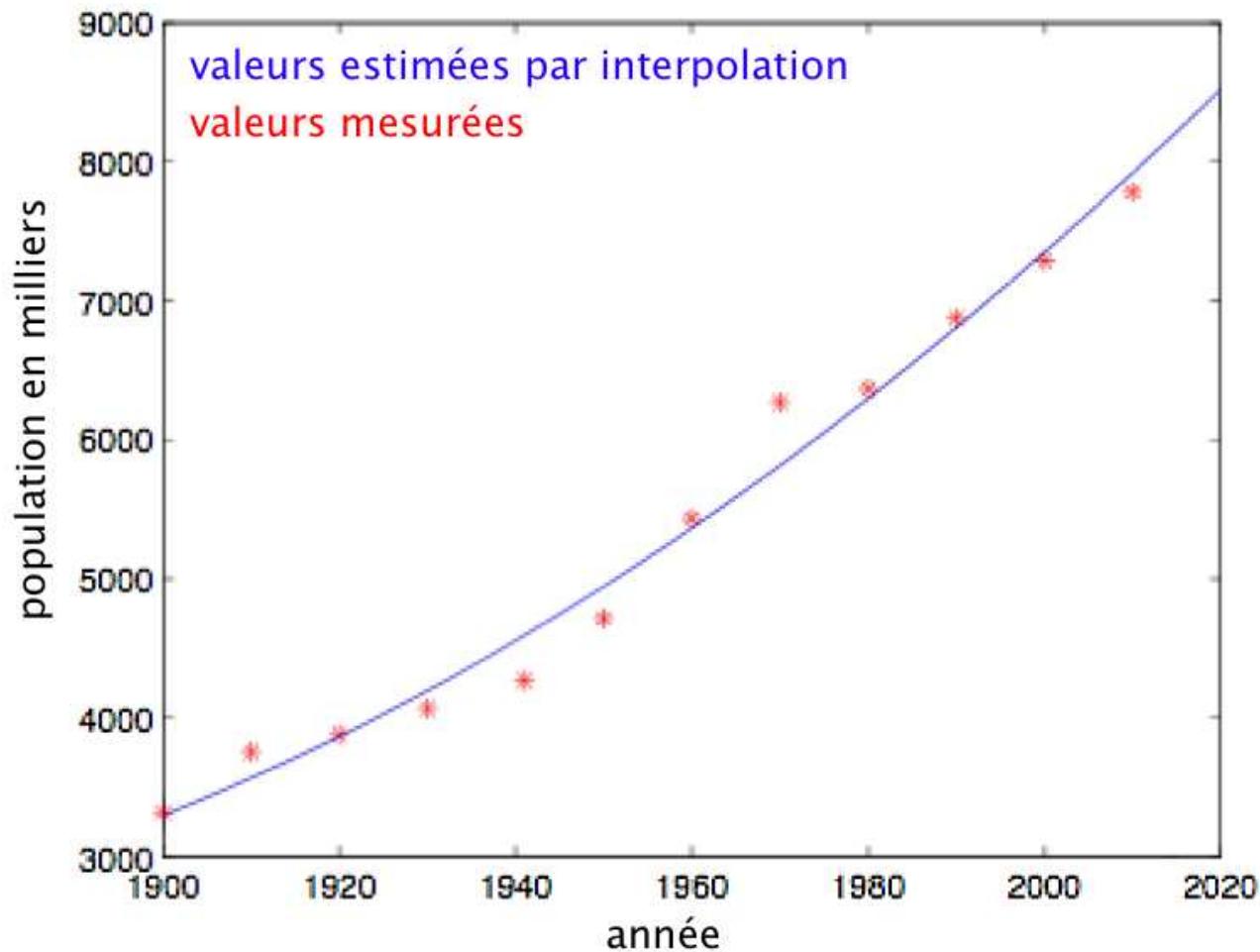
- The value obtained by  $\Pi_7$  (172.21%) is unrealistic
- On the contrary, the value obtained with the regression line is more appropriate to compute the value at  $\sigma = 0.9$ .

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');
```

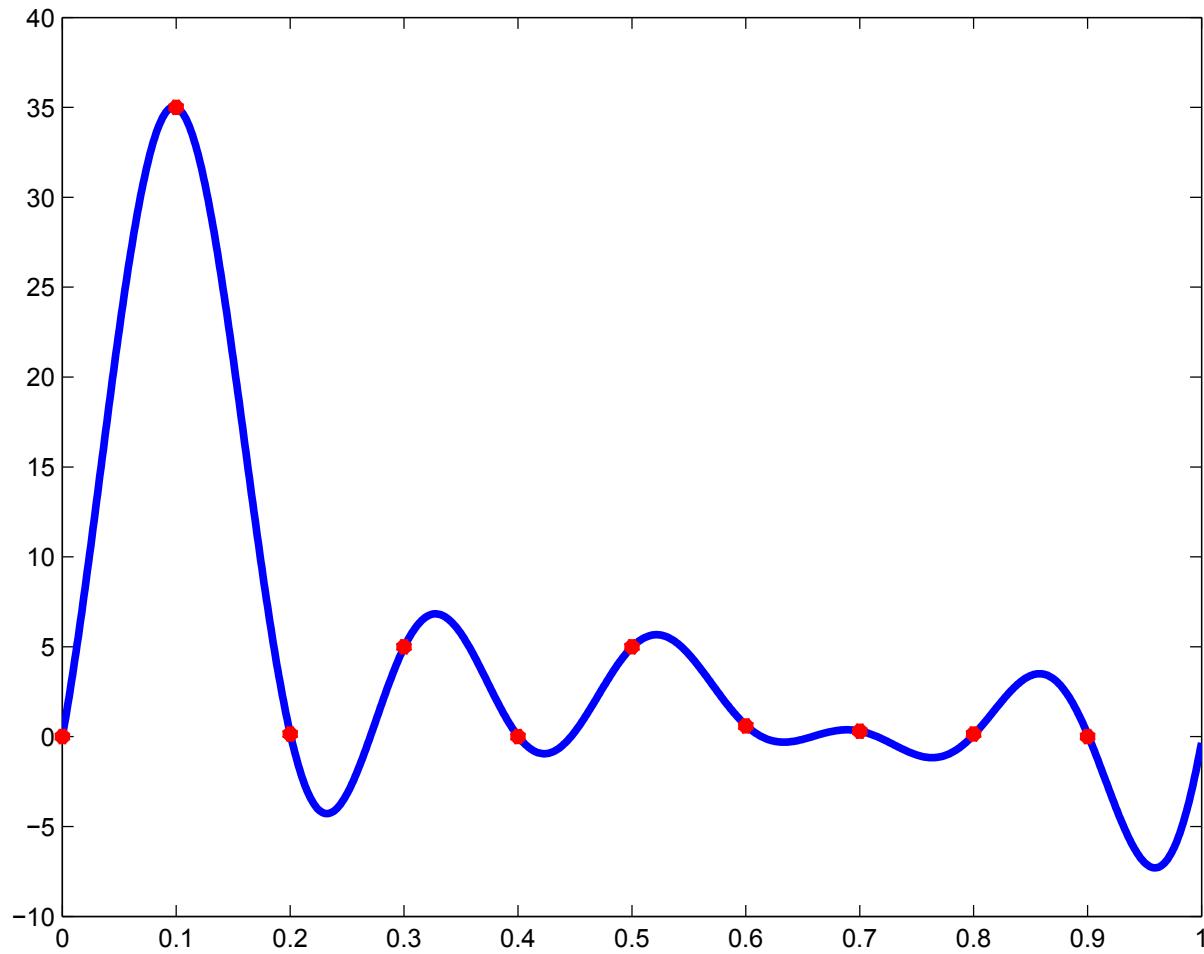


**Example. Carotide flow rate** We define 10 couples of data (time vs flowrate) with the vectors **t** and **deb**. With the help of the command **interpft**, we compute the value of the interpolating function at 1000 equidistributed points in the periodicity interval [0,1]:

```

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

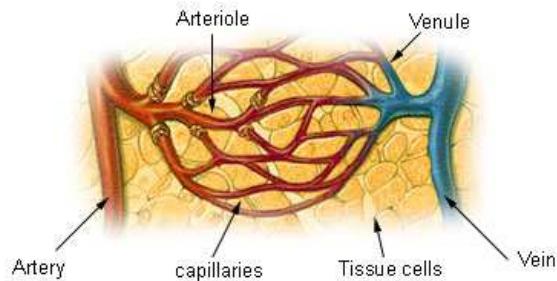
Profs. Gianluigi Rozza - Luca Heltai

2019-SISSA mathLab Trieste

# Examples and motivation

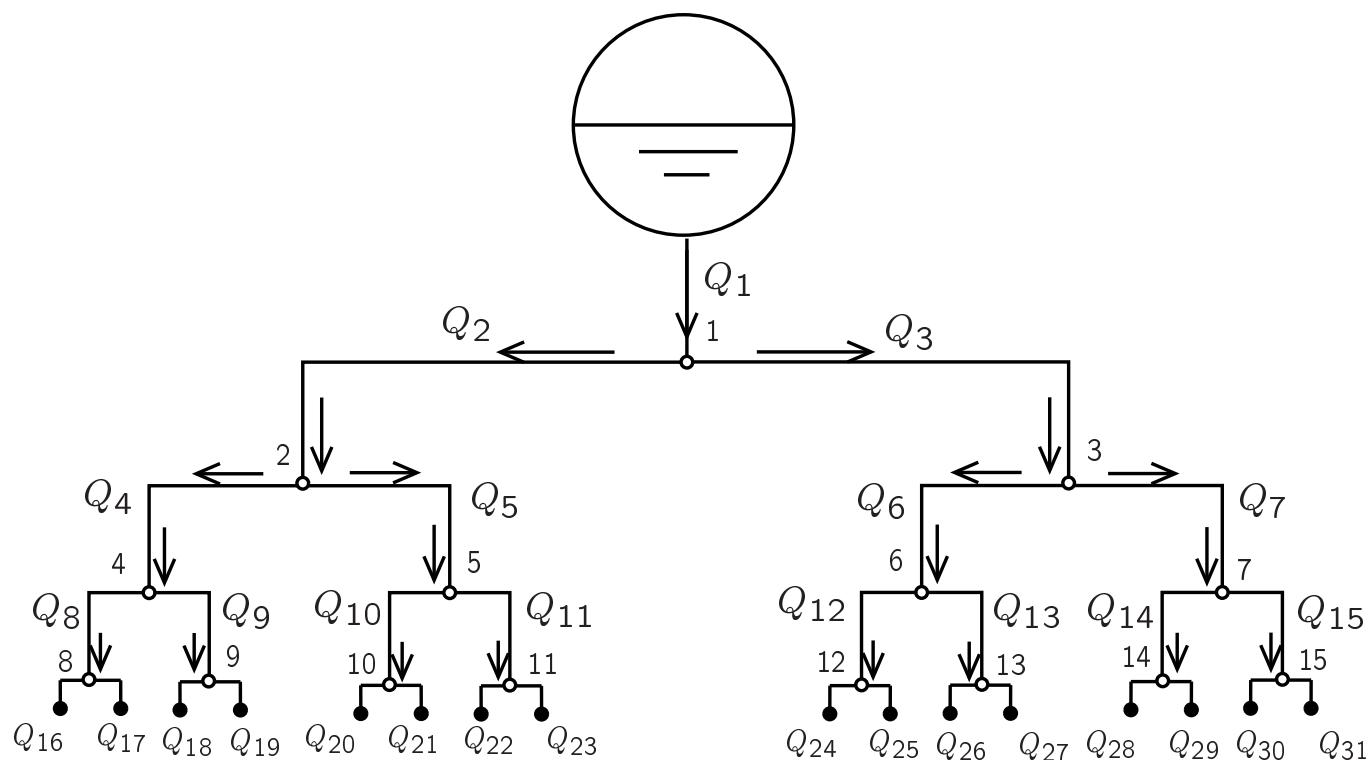
Capillaries are microscopic blood vessels, the smallest part of the circulatory system. The capillaries are grouped in networks called capillary beds, made of 10 to 100 capillaries, depending on the tissue.

- 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.

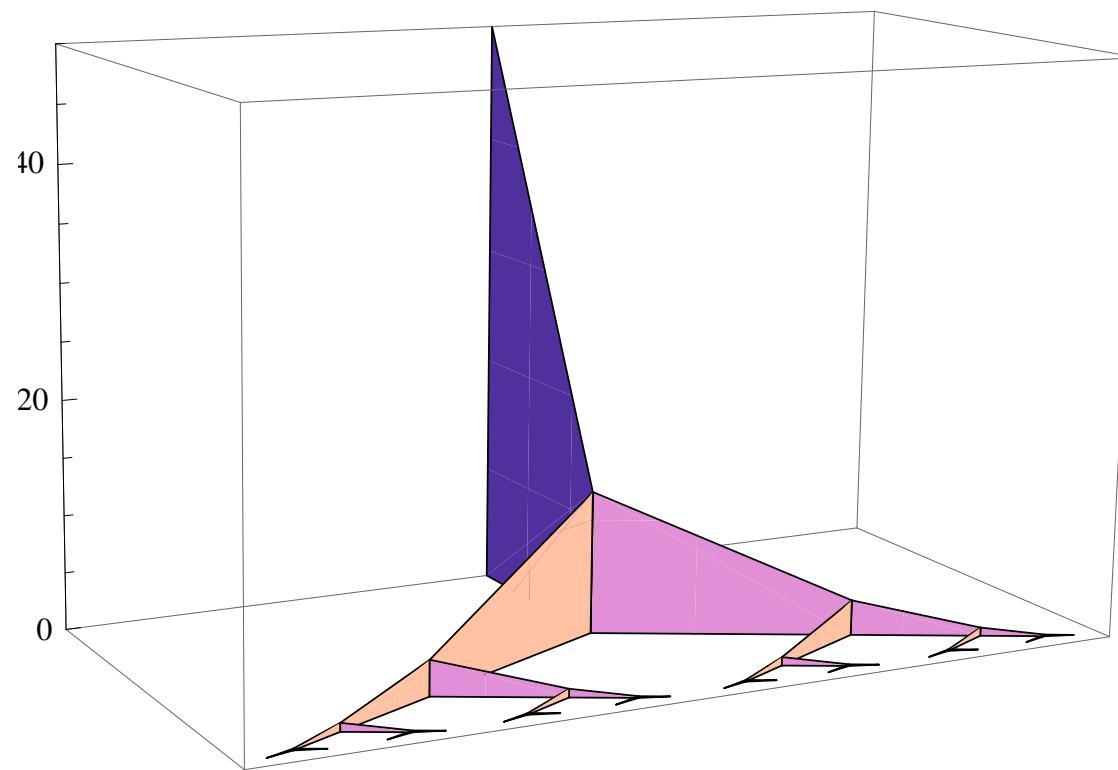
The solution is given by the vector

$$\mathbf{p} = [12.46, 3.07, 3.07, 0.73, 0.73, 0.73, 0.73, 0.15, 0.15, 0.15, 0.15, 0.15, 0.15, 0.15, 0.15]^\top.$$

Once we have the pressures, we can compute, using the relation (I), the flows:

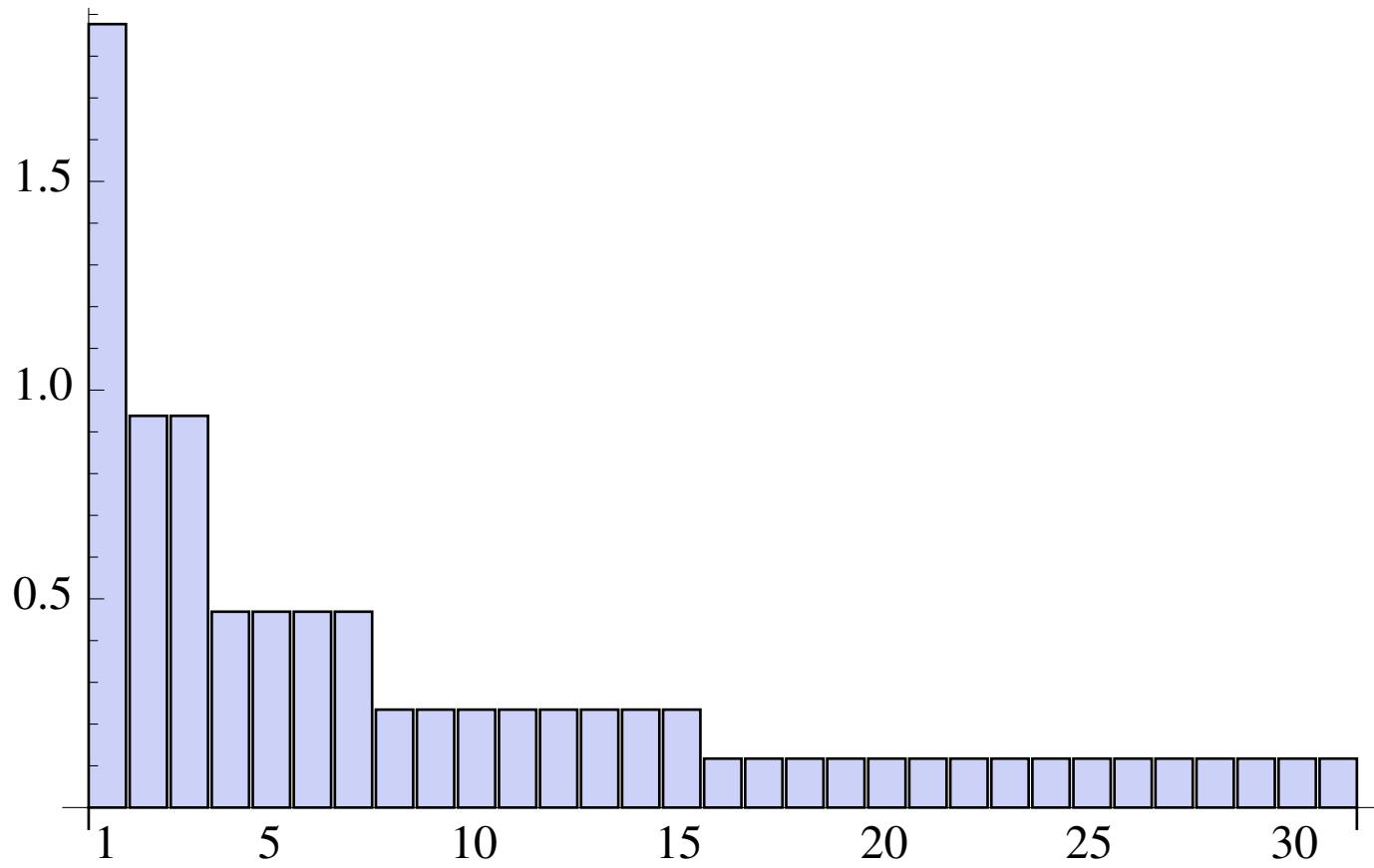
$$\begin{aligned}
 Q_1 &= 1.88 \\
 Q_{2,3} &= 0.94 \\
 Q_{4,\dots,7} &= 0.47 \\
 Q_{8,\dots,15} &= 0.23 \\
 Q_{16,\dots,31} &= 0.12
 \end{aligned}$$

Linear distribution of pressures in the capillary bed computed from the pressures at every node (solution of the system)

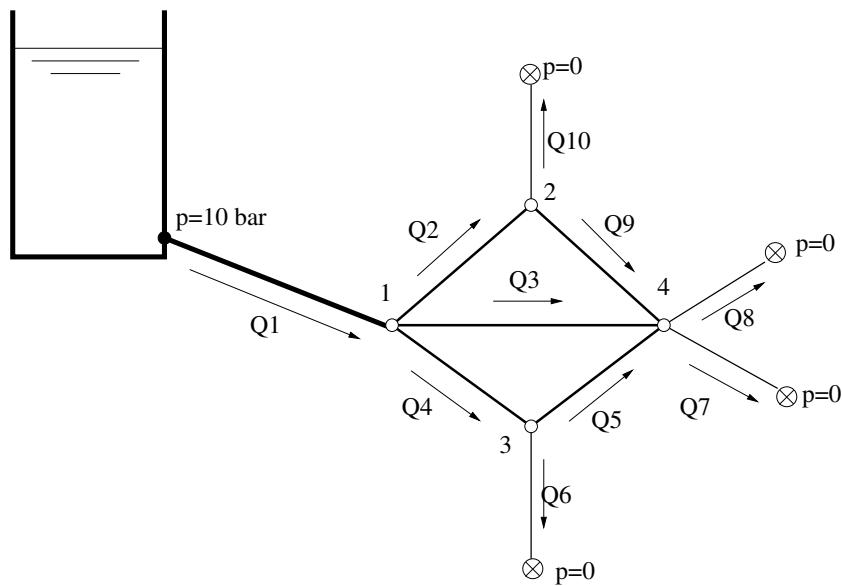


$\leftarrow$  system is dumping energy  
 we have a drop in pressure

## 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  $\text{m}^3/\text{s}$ ):

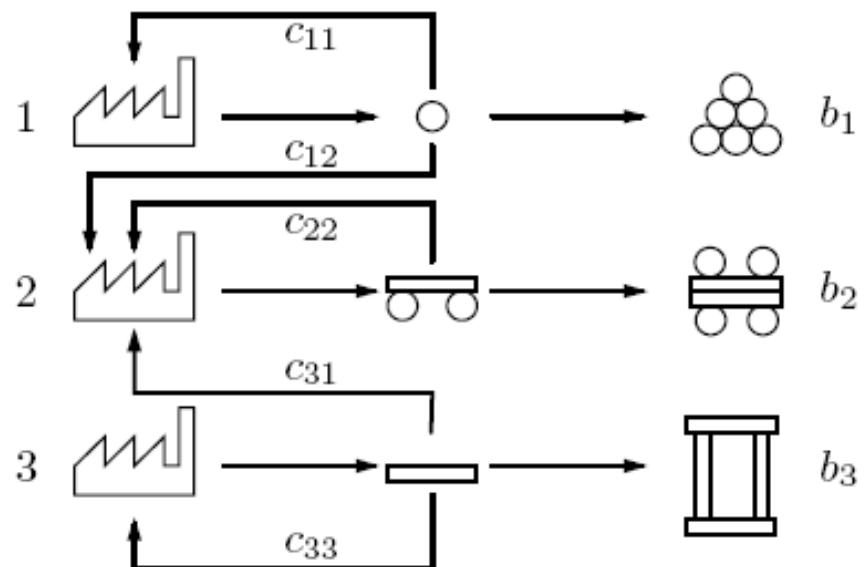
Q1	0.3766	Q2	1.0640	Q3	0.1638
Q4	0.1064	Q5	0.0106	Q6	0.0958
Q7	0.0924	Q8	0.0924	Q9	0.0106
Q10	0.0958				

**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  $i = n - 1, n - 2, \dots, 2, 1$

$$x_i = \frac{1}{u_{ii}} \left( y_i - \sum_{j=i+1}^n u_{ij}x_j \right).$$

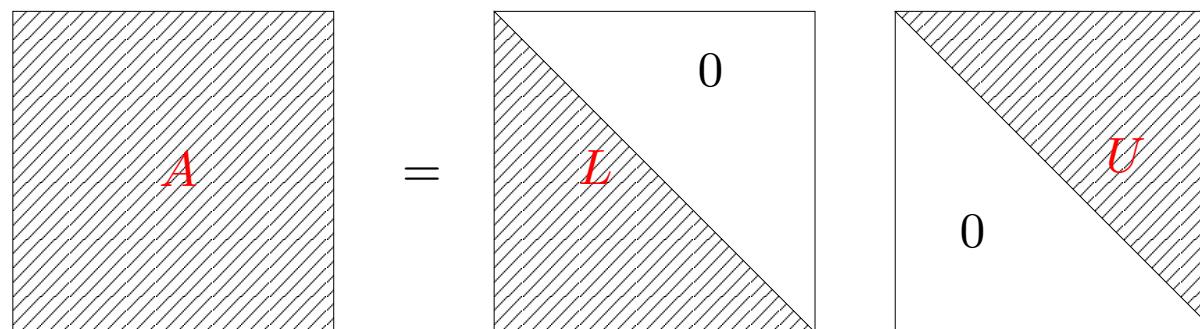
This algorithm is called *backward substitutions algorithm*. The cost is, once again,  $n^2$  operations.

# 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

$$A^{(k)} \mathbf{x} = \mathbf{b}^{(k)}, \quad 1 \leq k \leq n, \quad (7)$$

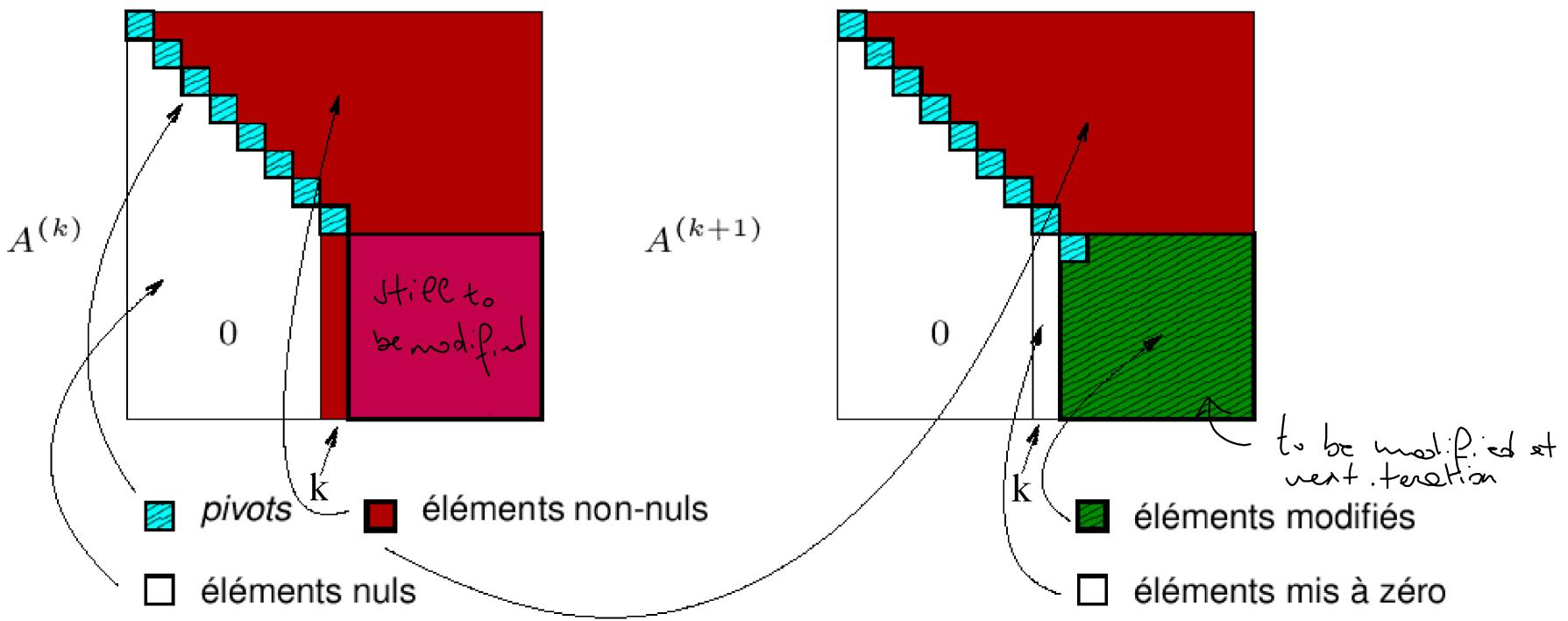
where, for  $k \geq 2$ , the matrix  $A^{(k)}$  is of the form

$$A^{(k)} = \begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & \dots & \dots & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & & & & a_{2n}^{(2)} \\ \vdots & & \ddots & & & \vdots \\ 0 & \dots & 0 & a_{kk}^{(k)} & \dots & a_{kn}^{(k)} \\ \vdots & & \vdots & \vdots & & \vdots \\ 0 & \dots & 0 & a_{nk}^{(k)} & \dots & a_{nn}^{(k)} \end{bmatrix},$$

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:

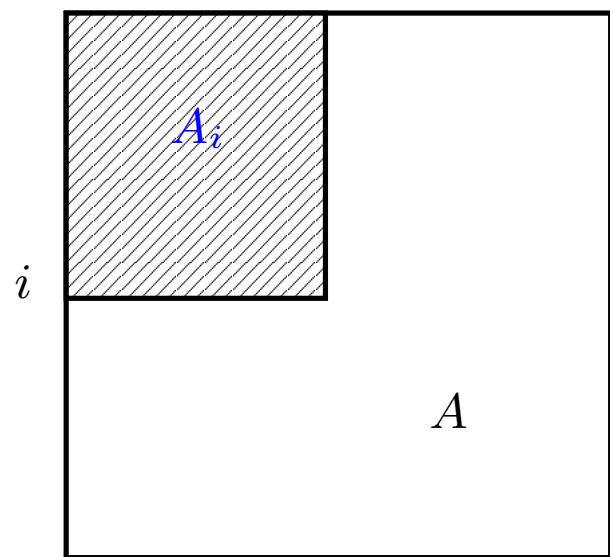
$n$	Number of flops of the computer		
	$10^9$ (Giga)	$10^{12}$ (Tera)	$10^{15}$ (Peta)
$10^2$	$7 \cdot 10^{-4}$ sec	negligible	negligible
$10^4$	11 min	0.7 sec	$7 \cdot 10^{-4}$ sec
$10^6$	21 years	7.7 months	11 min
$10^8$	o.r.	o.r.	21 years

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.

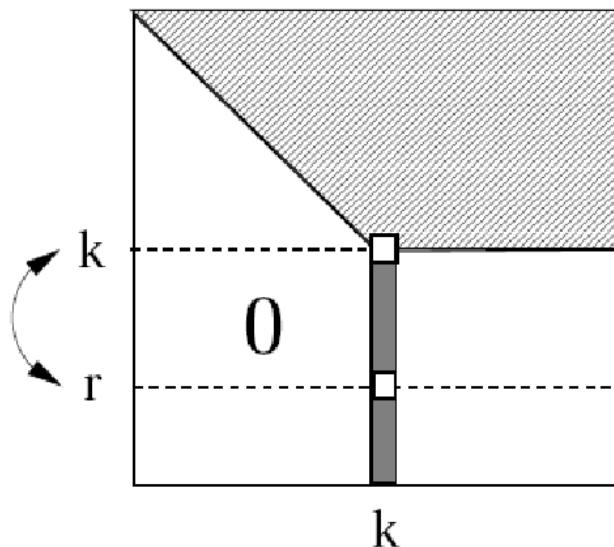
**Example 5.** Let us go back to the matrix (10) for which the Gauss method gives a null pivot at the second step. By just exchanging the second and the third rows, we get a non-zero pivot and can execute one step further. Indeed,

$$A^{(2)} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & \boxed{0} & -1 \\ 0 & -6 & -12 \end{bmatrix} \implies P_2 A^{(2)} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & \boxed{-6} & -12 \\ 0 & 0 & -1 \end{bmatrix},$$

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).

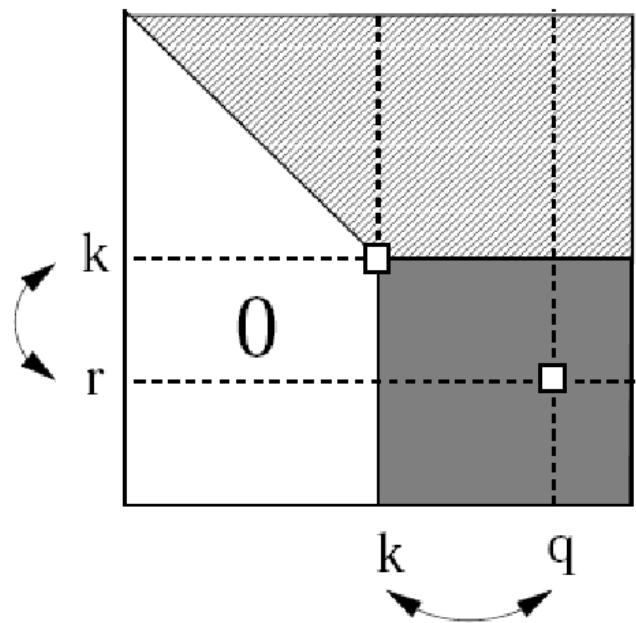


! Avoid big multipliers  
 $\Rightarrow$  avoid division by values close to zero  
 (you could have AMPLIFICATION OF ROUNDING ERRORS)

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
```

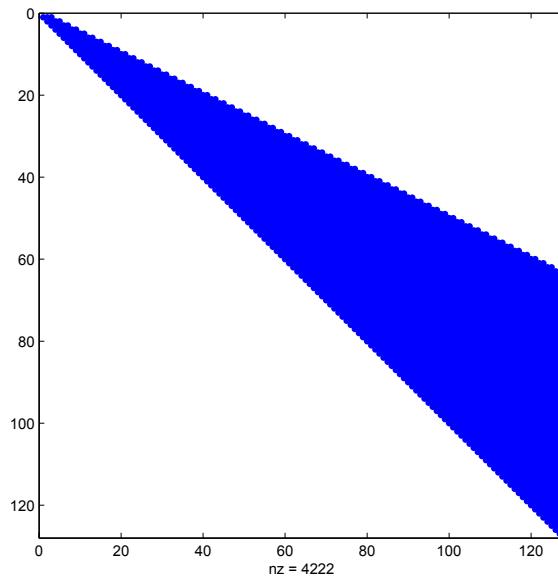
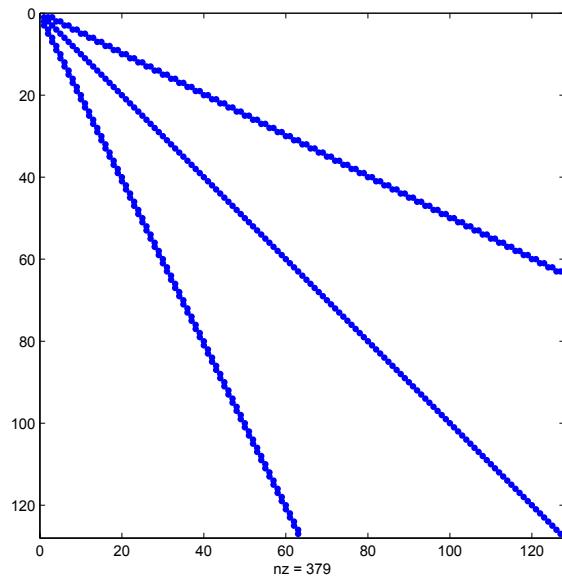
that grow proportionally to  $n^3$ . Indeed, if we calculate the ratio  $fl(n)/n^3$  we find a “quasi” constant value.

```
>> 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.

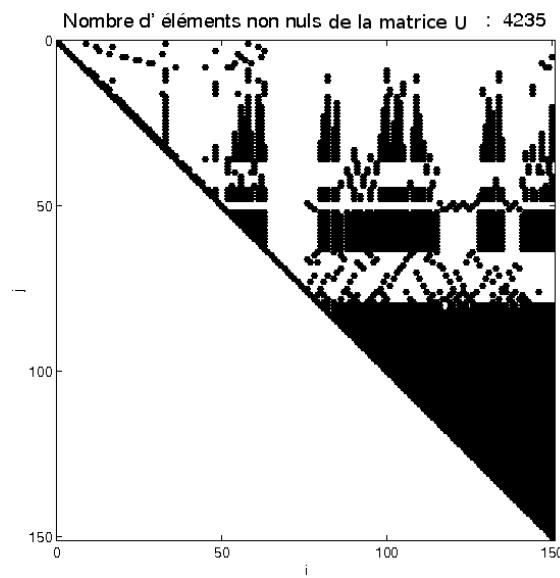
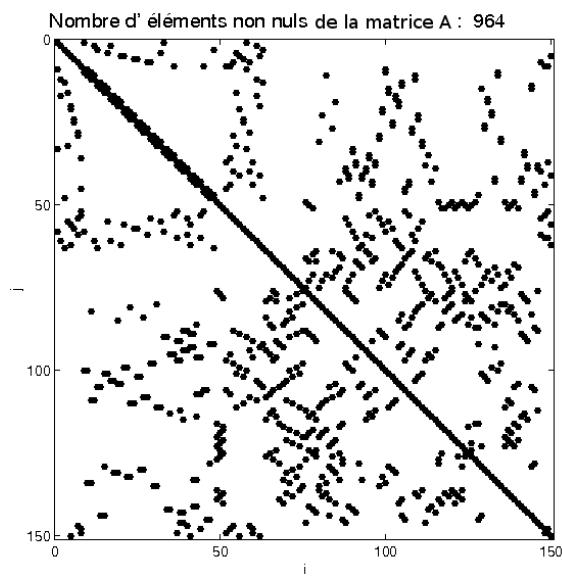
# Memory space limitations

**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  $R$ .



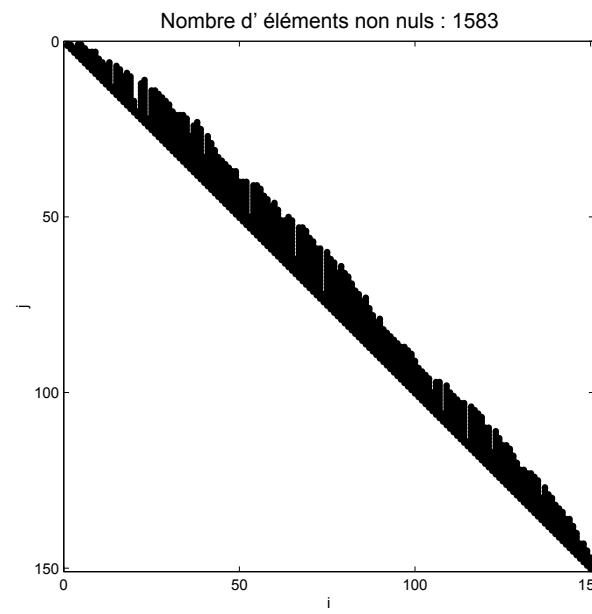
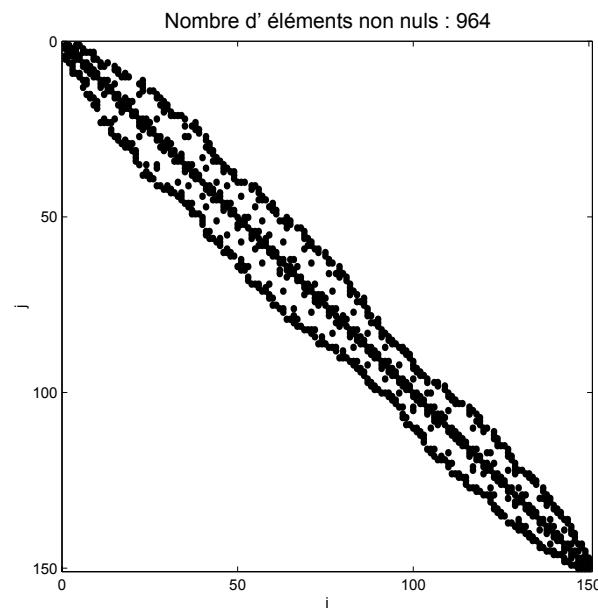
**Example 7.** Let us consider the problem of calculating the deformations in a structure subject to a given set of forces. The discretization using the finite elements method generates a matrix  $A$  of size  $150 \times 150$ . (The same matrix would have been produced by the approximation of an electric 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:

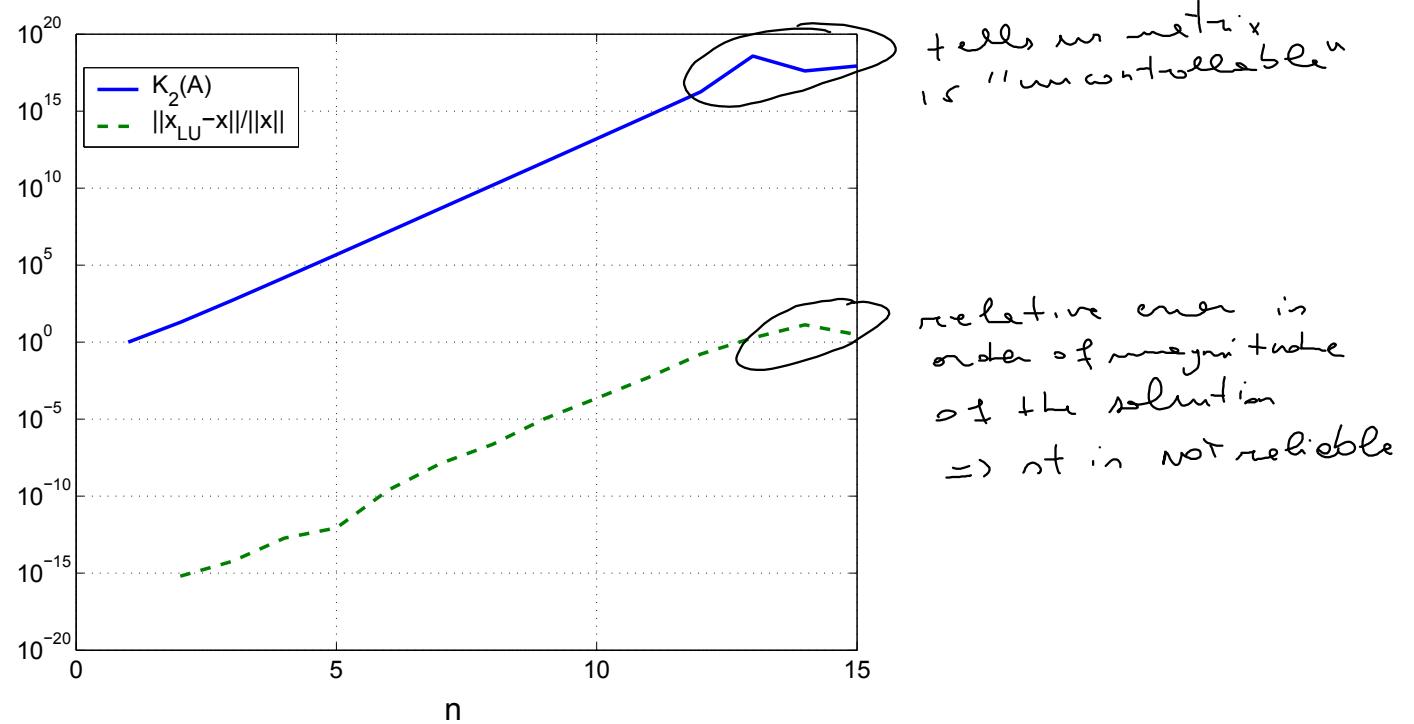
$$A_{ij} = \frac{1}{i+j-1}, \quad i, j = 1, \dots, n$$

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

$$\downarrow \quad \downarrow$$

$$x^{(k+1)} \quad x^{(k)}$$

*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  $\mathbf{b}$  and  $A\mathbf{x}^{(k)}$  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 be 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

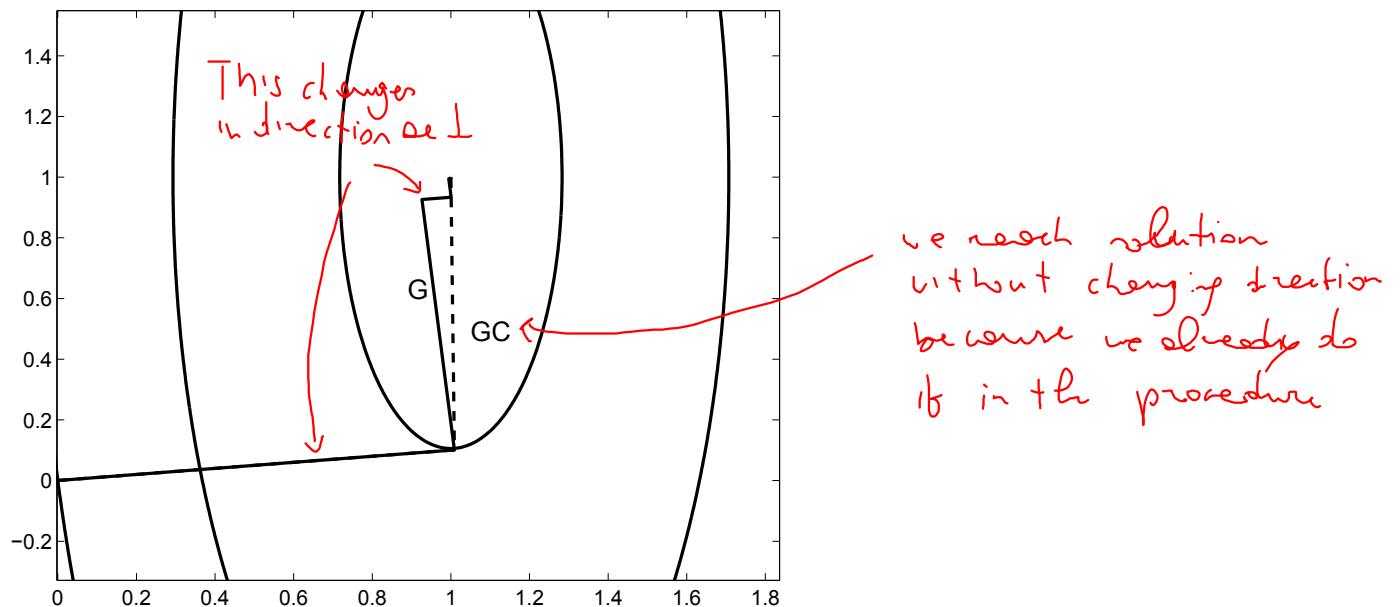
- Not so, but min number = to order of  $A$  matrix

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:

$$P\mathbf{z}^{(1)} = \mathbf{r}^{(1)} \iff \mathbf{z}^{(1)} = P^{-1}\mathbf{r}^{(1)} = \begin{pmatrix} 0.0896 \\ -0.0537 \end{pmatrix}$$

$$\beta_0 = \frac{(A\mathbf{p}^{(0)})^T \mathbf{z}^{(1)}}{(A\mathbf{p}^{(0)})^T A\mathbf{p}^{(0)}} = \frac{(A\mathbf{z}^{(0)})^T \mathbf{z}^{(1)}}{(A\mathbf{z}^{(0)})^T \mathbf{z}^{(0)}} = -0.0077$$

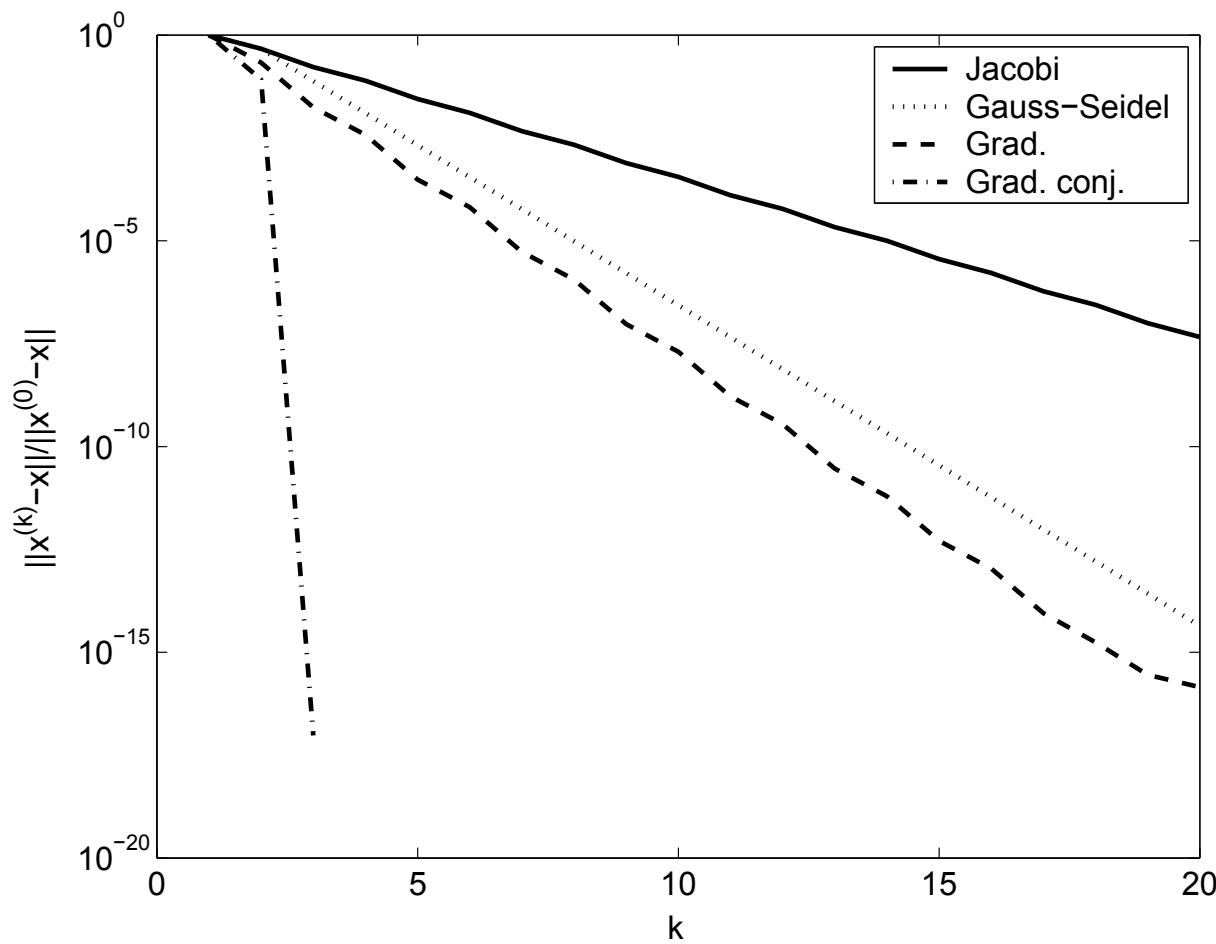
$$\mathbf{p}^{(1)} = \mathbf{z}^{(1)} - \beta_0 \mathbf{p}^{(0)} = \mathbf{z}^{(1)} - \beta_0 \mathbf{z}^{(0)} = \begin{pmatrix} 0.0838 \\ -0.0602 \end{pmatrix}.$$

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.

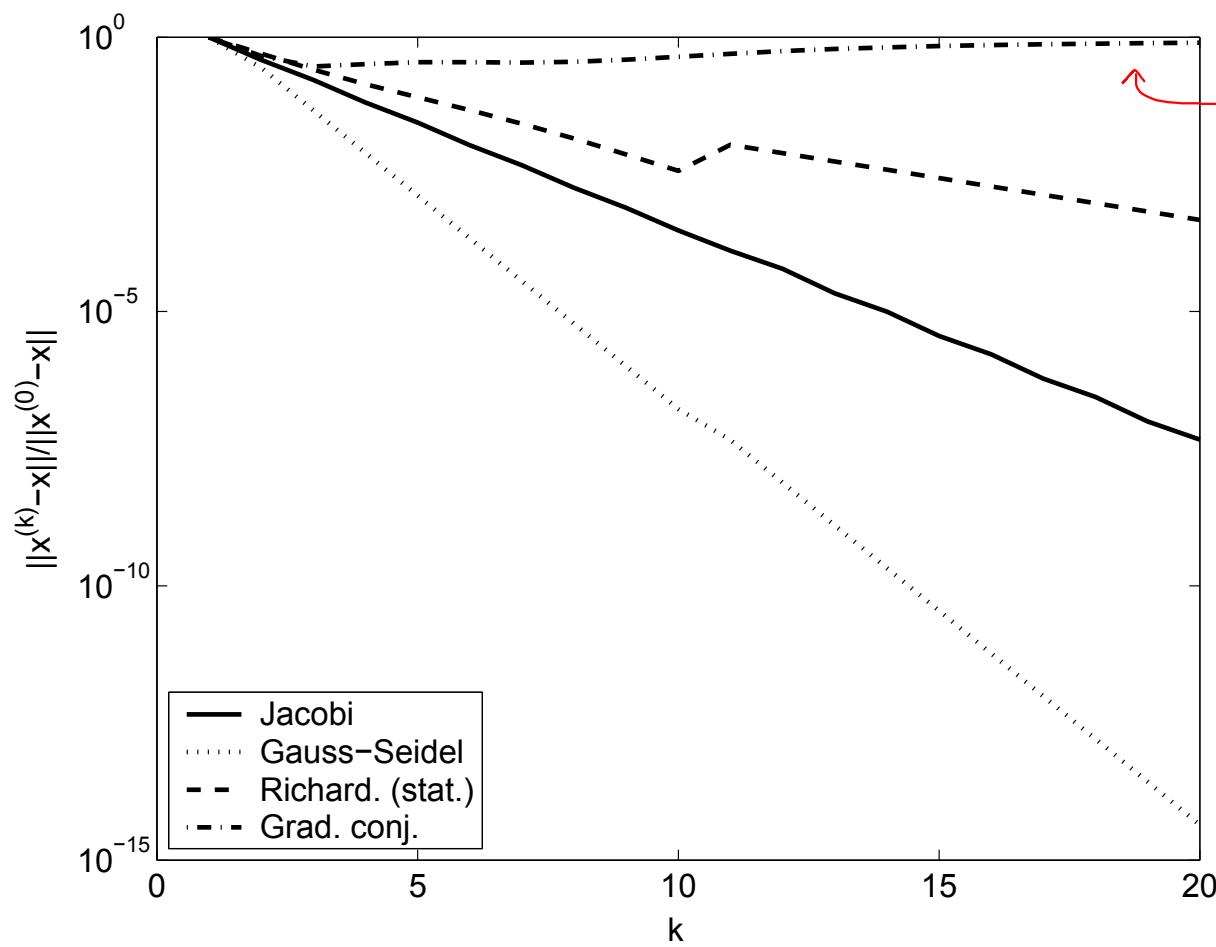
We approximate the solution with an iterative method starting from

$$\mathbf{x}^{(0)} = \begin{pmatrix} x_1^{(0)} \\ x_2^{(0)} \end{pmatrix} = \begin{pmatrix} 1 \\ \frac{1}{2} \end{pmatrix}.$$

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, like 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:

$$\frac{\|\mathbf{x}^{(k)} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq K(P^{-1}A) \frac{\|P^{-1}\mathbf{r}^{(k)}\|}{\|P^{-1}\mathbf{b}\|}.$$

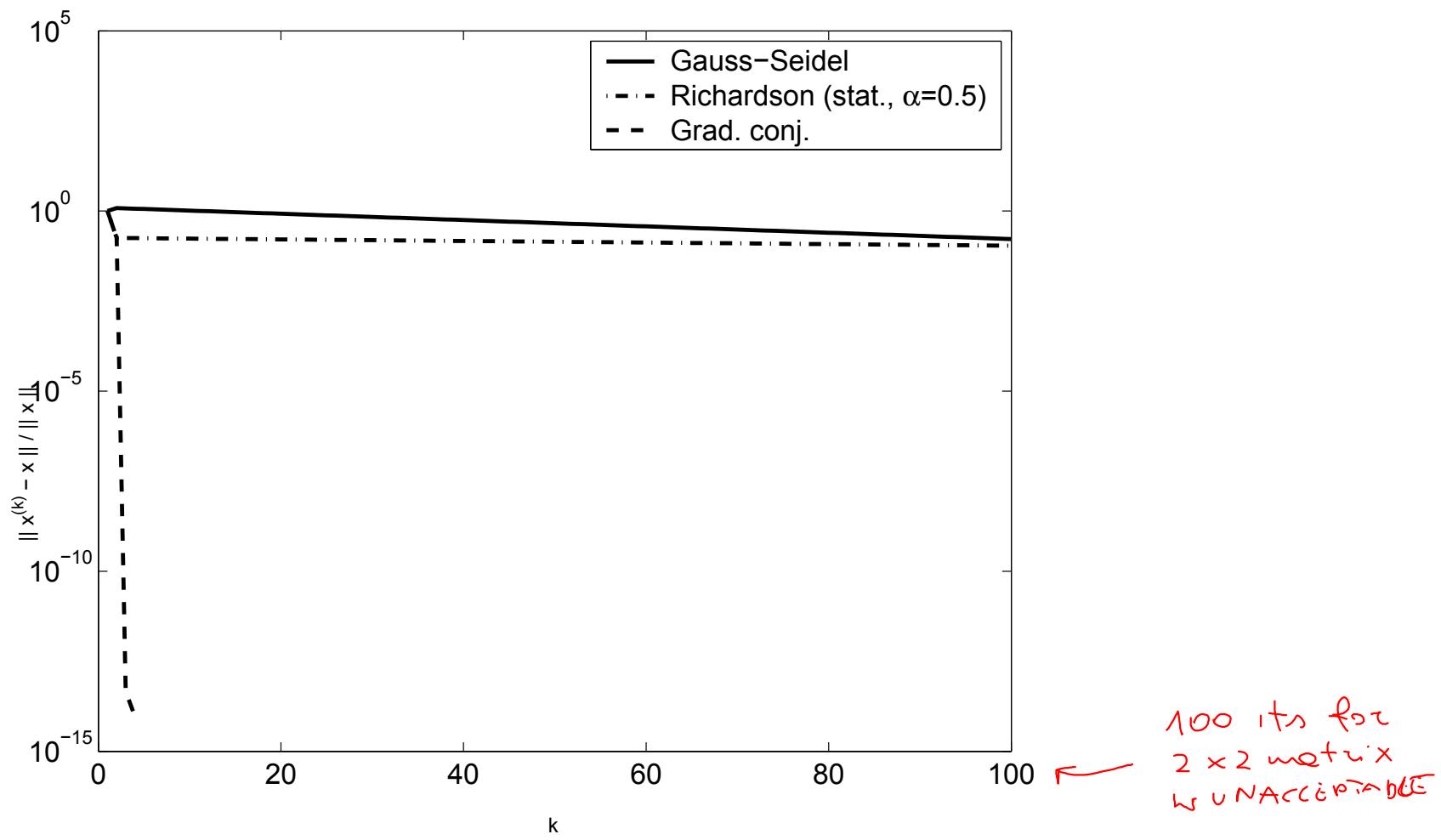
Some examples of convergence of iterative methods applied to linear systems of the kind  $A\mathbf{x} = \mathbf{b}$ .

**Example 4.** Lets start with the matrix

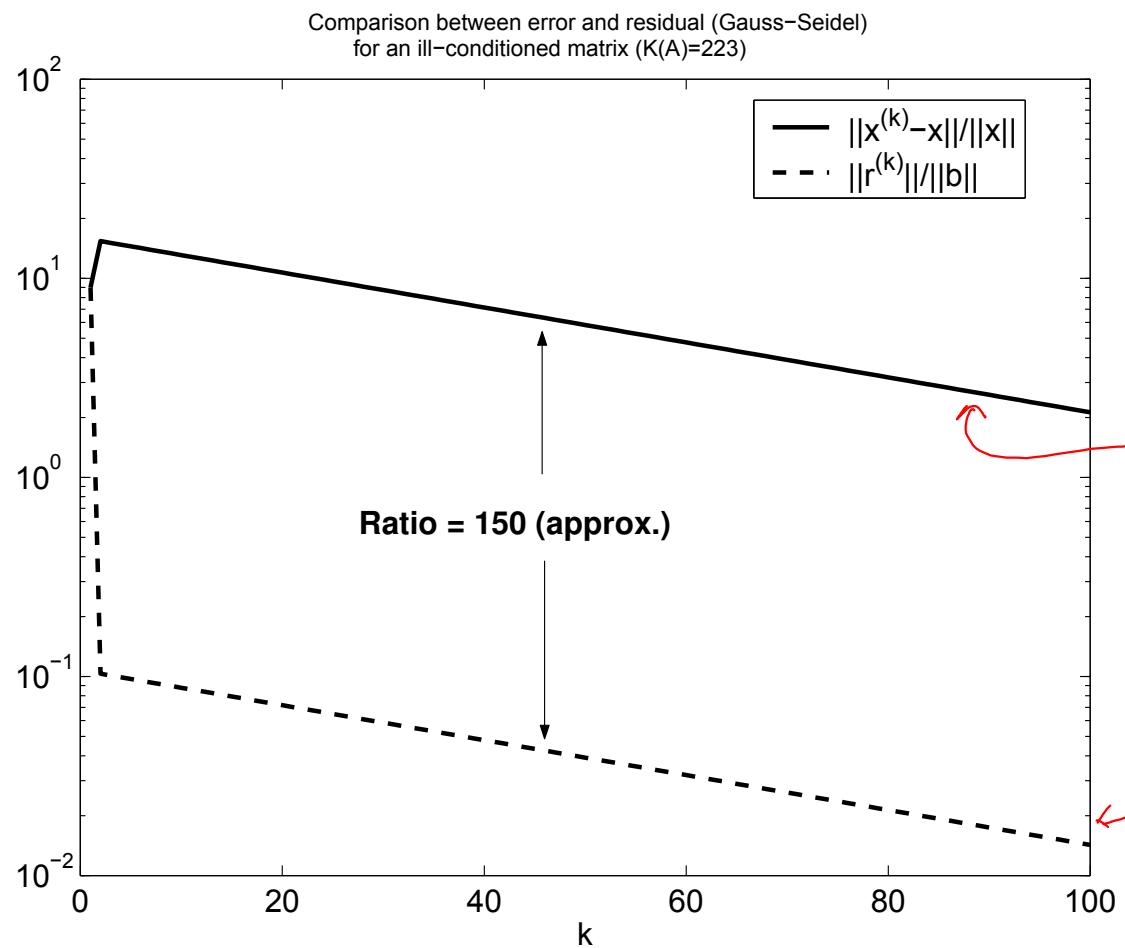
$$A = \begin{pmatrix} 5 & 7 \\ 7 & 10 \end{pmatrix} \quad (19)$$

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



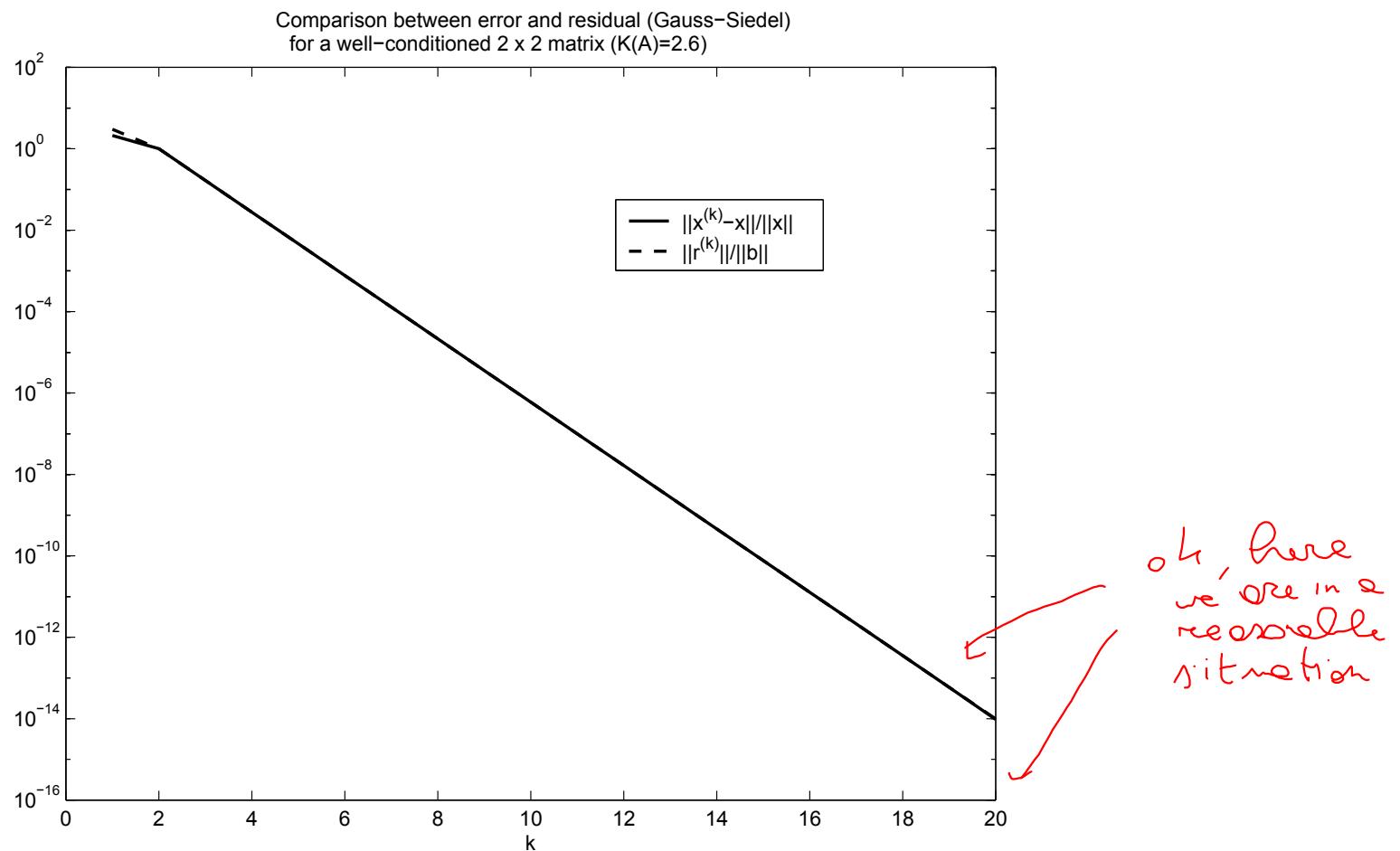
Comparison between error and residual (Gauss-Seidel) (19). Recall that  $K(A) \approx 223$



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

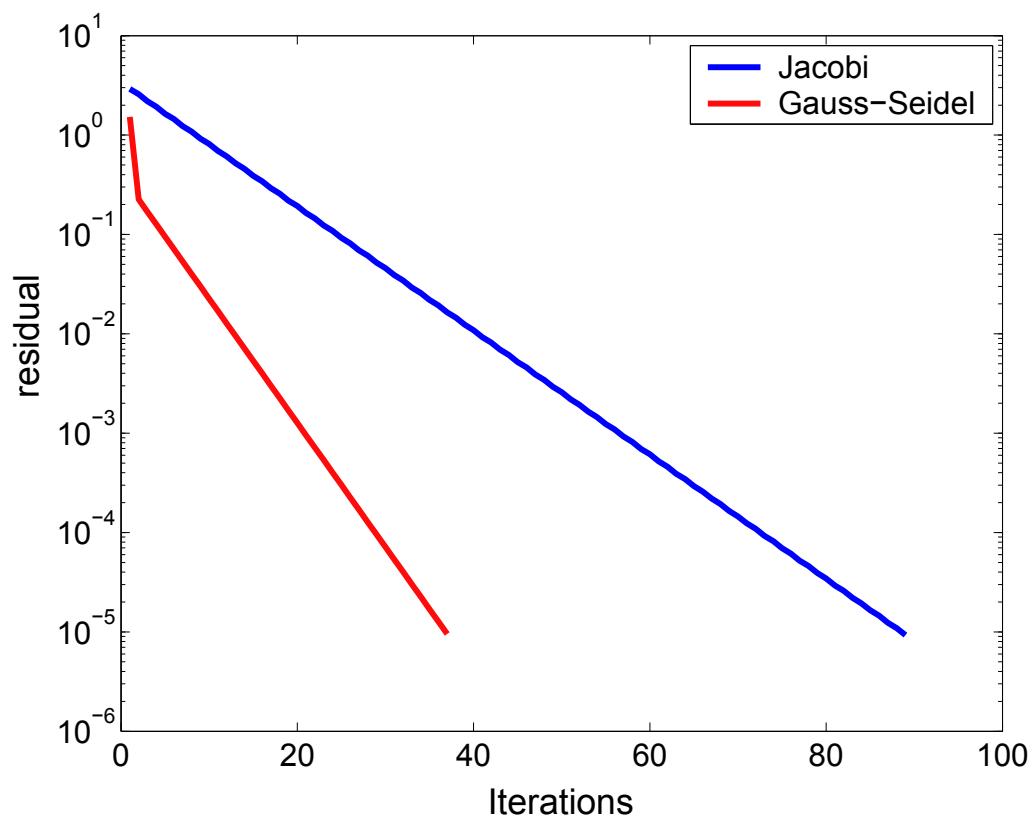
Computational cost (*flops*) and used memory (*bytes*): we consider the Cholesky and the conjugate gradient methods for sparse matrices of size  $n$  (originated in the approximation of solutions to the Poisson equation in the finite elements method) with  $m$  non zero elements.

		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

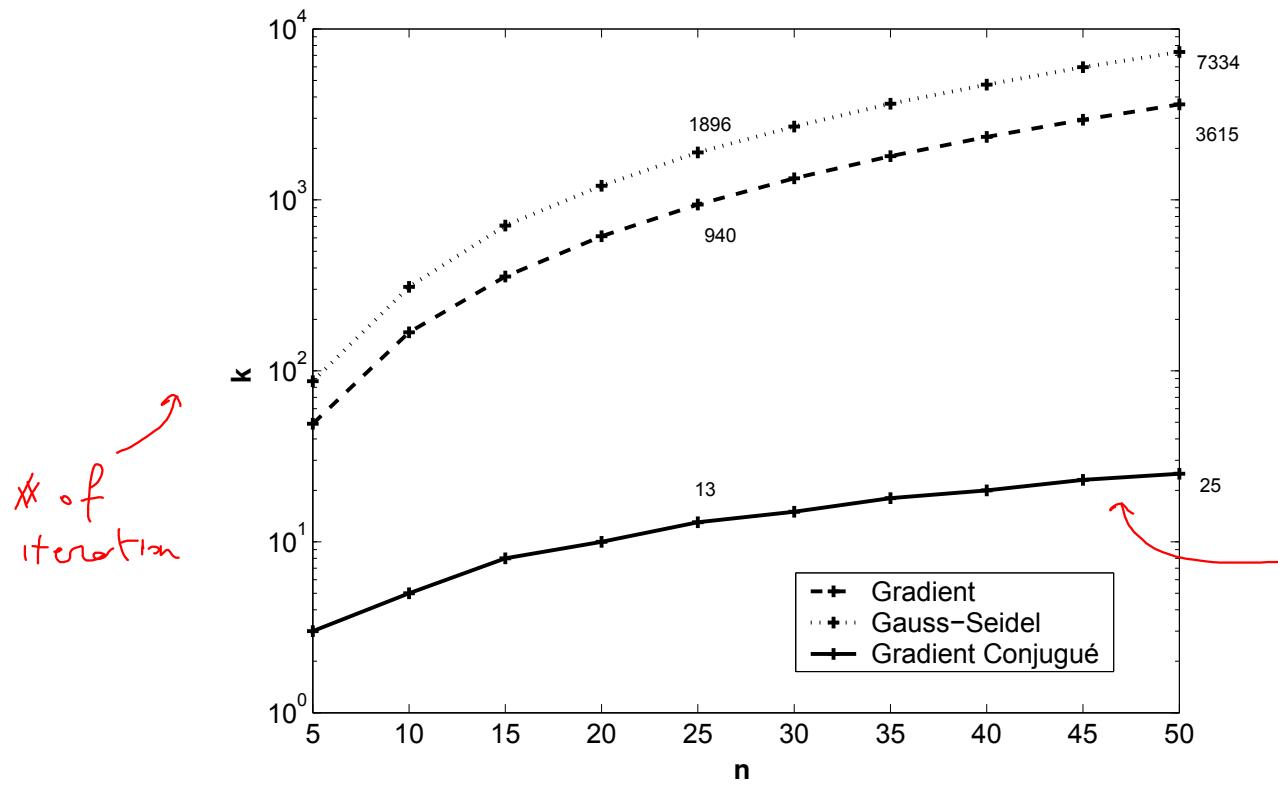
# Iterative methods: Jacobi, Gauss-Seidel

Behaviour of the error for a well conditioned matrix ( $K = 20$ )



# Iterative methods: G-S, Gradient, Gradient Conjugate

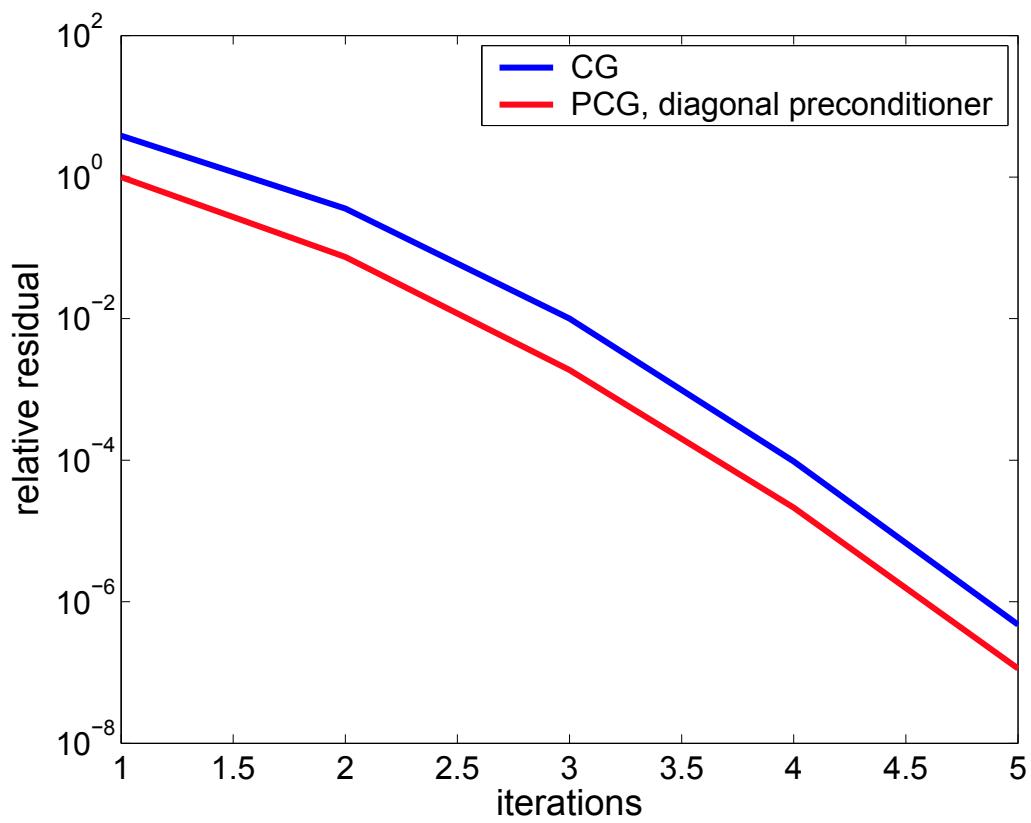
Number of iterations as function of the size  $n$  of a stiffness matrix, tolerance  $10^{-6}$



conjugate gradient  
requires less  
time than the  $\times$  of  
A orthogonal direction  
is equal to the order  
of the matrix

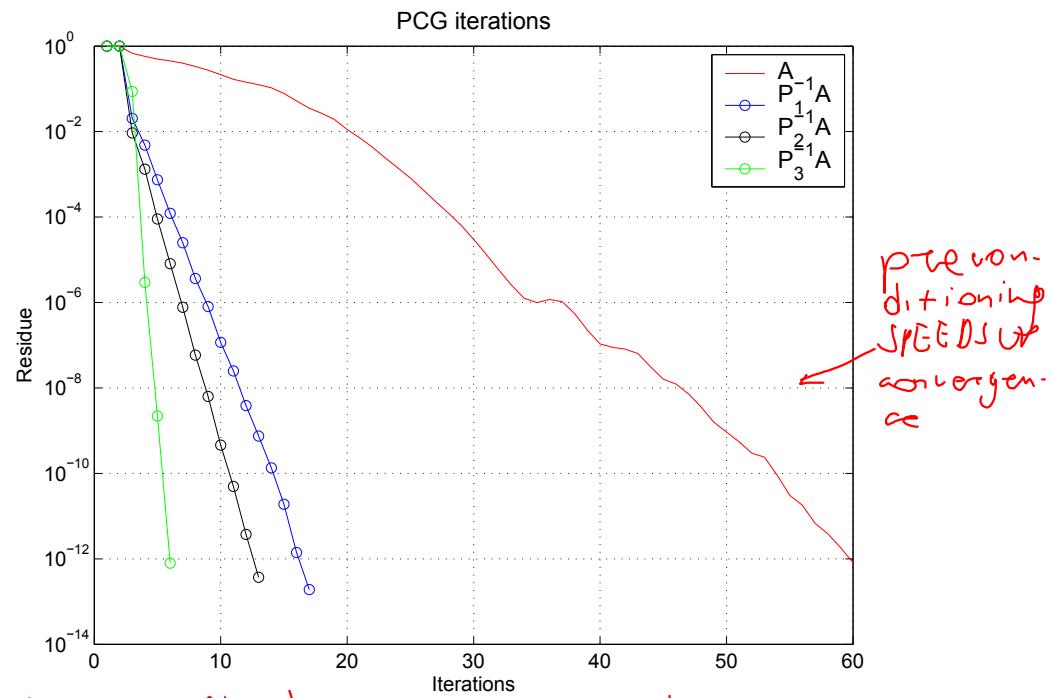
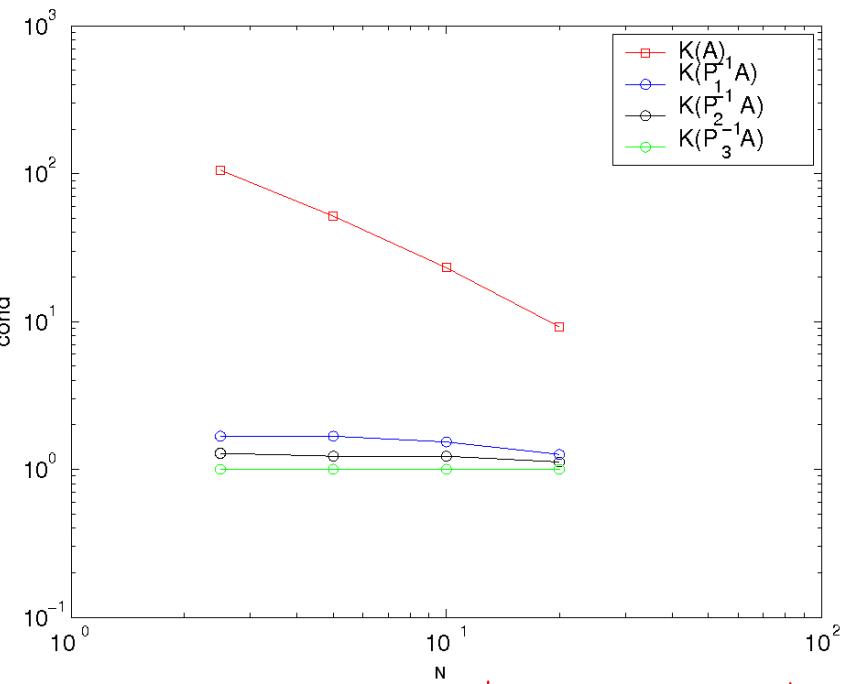
# Preconditioning

The convergence conjugate gradient and the preconditioned conjugate gradient methods with a diagonal preconditioner ( $K = 4 \cdot 10^8$ )



# Preconditioning

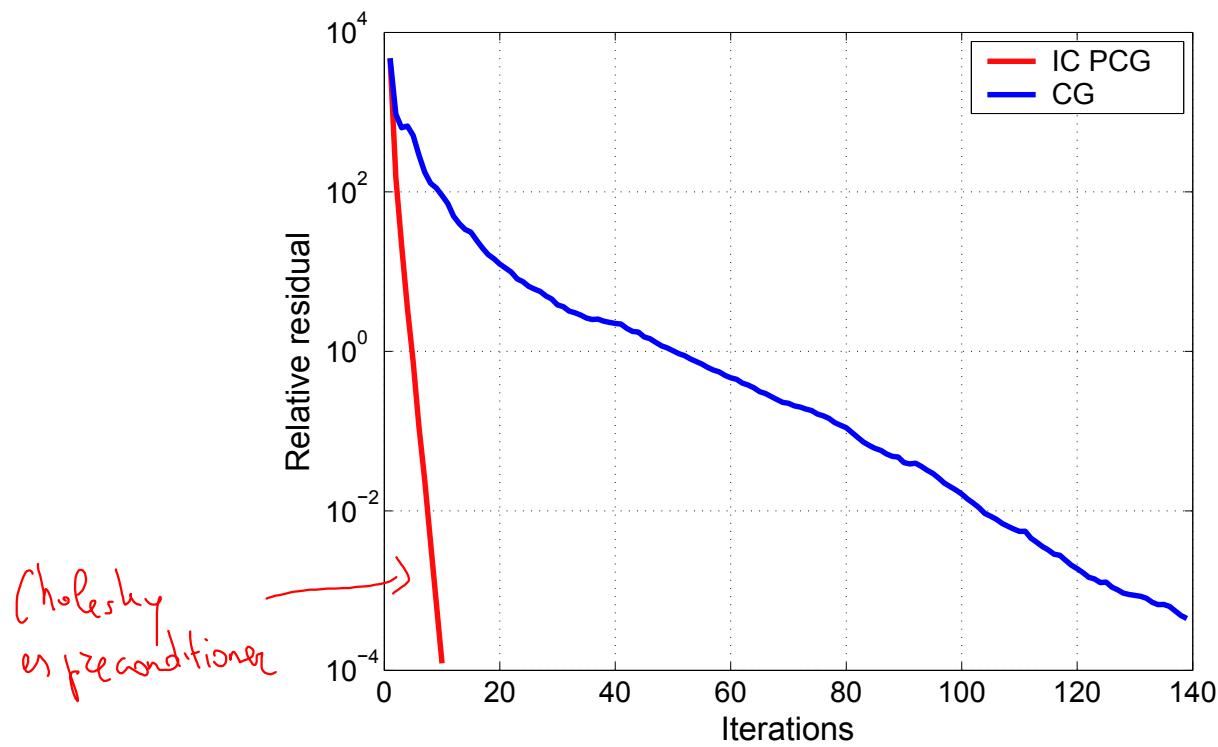
Role of the preconditioning over the condition number of a matrix  
 (approximation by the finite element method of the Laplacian operator)



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

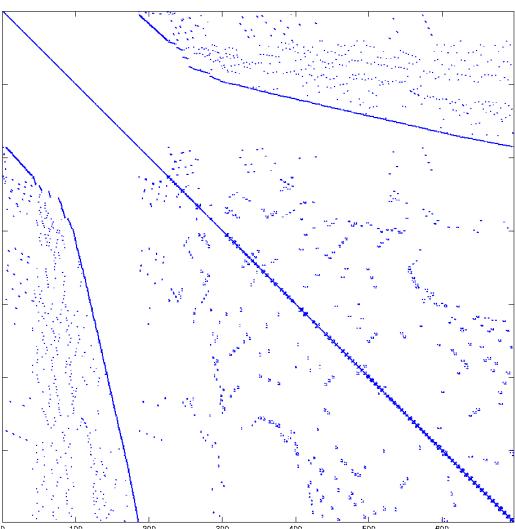
# Preconditioning

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 ( $K = 1.5 \cdot 10^3$ )

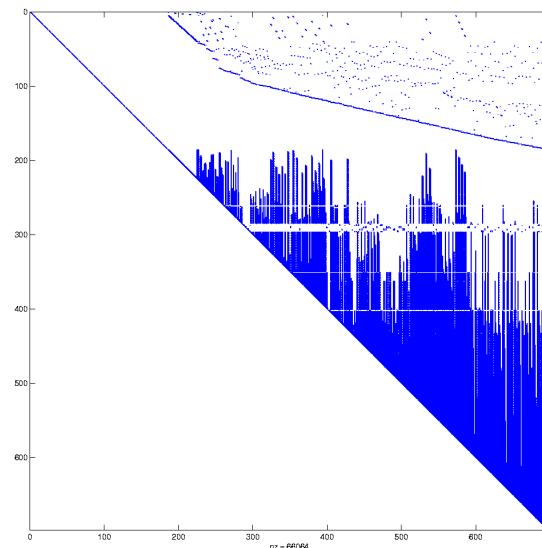


# Pre-conditioning

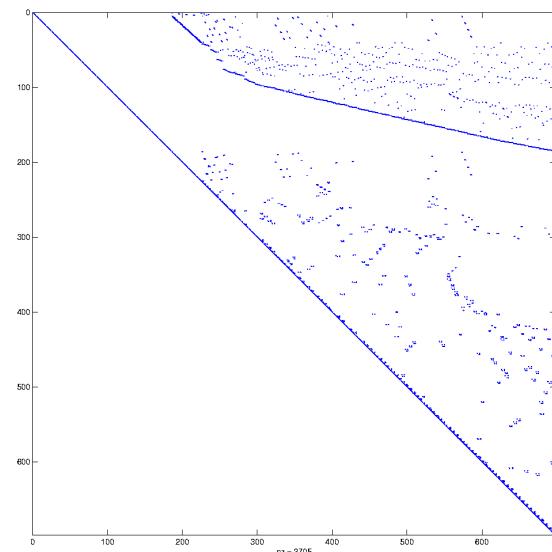
Comparison between the non zero elements of the sparse matrix  $A$  from the previous example, its Cholesky factor  $R$  and the matrix  $\tilde{R}$  obtained by the Cholesky incomplete decomposition:



$A$



$R$



$\tilde{R}$

# Non-linear systems

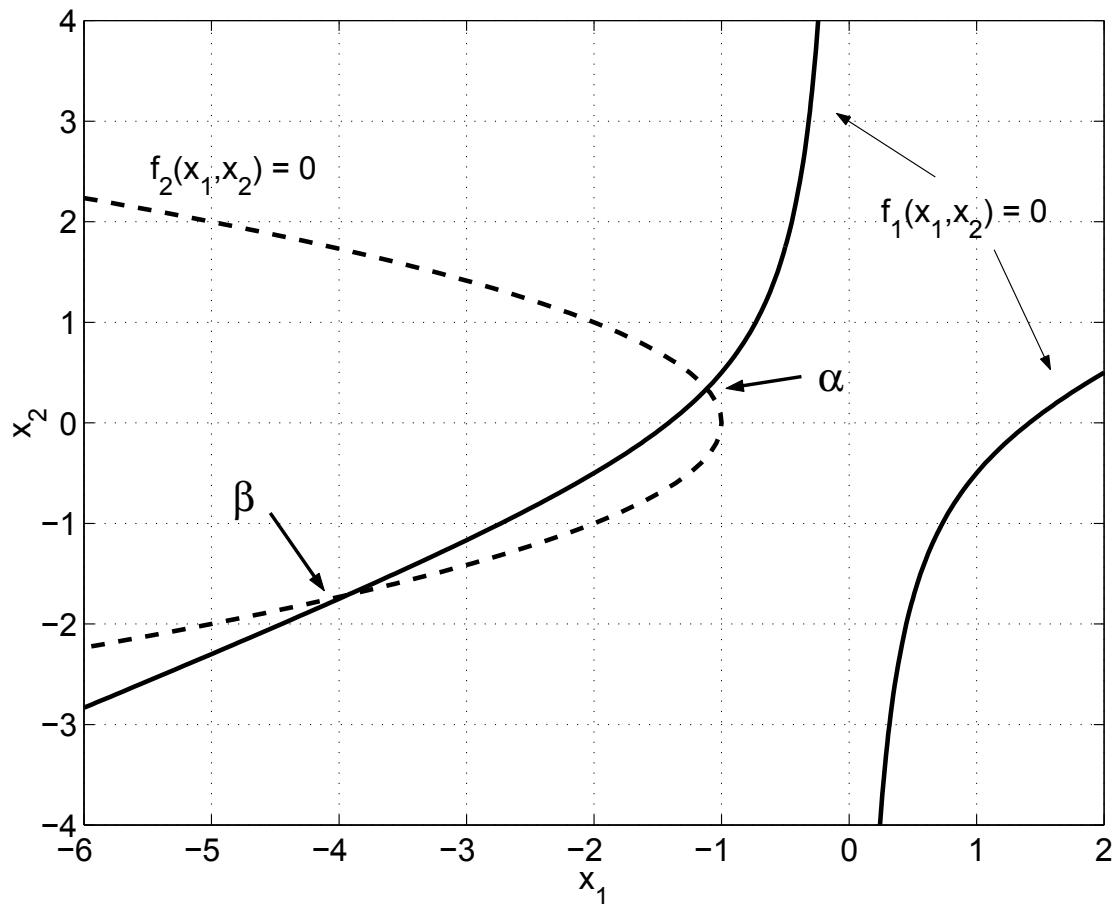
**Example 6.** Lets consider the following system of non-linear equations:

$$\begin{cases} x_1^2 - 2x_1x_2 = 2 \\ x_1 + x_2^2 = -1. \end{cases} \quad (20)$$

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 +ve / -ve sign at  $-1$  ( 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 iteration 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  $h \rightarrow 0$ )

$$u'(x) := \lim_{h \rightarrow 0} \frac{u(x) - u(x-h)}{h}$$

$$u'(x) := \lim_{h \rightarrow 0} \frac{u(x+h) - u(x-h)}{2h}$$

backward  
BD  
/

central  
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 turn any 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:

$$\frac{1}{h^2} \left[ -u^{i-1} + 2u^i - u^{i+1} \right] \approx -u'' = f_i \quad \begin{matrix} \text{for all } i \\ \text{system of equations} \end{matrix}$$

$$\sum_{j=0}^{N-1} A_{ij} u^j = f_i$$

$$A_{ij} = \frac{1}{h^2} \begin{cases} -1 & \text{when } j = i-1 \text{ or } i+1 \\ 2 & \text{when } j = i \\ 0 & \text{otherwise} \end{cases} \quad \text{and } i \in \{1, N-2\}$$

$$\begin{pmatrix} 0 & -1 & 2 & -1 & & & & \\ -1 & 2 & -1 & & & & & \\ & \ddots & \ddots & \ddots & & & & \\ & & -1 & 2 & -1 & & & \\ & & & & \ddots & \ddots & & \\ & & & & & 0 & & \\ & & & & & & \ddots & \\ N-2 & & & & & & & 0 \end{pmatrix} \cdot \frac{1}{h^2} \quad f_i = f(x_i) \quad i = 1, \dots, N-2$$

What about 1st and last columns?

For  $i=0, i=N-1$  we impose

We set the BOUNDARY CONDITION

$$u^0 = 0$$

$$u^{N-1} = 0$$



$$A_{00} = 1$$

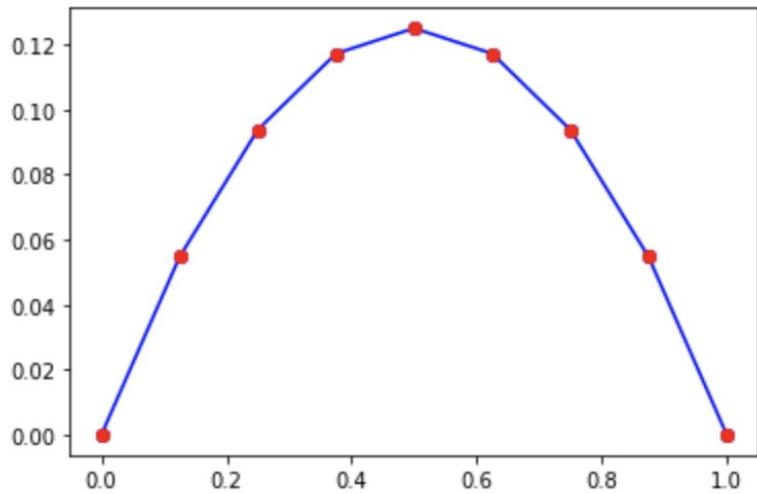
$$A_{N-1, N-1} = 1$$

$$f_0 = 0$$

$$f_{N-1} = 0$$

! The more the points the higher the condition number

$\Rightarrow$  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

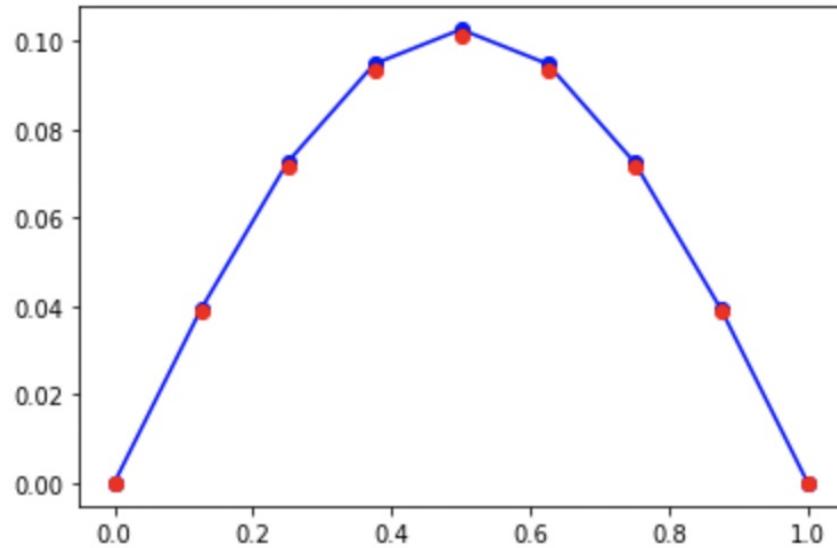
We have errors (blue overlaps red)  
because error is up to  $u'''$ , but  $u'''$  of  $e$   
parallel is zero!

Same thing with

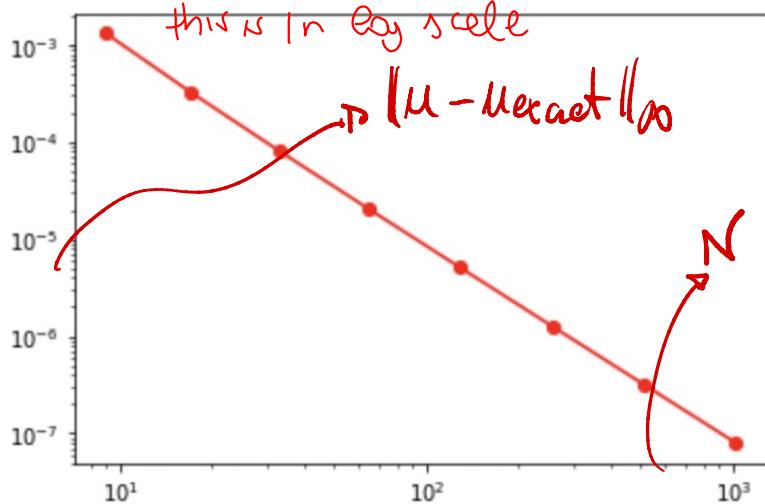
$$f(x) = \sin(\pi x)$$

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

$$\textcircled{1} \quad \begin{cases} -u'' = f & \text{in } [0,1] \\ u(0) = u(1) = 0 \end{cases}$$

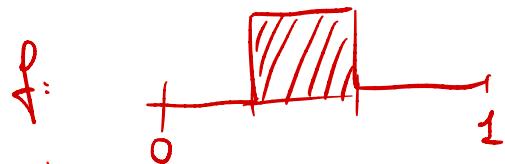
We tried to:  
approximate  $u''$

- 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 I 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

$$\int_0^1 Wv' = \int_0^1 W'v + \int_0^1 Wv'$$

$$\Rightarrow \int_0^1 W'v = - \int_0^1 Wv' + Wv \Big|_0^1$$

$$\boxed{\int_0^1 -u''v = \int_0^1 u'v' + v'v \Big|_0^1}$$

for the because  
 $v(0) = v(1) = 0$

$$-u'' = f \Rightarrow \int_0^1 u'v' = \int_0^1 fv \quad \forall v \in V$$

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

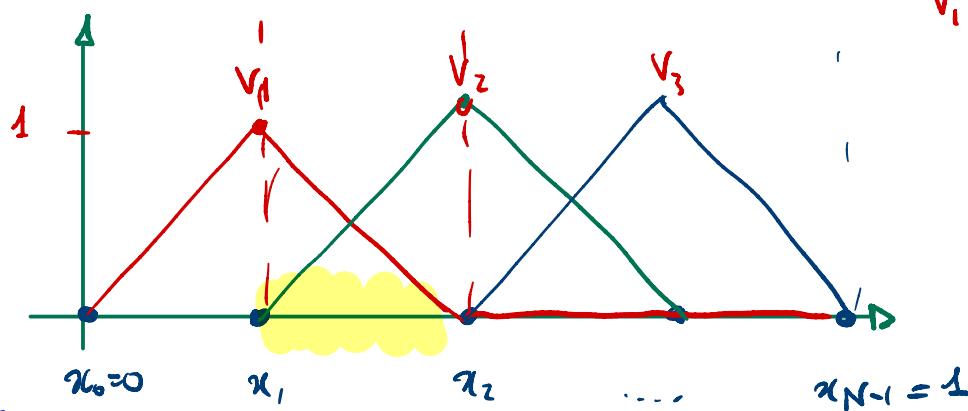
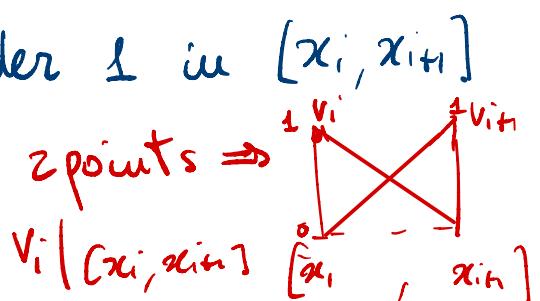
$$\{x_i\}_{i=0}^{N-1} = ih = \frac{i}{(N-1)} \quad \text{the intervals are } [x_i, x_{i+1}] \quad i=0, \dots, N-2$$

$$V_h^k := \left\{ v \in \mathbb{R}^k([x_i, x_{i+1}]) \mid v(0) = v(1) = 0 \right. \\ \left. v \in C^0([0, 1]) \right\}$$

Continuous piecewise polynomial of order  $k$ .

Example  $k=1 \Rightarrow P^1([x_i, x_{i+1}])$

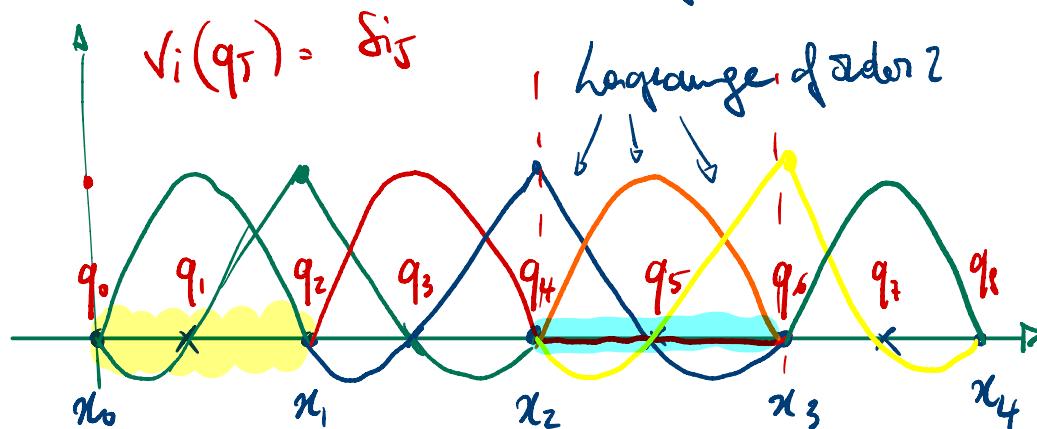
We choose Lagrange basis of order 1 in  $[x_i, x_{i+1}]$   
with support points  $\{x_i, x_{i+1}\}$



$q_i = x_i$   
for linear  
polynomials.

For example, choose

$$v_i(x_j) = \delta_{ij} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}, \quad v_i \in V_h^1$$

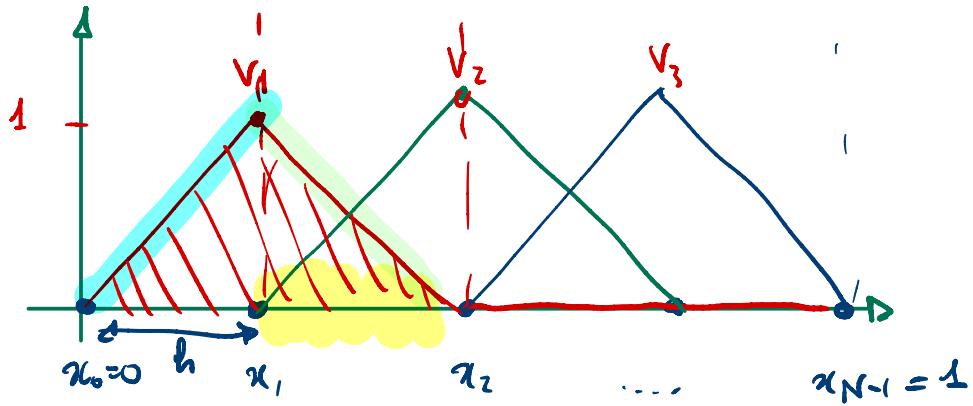


For polynomials of

order 2, choose

$$\{q_i\}_{i=0}^{(N-1)(N-2)} = \{x_i\}_{i=0}^{N-1} \cup \left\{ \frac{x_i + x_{i+1}}{2} \right\}_{i=0}^{N-2}$$

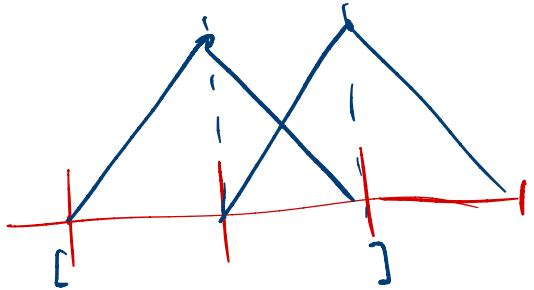
1D case linear



$$\dim(V_h) = N-2$$

$$v_i = \begin{cases} \frac{(x - x_{i-1})}{h} & x \in [x_{i-1}, x_i] \\ \frac{(x_{i+1} - x)}{h} & x \in [x_i, x_{i+1}] \\ 0 & \text{elsewhere} \end{cases} \quad i = 1, \dots, N-2$$

$$v_i^l = \begin{cases} \frac{1}{h} & \text{in } [x_{i-1}, x_i] \\ -\frac{1}{h} & \text{in } [x_i, x_{i+1}] \\ 0 & \text{elsewhere} \end{cases}$$



$$A_{ij} := \int_0^1 v_i^l v_j^l \neq 0 \quad \begin{matrix} \text{if } i = j, j-1, j+1 \\ \text{elsewhere} \end{matrix}$$

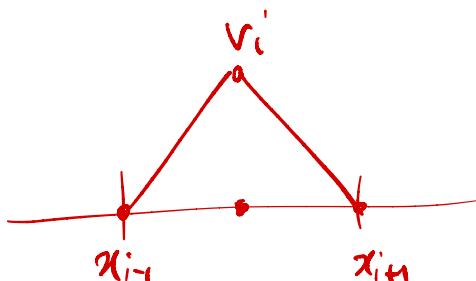
$$A_{jj} = \int_{x_{i-1}}^{x_{i+1}} \left(\frac{1}{h}\right)^2 = 2 \frac{1}{h} \quad \text{if } i=j$$

$$A_{ij} = \int_{x_i}^{x_{i+1}} \left(-\frac{1}{h}\right)^2 = -\frac{1}{h} \quad \begin{matrix} \text{if } i = j-1 \\ \text{or } i = j+1 \end{matrix}$$

$$A_{ij} := \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & & -1 & 2 & -1 \end{pmatrix} \frac{1}{h}$$

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  $1/h$

$$h \left[ \frac{1}{2} v_i(x_{i+1}) f(x_{i+1}) + \frac{1}{2} v_i(x_i) f(x_i) \right] + h \left[ \frac{1}{2} v_i(x_i) f(x_i) + \frac{1}{2} (v_i(x_{i+1}) f(x_{i+1})) \right]$$

$$= f(x_i) h$$

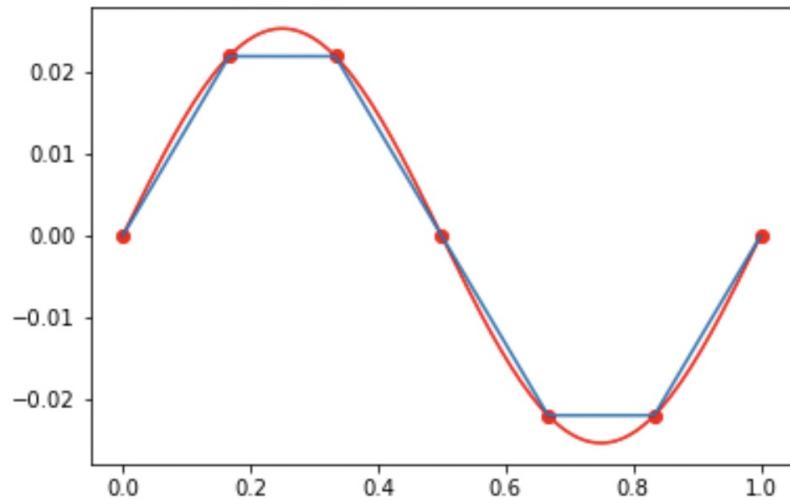
Approximation property of FEM:

$$\int_0^1 u' v_h = \int_0^1 f v_h \quad \forall v_h \in V_h$$

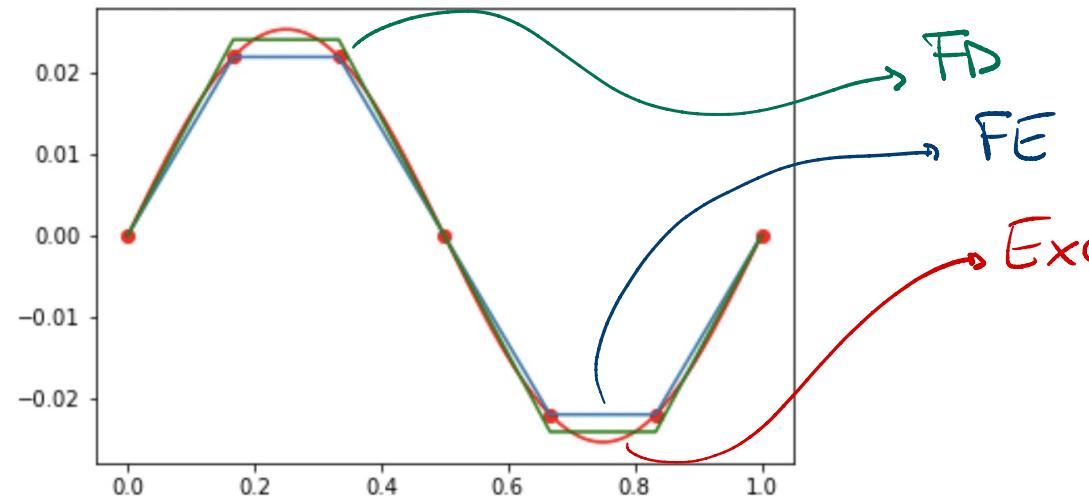
$$\int_0^1 u_h' v_h = \int_0^1 f v_h \quad \forall v_h \in V_h \subset V$$

$$\int_0^1 (u' - u_h') v_h = 0 \quad \forall v_h \in V_h \subset V$$

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 "petticoat" 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!

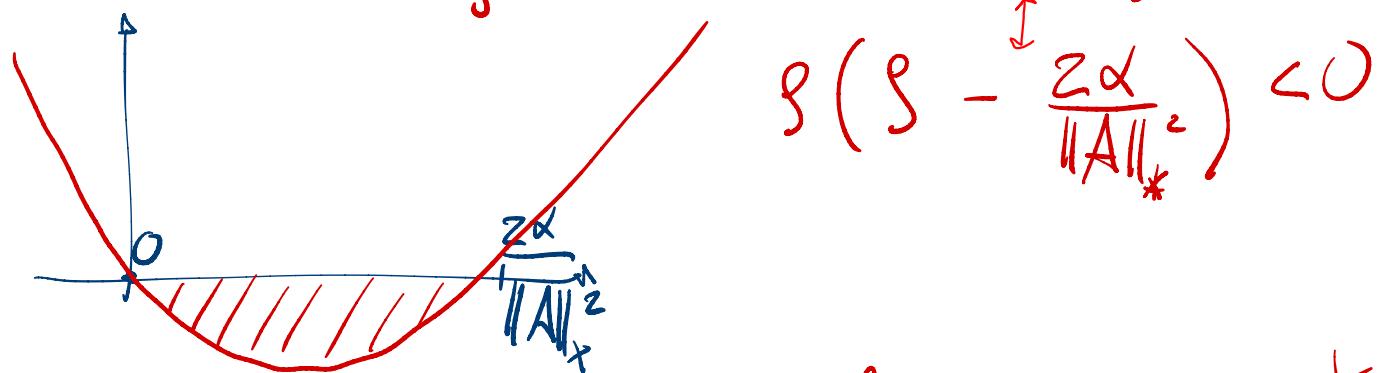
$$\|T_\rho(v) - T_\rho(w)\|_V^2 = \|v-w - \rho z A(v-w)\|_V^2$$

$$\begin{aligned}
 \|T_p(v) - T_p(w)\|_V^2 &= \|v-w\|_V^2 + g^2 \left\| z A(v-w) \right\|_V^2 - 2g \langle z A(v-w), v-w \rangle_V \\
 &= \|v-w\|_V^2 + g^2 \|A(v-w)\|_{V^*}^2 - 2g \langle A(v-w), v-w \rangle \\
 &\leq \|v-w\|_V^2 + g^2 \|A\|_*^2 \|v-w\|_V^2 - 2g \alpha \|v-w\|_V^2
 \end{aligned}$$

$$\|T_g(v) - T_g(w)\|^2 \leq \left( g^2 \|A\|_*^2 - 2g\alpha + 1 \right) \|v-w\|_V^2$$

$\underbrace{\phantom{\dots}}_{< 1}$

Can we choose  $g$  s.t.  $\underbrace{g^2 \|A\|_*^2 - 2g\alpha + 1}_{< 1} < 1$ ?



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:

$$\alpha \|u - u_h\|^2 \leq \langle A(u - u_h), u - u_h \rangle = \langle A(u - u_h), u - v_h \rangle \leq \|A\|^* \|u - u_h\| \|u - v_h\| \quad \forall v_h \in V_h$$

$$\|u - u_h\|_V \leq \frac{\|A\|_*}{\alpha} \inf_{v_h \in V_h} \|u - v_h\|_V$$

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

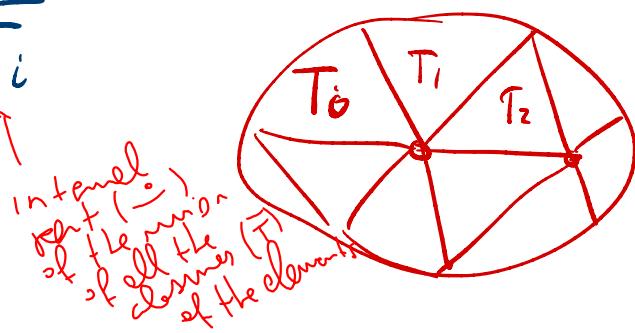
$$\|u - I(u)\|_{m,T} \leq C h_T^{k+1-m} |u|_{k+1,T}$$

measure  
of the  
element

on each "element"  $T$  of  $\mathcal{S}_h$ , discretization of  $S$ . 1D:  $S = [0, 1]$   $T = [x_i, x_{i+1}]$

$$S_h = \bigcup_{i=1}^n T_i$$

$$h_T = (x_{i+1} - x_i)$$



$$\|u\|_{\alpha,T} := \left( \sum_{|\alpha|=d} \int_T ((D^\alpha u)^2) \right)^{\frac{1}{2}}$$

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

$$u \approx V V^T f$$

$$Vv \approx u \quad u = 0$$

$$\text{Start with } q_0 = r_0 = A u_0 - f = -f$$

$$q_0^T A q_0 := \beta_0 \rightarrow v_0 = \frac{f_0}{\beta_0}$$

$$\rightarrow r_0^T A v_0 = 1$$

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

$$1) \begin{aligned} -\Delta u &= f && \text{in } \Omega \\ u &= 0 && \text{on } \partial\Omega \end{aligned}$$

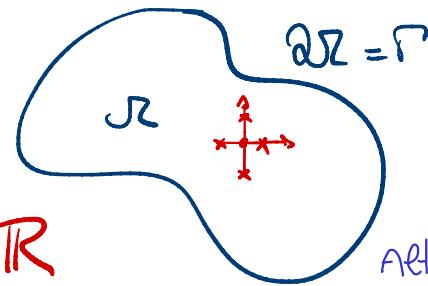
LAPLACIAN

$$\downarrow \quad -\Delta u := -\frac{\partial^2 u}{\partial x_0^2} - \frac{\partial^2 u}{\partial x_1^2} = \sum_{d=0}^{N-1} -\frac{\partial^2 u}{\partial x_d^2}$$

$$u : \Omega \rightarrow \mathbb{R}$$

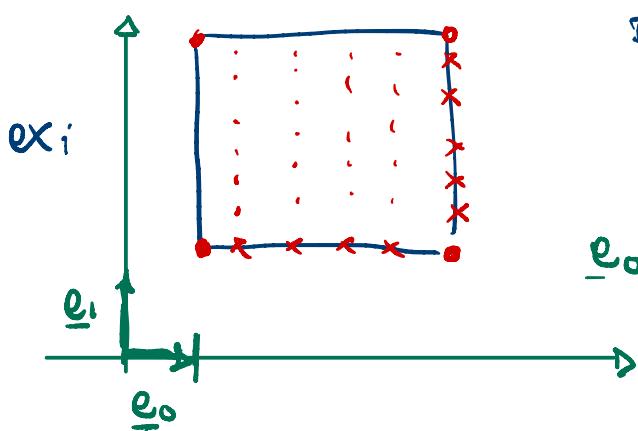
sometimes:

$$\begin{aligned} \text{Alternatively,} \\ \text{another problem} \\ \text{that can appear is} \\ -\Delta u + u = f \\ 2) \frac{\partial u}{\partial n} = 0 \quad \text{on } \partial\Omega \end{aligned}$$



$N=2$

For FD  $\Omega$  has to be "simple" and "regular"



Fix a grid of points  $X_i \in \Omega \subset \mathbb{R}^2$

$$\Omega = [a, b] \times [c, d]$$

$$\underline{e}_0 := (1, 0) \quad \underline{e}_1 := (0, 1)$$

$$\underline{x} = \underline{e}_0 x^0 + \underline{e}_1 x^1 = \underline{e}_i x^i$$

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)  $d \in \{0, \dots, m-1\}$

$(X_{di})_d \rightarrow d$  comp. of  $X_{di}$   $i \in \{0, \dots, n-1\}$

$$\underline{X}_{di} := (\varphi_i, \psi_d)$$

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

$$X_{did} = \begin{cases} \varphi_i & \text{if } d=0 \\ \psi_d & \text{if } d=1 \end{cases}$$

$$\left\{ X_{did} \right\}_{\substack{d=0, \dots, m-1 \\ i=0, \dots, n-1}}$$

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  $d \rightarrow e_1$

$$v_i := u_{di} \quad -\Delta u_{di} \approx \frac{1}{h^2} [4u_{di} - u_{d-1,i} - u_{d+1,i} - u_{d-1,i} - u_{d,i+1}]$$

$$w_d := u_{di}$$

$$-\tilde{\Delta} u_{di} = f_{di} = f(X_{di})$$

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

*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

$$\sum_j \tilde{A}_{Ij} \tilde{\mu}_j = \tilde{f}_I$$

$$I, j \in [0, \dots, m \cdot n - 1]$$

$$\tilde{\mu}_j = \sum_i (\tilde{A}^{-1})_{ji} \tilde{f}_i$$